# FEDAGENTBENCH: TOWARDS AUTOMATING REAL-WORLD FEDERATED MEDICAL IMAGE ANALYSIS WITH SERVER-CLIENT LLM AGENTS

#### **Anonymous authors**

000

001

002

004

006

008

010 011

012

013

014

015

016

017

018

019

021

025

026

027

028

029

031

033

034

037

040

041

042

043

044

045

046

047

048

051

052

Paper under double-blind review

#### **ABSTRACT**

Federated learning (FL) allows collaborative model training across healthcare sites without sharing sensitive patient data. However, real-world FL deployment is often hindered by complex operational challenges that demand substantial human efforts in cross-client coordination and data engineering. This includes: (a) selecting appropriate clients (hospitals), (b) coordinating between the central server and clients, (c) client-level data pre-processing, (d) harmonizing non-standardized data and labels across clients, and (e) selecting FL algorithms based on user instructions and cross-client data characteristics. However, the existing FL works overlook these practical orchestration challenges. These operational bottlenecks motivate the need for autonomous, agent-driven FL systems, where intelligent agents at each hospital client and the central server agent collaboratively manage FL setup and model training with minimal human intervention. To this end, we first introduce: (i) an agent-driven FL framework that captures key phases of real-world FL workflows from client selection to training completion, and (ii) a benchmark dubbed FedAgentBench that evaluates the ability of LLM agents to autonomously coordinate healthcare FL. Our framework incorporates 40 FL algorithms, each tailored to address diverse task-specific requirements and crossclient characteristics. Furthermore, we introduce a diverse set of complex tasks across 201 carefully curated datasets, simulating 6 modality-specific real-world healthcare environments, viz., Dermatoscopy, Ultrasound, Fundus, Histopathology, MRI, and X-Ray. We assess the agentic performance of 14 open-source and 10 proprietary LLMs spanning small, medium, and large model scales. While some agent cores such as GPT-4.1 and DeepSeek V3 can automate various stages of the FL pipeline, our results reveal that more complex, interdependent tasks based on implicit goals remain challenging for even the strongest models.

#### 1 Introduction and Background

Federated Learning (FL) Li et al. (2021b); McMahan et al. (2017); Li et al. (2020a) allows collaborative model training across multiple healthcare institutions (e.g., hospitals) without sharing raw medical data. A typical FL workflow involves several tightly coupled components: selecting suitable clients for training, preprocessing heterogeneous data locally, harmonizing labels and datasets across clients, coordinating client-server communication, selecting optimal FL algorithm, and aggregating model updates in the server. These components must be executed in a precise and orchestrated manner across multiple clients. Real-world execution of an FL pipeline necessitates close coordination by human data scientists and machine learning engineers in server and client locations, who are tasked with managing a range of demanding communicational and technical operations. These include selecting appropriate client nodes based on task relevance and resource availability, implementing local data preprocessing pipelines (e.g., normalization, filtering, schema mapping), and harmonizing cross-site inconsistencies of data and label spaces. Additionally, they must determine the most suitable FL algorithms, and manage training schedules and aggregation strategies. This manual orchestration poses a significant barrier to scalable and real-time deployment of FL, particularly in sensitive domains like healthcare, where institutions store diverse yet complementary datasets that cannot be centralized due to privacy and compliance constraints. Moreover, many healthcare facilities, especially in low- and middle-income countries (LMICs) and rural areas, lack the resources to hire dedicated data scientists or machine learning engineers, further limiting their ability to participate

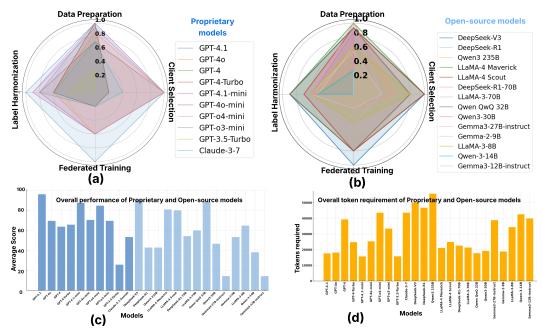


Figure 1: Performance of 24 LLM Agents on 4 FL sub-tasks over 6 healthcare environments. (a) and (b) show the performance of proprietary and open-source models respectively on four subtasks each, *viz.*, Client Selection, Data preprocessing, Label Harmonization, and Federated Training. (c) and (d) show the average score and mean overall token requirement of all models across all tasks.

in FL initiatives despite having valuable local data. To this end, in this paper, we investigate the capabilities of LLM Agents in tackling these issues with minimal human intervention.

The rapid advancement of LLMs has led to the emergence of autonomous AI agents capable of executing complex, multi-step tasks across various domains Gur et al.; Gou et al.; Cai et al.; Li et al. (2023a); Wang et al. (2023); Wu et al.; Mei et al. (2024); Chu et al. (2025); Qiu et al. (2024); Luo et al. (2025). This capability can be particularly transformative for real-world healthcare FL, where agent-based automation can reduce the operational burden on healthcare sites and enable broader participation in collaborative AI development. There are no existing works on agent-driven FL workflow; for general-purpose agents or agentic FL works, refer to **Related Works in Appendix A**.

To this end, we introduce an agentic FL framework (see Figs. 2 & 3) along with a benchmark **FedAgentBench** (see Fig. 1), designed to systematically evaluate the performance of LLM-driven agents in orchestrating FL workflows. To ensure comprehensive coverage, we incorporate 201 datasets, 6 major medical imaging modalities, and 40 representative FL algorithms designed for diverse real-world healthcare objectives and cross-client data compositions. To the best of our knowledge, this is the first work addressing FL problem-solving capabilities of LLM Agents directly dealing with server and client interactions. Our benchmark makes the following key contributions:

- (1) Technical contribution: We first present a plug-and-play modular agentic FL framework supporting 40 FL algorithms and 24 LLM agents. It also allows for easy integration of new FL algorithms, agents and tasks with minimal adaptation. It is a unified FL framework with multi-faceted scenarios, multiple imaging modalities, and complex FL workflow structures. It encompasses four realistic and interlinked agent-driven FL phases: (i) Client Selection, where server and client agents communicate dataset suitability, (ii) Data Preprocessing, involving data restructuring, cleaning, and standardization using learned tools, (iii) Label Harmonization, where agents align inconsistent label taxonomies across clients, and (iv) Federated Model Training, where selected algorithms are deployed in a decentralized setup. It is worth noting that while we simulate healthcare environments in this work, the framework can be readily extended to other FL settings such as finance, IoT, etc.
- (2) Dataset and Task contribution: To evaluate the effectiveness of LLM agents in real-world healthcare tasks, we construct a realistic simulation of inter-hospital collaboration within a FL framework in representative clinical scenarios. Specifically, we curate and publicly release six medical imaging FL agentic environments comprising a total of 201 datasets and a diverse collection of tasks spanning a range of difficulties. To introduce greater variability across clients, we systematically modify the original image resolutions, file format extensions, and intensity distributions



Figure 2: Overview of our agent-driven FL setup. First, user defines task specification. Accordingly, LLM agents perform server-client coordination and complete required tasks using available tools and FL algorithms in any of the 6 modality-specific healthcare environments.

of the client datasets. Additionally, we carefully inject noisy and irrelevant samples spanning images from other modalities, text files, and other extraneous formats into client data directories to simulate realistic uncurated data environments and reflect the challenges of real-world clinical settings.

(3) Empirical contribution: As a part of FedAgentBench, we evaluate the performance of 24 LLM agents across diverse FL tasks based on task completion rate (*i.e.*, success rate), token efficiency, and overall time required. We investigate how varying levels of prompt granularity affect task execution and systematically compare agent performance across different autonomy tiers: guided tool invocation, autonomous planning, and fully independent script generation. Our analysis provides a comprehensive assessment of agentic capabilities and limitations in supporting real-world collaborative healthcare workflows. We will open-source and continuously update the benchmark on Github repository to support FL research and help healthcare data holders fully realize the value of cross-silo data.

# 2 FEDAGENTBENCH FRAMEWORK

#### 2.1 PROBLEM FORMULATION AND OVERVIEW

Given a user-defined task specification for federated medical image analysis, denoted as  $\mathcal{T}$ , our objective is to construct and execute a complete FL pipeline through collaborative decision-making by a set of autonomous agents. As outlined in Fig. 3, FedAgentBench consists of two main components: (i) Federated medical imaging workspace  $\mathcal{W}$  which can be sub-categorized to server workspace  $\mathcal{W}_s$  and client workspace  $\mathcal{W}_c$  as well as (ii) Multi-agent coordination system  $\mathcal{A}$ . The workspace  $\mathcal{W}$  encapsulates the critical resources required for FL pipeline construction and includes: (1) client metadata files (data cards) containing natural language descriptions of local datasets (in  $\mathcal{W}_c$ ), (2) FL algorithm specifications (in  $\mathcal{W}_s$ ) and tool usage descriptions (in  $\mathcal{W}_c$  and  $\mathcal{W}_s$ ) and (3) structured code templates for each phase of the FL workflow (in  $\mathcal{W}_c$  and  $\mathcal{W}_s$ ).

Built on top of this workspace, the agents operate under a divide-and-conquer strategy to address the complexity and modularity of the entire FL process. The server-client agent system  $\mathcal{A} = \{S_1, S_2, S_3, S_4, C_1, C_2, C_3\}$  comprises 7 role-specialized LLM agents (see Fig. 3) responsible for: (1) client selection and server-client communication or orchestration  $(S_1, S_2, C_1)$ , (2) data preprocessing and cleaning  $(C_2)$ , (3) label harmonization  $(C_3)$ , and (4) federated model selection and training  $(S_3, S_4)$ . The collaborative pipeline proceeds iteratively as agents can invoke tools, write scripts, or reason over workspace content to solve subtasks, with execution feedback enabling adaptation. This process can be formally represented as:  $\{D_i, R_i\} = \mathcal{A}(D_{i-1}, R_{i-1}, \mathcal{T} \mid \mathcal{W})$  where  $D_i$  denotes the code, decisions, or configurations generated or modified in the *i*-th iteration, and  $R_i$  represents execution results or tool feedback (e.g., logs, errors, evaluation metrics), with  $D_0 = R_0 = \emptyset$ . The goal is to produce a complete, executable FL pipeline satisfying task specification  $\mathcal{T}$ , measured in terms of success and efficiency under real-world constraints simulated by  $\mathcal{W}$ .

### 2.2 CLIENT DATASET CURATION AND FL ALGORITHM INTEGRATION

Broad coverage of real-world medical specialties and data sets: We construct FedAgentBench clients by adapting 201 publicly available datasets with 2D and 3D dimensionality across 6 different medical imaging modalities *viz.* 25 Dermatology, 33 Ultrasound, 63 Fundus, 32 X-Ray, 28 MRI, and 20 Histopathology datasets. It spans a broad range of tasks, including disease classification (e.g., tumor detection, cancer subtype identification), disease staging or grading (e.g., cancer and diabetic retinopathy severity levels), anatomical or pathological region segmentation (e.g., tumor or stroke localization), object detection, regression, reconstruction, *etc.* Each client is simulated to

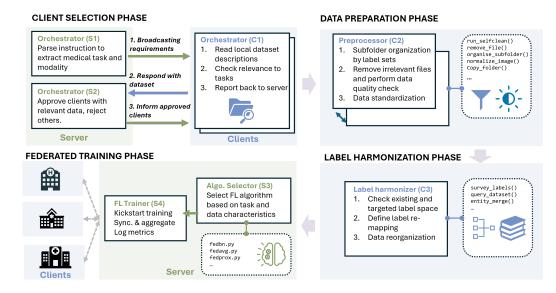


Figure 3: An overview of the FedAgentBench Framework. It comprises 7 role-specialized LLM agents  $(S_1, S_2, S_3, S_4, C_1, C_2, C_3)$  for completing 4 distinct phases of the FL workflow (see §2.3)

comprise one or more of these datasets, reflecting the diversity and heterogeneity typical of real-world healthcare institutions. We construct a datacard accompanying each client based on the metadata sourced from its original publication, repository or website. See Appendix C.1 & Listings 6-8.

**Cross-client data heterogeneity beyond distribution shifts:** In order to introduce greater variability across clients and better emulate the heterogeneity found in real-world clinical data silos, we systematically modify several aspects of the original datasets:

- (i) Structured Dataset Perturbations: We introduce systematic modifications to dataset characteristics, such as varying image resolutions (e.g., downsampling images), altering file format extensions (e.g., converting .png files to .jpeg, .bmp, or .tiff), and modifying intensity distributions to reflect differences in scanner settings or preprocessing pipelines.
- (ii) Inclusion of Uncurated and Irrelevant Files: To reflect the messiness of real-world clinical storage, we inject non-image and unrelated files into client directories. These include textual notes (.txt,.doc,.pdf), spreadsheets (.csv,.xls), and miscellaneous files (e.g., .log,.xml,.ini). For example, our dermatoscopy dataset contains lesion images mixed with dermatologist notes in .pdf format and other unrelated documents.
- (iii) Simulation of Label and Modality Noise: We simulate common data quality issues by introducing random duplication of 2-5 samples, injecting 2-5 anatomically or modality-inconsistent images, and deliberately corrupting labels of 2-5 samples to model annotation noise in each dataset.

These artifacts challenge the robustness of agent-based preprocessing and reflect the complexities encountered in real hospital PACS or data repositories. See Appendix C for more details.

**Algorithm suite for a wide spectrum of FL settings:** As a part of the benchmark design, we also curate a comprehensive suite of **40 FL algorithms** by integrating and adapting existing implementations. This algorithm collection spans a broad spectrum of FL paradigms enabling standardized and reproducible evaluation across diverse medical imaging tasks (**See Appendix §C.4**). This includes:

(i) Classical FL algorithms such as FedAvg, FedProx, and Scaffold, which address baseline aggregation and client drift; (ii) Personalized FL algorithms like Per-FedAvg, pFedMe, and FedRep, which tailor models to heterogeneous client data distributions; (iii) Regularization-based approaches like Ditto which impose constraints to preserve global knowledge during local updates; (iv) Knowledge Distillation-based methods such as FedDF, enabling model-agnostic communication via logits; (v) Domain generalization techniques like FedSR, FedDG, and FedIRM, which aim to learn invariant representations across non-IID clients; and (vi) Optimization and scheduling variants, such as FedNova which address stability, and convergence rate.

#### 2.3 FEDERATED AGENTIC FRAMEWORK CONSTRUCTION

FL workflows typically follow a common set of phases, from which we abstract the key human roles and tasks fundamental to their execution as discussed below (See Appendix B.2 for more details):

1. Client orchestrator agents: These agents act as the coordinators of the framework by communicating between the server and clients as well as by selecting the most suitable clients for the task based on the user requirements and individual client responses (see Fig. 4).

Server agent  $S_1$  interprets the user-defined task  $\mathcal{T}$  and communicates imaging modality/task requirements to initiate client selection. For this, it first parses  $\mathcal{T}$  and broadcasts a query to all Client Agents (i.e., healthcare sites). Each Client Agent  $C_1$  reads local dataset description file, which contains metadata about available datasets, including label sets/imaging

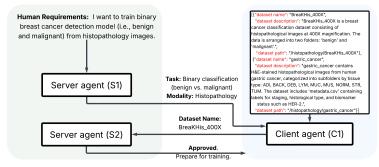


Figure 4: Client orchestrator agents  $S_1$ ,  $C_1$ , and  $S_2$  in a histopathology-based breast cancer classification task

types. Based on semantic and modality matching,  $C_1$  evaluates relevance of its datasets to  $\mathcal{T}$ , returning only matching datasets (if any). Server Agent  $S_2$  collects these responses and selects a subset of relevant clients  $\mathcal{C}_{\text{sel}}$ , which are then approved for further processing (see Figs. 9-14 in Appendix D).

**2.** Data pre-processor agent: It is responsible for preparing selected client datasets for effective participation in the FL pipeline. Given the diversity of data storage formats and quality issues across real-world sites, Data pre-processor agent  $C_2$  at each client ensures that the dataset adheres to a standardized structure and meets minimum quality criteria.

Concretely, it is responsible for standardizing and cleaning datasets at each selected client (see Fig. 5). This includes:

(i) Subfolder Organization: Verifies whether datasets are organized into class-specific subfolders. If disorganized,  $C_2$  restructures the folder hierarchy.

(ii) File Cleaning: Removes irrelevant files (non-image formats .txt, .csv etc.) to ensure format consistency.

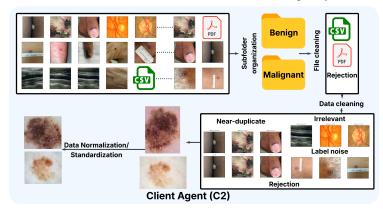


Figure 5: Data pre-processor agent  $\mathcal{C}_2$  in skin cancer detection task

(iii) **Data Cleaning**: Detects and flags duplicates, off-topic samples, and noisy labels, which are then removed. This ensures all selected clients have curated structurally consistent data, enabling downstream harmonization and consequent training (see Figs. 34-35 in Appendix D).

(iv) Data Normalization/Standardization: Standardizes images across clients based on resolution, intensity, and file extension. This agent thus plays an essential role in bridging the gap between raw, heterogeneous clinical data and the clean, harmonized inputs required for FL. Its operations ensure that all participating clients contribute structurally consistent, high-quality data harmonized across clients, which is crucial for the success of the overall FL system.

3. Task-conditioned label harmonizer agent: This agent  $(C_3)$  addresses one of the most critical challenges in multi-institutional FL, *i.e.*, the inconsistency in label nomenclature and granularity across client datasets (see Fig. 6). Due to variations in annotation protocols, terminologies, and domain-specific taxonomies, class labels across clients may not align semantically or structurally.  $C_3$  plays a pivotal role in reconciling these differences based on the user requirements: (i) Class Inspection: Enumerates all class labels present in client datasets.

(ii) Label Mapping: Converts fine-grained labels (e.g., "melanoma", "nevus") to harmonized classes (e.g., "malignant", "benign") according to a self-developed mapping schema.

(iii) Data Reorgani-zation: Reorganizes the dataset structure to reflect harmonized labels, aligning image samples with their mapped class defini-tions. This standardization enables cross-client training without semantic conflicts in label interpretation. 

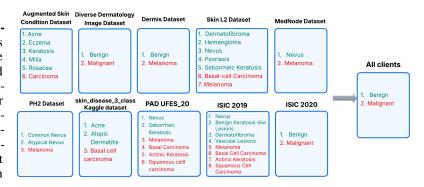


Figure 6: Label harmonization by agent  $C_3$  in dermatology-based skin Through these actions, cancer detection (benign/malignant classes color-coded in green/red) the agent guarantees that all clients adhere to a shared label vocabulary.

- **4. Federated trainer agents:** These agents are responsible for initiating the actual federated training process across the selected set of clients and play a central role in converting the prepared environment into a functioning FL system. They initiate and coordinate federated training in 2 steps:
- (i) Based on  $\mathcal{T}$ , FL Algorithm Selector Agent ( $S_3$ ) queries a registry of 40 FL algorithms containing the algorithmic descriptions and then selects a suitable method (e.g., FedAvg, pFedSim, FedSR) based on user requirements.
- (ii) **Trainer Agent**  $(S_4)$  then distributes training details to approved clients and executes Federated Training. During training,  $S_4$  logs per-client and global metrics (e.g., accuracy) and performs model aggregation. Its modular structure supports plug-and-play experimentation with different FL algorithms and training configurations.

#### 2.4 PRIVACY PRESERVING AND MODULAR DESIGN

A key advantage of our framework is its modular design across phases and agent specializations: Each agent component and phase can be independently evaluated, replaced, or extended. More importantly, this modularity enables future expansion of the benchmark and adaptation to diverse real-world scenarios. For instance, additional components simulating privacy/safety audits conducted by humans or AI can be seamlessly inserted between server and client agents or workflow phases, without the need for altering the existing workflow.

It is to be noted that our framework enforces data privacy by design, aligning fully with FL principles. We explicitly prevent agents from ever accessing or transmitting raw data, model weights, or sensitive metadata. The server receives approvals/configuration signals only, not images, so the agent layer never handles patient data. Instead, agents operate at orchestration layer only and exchange only predefined information (JSON configs, file paths, status signals). They do not have direct access to raw client data (e.g., patient images) or sensitive metadata and never transmit patient data or intermediate outputs externally. Training is invoked via a tool wrapper that runs locally per client; no raw data leaves clients at the agent layer, *i.e.*, federated training is triggered by the agent, but executed on local clients via tools. All data preprocessing and label harmonization also happen locally at clients. Eg: In label harmonization, the agent creates mapping logic, but the mapping execution and label replacement are performed entirely on the local client side.

#### 3 EXPERIMENTS AND RESULTS

#### 3.1 IMPLEMENTATION AND EVALUATION DETAILS

We utilize the LangGraph architecture Langgraph (2025) for agent construction and workflow graph compilation. Each agent is assigned a tailored toolset, drawn from our proposed suite of 16 tools (see Appendix B.1), with the selection guided by the agent's specific role and the need to omit redundant or irrelevant functionalities. In order to assess the capabilities of existing LLM agents, we validate a total number of 24 models on the FedAgentBench datasets, including: (1) 10 representative proprietary LLMs: GPT 4.1, GPT-40, GPT-4, GPT-4-Turbo, GPT 4.1-mini, GPT-40-mini, GPT o4-mini, GPT o3-mini, GPT-3.5 Turbo, and Claude-3.7 Sonnet. (2) 14 state-of-the-art open-sourced LLMs ranging from 9B to 685B: LLaMA series models (LLaMA-4 Maverick, LLaMA-4 Scout, LLaMA-3 70B, LLaMA-3 8B), DeepSeek series models (DeepSeek-V3, Deepseek-R1, DeepSeek-R1-Distill-Llama-70B), Qwen series models (Qwen 3 235B, Qwen QwQ 32B, Qwen 3 30B, Qwen

Table 1: Comparison of LLM agents in **Dermatology** environment based on skin cancer detection task. Here **P**, **R**, **F1** indicate Precision, Recall, and F1 score. **S**, **D**, **F** indicate Schema Compliance Rate, Duplicate Removal Rate, and Format Normalization Rate. **E**, **C**, **Co** indicate Exact-match Accuracy, Coverage Rate, and Conflict Rate. **T** indicates Training-start verification score.

Model	Fine-grained guidance Goal-oriented guidance								
	Client-Sel	Data-Pre	Label-Harm	Fed-Train	Client-Sel	Data-Pre	Label-Harm	Fed-Train	
	P, R, F1	S, D, F	E, C, Co	T	P, R, F1	S, D, F	E, C, Co	T	
Proprietary Models									
GPT-4.1	0.96, 1.00, 0.98	1.00, 0.97, 1.00	0.61, 0.65, 0.35	0.99	0.88, 0.86, 0.87	1.00, 0.96, 0.98	0.61, 0.61, 0.39	0.85	
GPT-40	0.88, 0.89, 0.88	1.00, 0.94, 0.95	0.18, 0.27, 0.73	0.21	0.79, 0.76, 0.77	0.96, 0.91, 0.92	0.16, 0.24, 0.76	0.18	
GPT-4	1.00, 0.92, 0.96	0.02, 0.01, 0.00	0.22, 0.29, 0.71	0.61	0.70, 0.68, 0.69	0.05, 0.00, 0.00	0.00, 0.01, 0.96	0.43	
GPT-4-Turbo	0.91, 0.89, 0.90	0.41, 0.33, 0.39	0.19, 0.24, 0.76	0.64	0.88, 0.79, 0.83	1.00, 0.98, 0.97	0.25, 0.29, 0.71	0.45	
GPT-4.1-mini	1.00, 1.00, 1.00	1.00, 0.93, 0.98	0.59, 0.65, 0.35	0.61	1.00, 0.97, 0.98	0.57, 0.53, 0.57	0.59, 0.60, 0.40	0.58	
GPT-4o-mini	0.64, 0.61, 0.62	1.00, 0.92, 1.00	0.60, 0.63, 0.37	0.61	0.50, 0.56, 0.53	1.00, 0.96, 0.98	0.23, 0.26, 0.74	0.40	
GPT-o4-mini	0.94, 0.91, 0.92	0.98, 0.95, 0.96	0.63, 0.71, 0.29	0.57	0.90, 0.80, 0.85	0.74, 0.70, 0.73	0.45, 0.50, 0.50	0.60	
GPT-o3-mini	0.86, 0.89, 0.87	0.00, 0.00, 0.00	0.45, 0.49, 0.51	0.58	0.71, 0.77, 0.74	0.05, 0.00, 0.00	0.44, 0.50, 0.50	0.63	
GPT-3.5-Turbo	0.32, 0.35, 0.33	0.04, 0.00, 0.00	0.00, 0.03, 0.97	0.18	0.41, 0.30, 0.35	0.43, 0.38, 0.38	0.00, 0.00, 1.00	0.21	
Claude-3-7-Sonnet	0.67, 0.68, 0.67	0.44, 0.42, 0.42	0.21, 0.27, 0.73	0.42	0.69, 0.69, 0.69	0.40, 0.38, 0.39	0.26, 0.32, 0.68	0.44	
			Open-sou	rce Models					
			Huge	Models					
DeepSeek-V3	0.79, 0.78, 0.78	0.97, 0.96, 0.94	1.00, 1.00, 0.00	0.78	0.76, 0.75, 0.75	0.77, 0.73, 0.75	0.81, 0.83, 0.17	0.82	
DeepSeek-R1	0.70, 0.65, 0.67	0.00, 0.00, 0.00	0.02, 0.08, 0.92	0.03	0.68, 0.63, 0.65	0.00, 0.00, 0.00	0.01, 0.01, 0.97	0.00	
Qwen3 235B	0.62, 0.68, 0.65	0.01, 0.00, 0.00	0.02, 0.09, 0.91	0.00	0.64, 0.69, 0.66	0.08, 0.00, 0.00	0.04, 0.08, 0.92	0.01	
LLaMA-4 Maverick	0.65, 0.69, 0.67	0.98, 0.90, 0.97	0.57, 0.66, 0.34	0.37	0.73, 0.64, 0.68	0.98, 0.95, 0.94	0.65, 0.68, 0.32	0.62	
LLaMA-4 Scout	0.75, 0.77, 0.76	1.00, 0.93, 0.95	0.66, 0.73, 0.27	0.41	0.79, 0.80, 0.79	1.00, 0.95, 0.97	0.56, 0.64, 0.36	0.44	
			Large	Models					
DeepSeek-R1-70B	0.71, 0.71, 0.71	0.00, 0.00, 0.00	0.02, 0.03, 0.95	0.19	0.64, 0.72, 0.68	0.00, 0.00, 0.00	0.03, 0.09, 0.91	0.00	
LLaMA-3-70B	0.72, 0.65, 0.68	0.17, 0.11, 0.12	0.17, 0.20, 0.80	0.43	0.70, 0.66, 0.68	0.41, 0.39, 0.39	0.48, 0.55, 0.45	0.20	
			Mediu	n Models					
Owen OwO 32B	0.94, 0.92, 0.93	1.00, 0.96, 1.00	0.87, 0.89, 0.11	0.84	0.86, 0.93, 0.89	1.00, 0.97, 1.00	0.57, 0.65, 0.35	0.64	
Owen3-30B	0.74, 0.68, 0.71	0.04, 0.04, 0.03	0.05, 0.06, 0.94	0.19	0.74, 0.62, 0.67	0.00, 0.00, 0.00	0.01, 0.04, 0.96	0.20	
Gemma3-27B	0.30, 0.38, 0.34	0.00, 0.00, 0.00	0.00, 0.03, 0.97	0.01	0.26, 0.34, 0.29	0.00, 0.00, 0.00	0.00, 0.02, 0.95	0.04	
Small Models									
Gemma-2-9B	0.69, 0.67, 0.68	0.24, 0.15, 0.19	0.19, 0.23, 0.77	0.24	0.60, 0.72, 0.65	0.24, 0.15, 0.17	0.17, 0.21, 0.79	0.19	
LLaMA-3-8B	0.72, 0.65, 0.68	1.00, 0.92, 0.98	0.38, 0.44, 0.56	0.20	0.71, 0.61, 0.66	0.98, 0.95, 0.97	0.45, 0.51, 0.49	0.19	
Owen-3-14B	0.70, 0.69, 0.69	0.04, 0.00, 0.04	0.06, 0.11, 0.89	0.02	0.59, 0.65, 0.62	0.00, 0.00, 0.00	0.03, 0.07, 0.93	0.04	
Gemma3-12B-instruct	0.38, 0.36, 0.37	0.00, 0.00, 0.00	0.00, 0.05, 0.95	0.05	0.34, 0.37, 0.35	0.00, 0.00, 0.00	0.06, 0.08, 0.92	0.04	
Qwen-3-14B	0.70, 0.69, 0.69	0.04, 0.00, 0.04	0.06, 0.11, 0.89	0.02	0.59, 0.65, 0.62	0.00, 0.00, 0.00	0.03, 0.07, 0.93	0.04	

3 14B) and Gemma series models (Gemma 3 27B Instruct, Gemma 3 12B Instruct, Gemma 2 9B Instruct). We utilize APIs from OpenAI (2025), Groq (2025), Deep Infra (2025).

Evaluation metrics: We evaluate the agentic performance using a total of 13 key metrics in different steps of the FL workflow: (1) For each step, we use Success Rate over 5 runs which is a binary indicator of task success/completion. It evaluates the ability of the multi-agent framework to generate executable outputs that satisfy the task requirements. (2) For client selection step, we use **Precision**, **Recall**, and **F1 score** of selected clients vs. the canonical eligible client set. (3) For data pre-processing step, we use (i) Schema Compliance Rate, i.e., proportion of correctly structured folders/files, (ii) Duplicate Removal Rate, i.e., proportion of duplicates removed, and (iii) **Format Normalization Rate**, *i.e.*, proportion of files correctly normalized (e.g., format, resolution). (4) For label harmonization step, we use: (i) Exact-match Accuracy of label mappings vs. the canonical schema, (ii) Coverage Rate, i.e., proportion of local classes successfully mapped, (iii) Conflict Rate, i.e., proportion of classes with ambiguous mappings. (5) For federated training step, we use Training Start Verification as the metric to determine whether the agent produces valid configuration files, initializes the training process, and logs the start signal. Besides, for each step, we also compute (6) Time Spent in seconds which denotes the duration required to complete the task (see Appendix D & Table 16 for comparison of average time); and (7) Token Requirement which indicates the number of tokens involved (see Fig. 1 (d) for comparison of token requirement).

#### 3.2 Main Results and Key Insights

We summarize the overall success scores of all agent cores over 6 modality specific environments with two types of guidance styles for prompting LLMs *viz.*, finegrained guidance (explicit step-by-step instructions) and goaloriented guidance (high-level task description) in Fig. 7. We also show detailed performance breakdown of Dermatology environment in Table 1 and

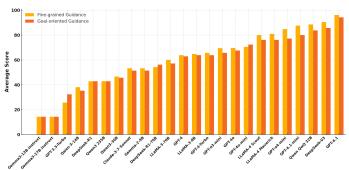


Figure 7: Overall performance of LLM agents in FedAgentBench

Histopathology in Table 2. For detailed results in all other environments, please see Appendix D. & Tables 10-15 Also, see Fig 1 (d) for overall token requirements of each model.

From the tables, we find proprietary models consistently outperform open-source ones across all FL stages. Besides, fine-grained guidance yields higher success rates than goal-oriented prompts for most models. Performance drops in more complex tasks like label harmonization compared to client selection. We also observe that model size alone does not guarantee performance (see Fig. 7). Instead, architectural design and instruction-following capability are more critical.

**Proprietary Model Performance:** GPT-4.1 and GPT-4.1-mini show top-tier performance (85–100%), especially under fine-grained guidance. GPT-40, although newer, struggles with label harmonization and federated training across all environments, leading to overall lower scores (62-71%). Claude-3.7-Sonnet achieves moderate performance (51–57%), inferior to GPT-4 variants. GPT-3.5-Turbo and older variants perform poorly, barely completing the complex stages.

Open-source Model Performance: We discuss agent performance based on model sizes below:

- (i) **Huge Models:** DeepSeek-V3 is the strongest open-source model contender with 80–94% success rate comparable to the best proprietary models. Qwen3 and DeepSeek-R1 perform inconsistently, often failing in more structured stages like data pre-processing and label harmonization.
- (ii) Medium and Large Models: Qwen QwQ 32B demonstrates strong performance (82–91%) and outperforms several proprietary models even under goal-oriented setups. LLaMA-4 Scout and Maverick also deliver competitive performance, especially in label harmonization and federated training, with scores in the 71–94% range Other large models such as LLaMA-3-70B, and Qwen3-30B struggle with most tasks except initial client communication or final training step. Gemma3-27B-instruct is unusable under almost all these settings.
- (iii) Small Models: Performance of 8-14B sized-models drops significantly. Most models (except LLaMA 3 8B) achieve less than or around 50% success. Particularly, Gemma 3-12B-instruct and Qwen 3 14B are observed to fail due to extreme hallucinations. These models are unable to perform any label-oriented reasoning and structured data operations, even under fine-grained instructions.
- **Impact of Task Complexity:** High success is observed in the initial and final steps of client orchestration and federated training across almost all agents, including weaker ones indicating that these tasks are relatively simpler. Data Pre-processing and Label Harmonization are seen to be major differentiators among agents. Weaker agents particularly fail to perform these tasks especially in goal-oriented scenarios, where planning and file structure comprehension are needed. Across almost all agents, label harmonization shows lowest success rates, regardless of guidance type. This suggests that aligning semantic labels across clients remains one of the hardest challenges. Among modalities, histopathology has the highest semantic complexity, potentially due to domain-specific terminology.
- **Granularity of guidance:** In fine-grained guidance, we provide explicit instruction to the models to follow a particular workflow whereas in goal-oriented guidance, we mention the overall objective of the agent without specifying the exact steps, thereby requiring autonomous planning or reasoning. Fine-grained guidance is seen to outperform goal-oriented guidance across almost every model, especially for weaker agents. More capable models like GPT-4.1 and DeepSeek-V3 close this gap, showing their capability to plan even based on implicit prompts.

#### 3.3 AGENT FAILURE ANALYSIS:

- We identify six key recurring failure modes of LLM agents across FL sub-tasks that highlight important limitations of current LLM capabilities in FL workflows (see Appendix D for more details):
- (i) Lack of Domain-Specific Reasoning: The agents frequently fail to apply relevant medical domain knowledge. Eg: In label harmonization (Fig 6), the agents often miss subtle mismatches between dermatology folder names and coarse class labels possibly due to the lack of domain grounding and inability to handle naming conventions specific to medical datasets.
- (ii) Failure in Multi-Step Planning: The agents are often unable to follow multi-step workflows, skipping essential operations where multiple sequential actions are required. Eg: Data pre-processor agents often overlook file/data cleaning steps of Fig. 5 due to multiple tasks in single execution cycle.
- (iii) Overconfidence and Shortcutting: The agents recurrently provide wrong solutions, by defaulting to plausible but incorrect logic when unsure, instead of expressing uncertainty. Eg: Assigning both "nevus" and "melanoma metastasis" to the 'benign' class to simplify label mapping.
- (iv) Hallucination in Structured Multi-Agent Tasks: The agents (particularly DeepSeek R1 and Gemma-based models) often generate irrelevant or unrelated outputs despite specific instructions due to misalignment with structured task formats and poor control over output scope (see Fig. 17-18

Table 2: Comparison in terms of success rate over 5 runs in **Histopathology** environment

Model	Fine-grained guidance				Goal-oriented guidance					
	Client-Sel			mFed-Train	Overall	Client-Sel			nFed-Train	Overall
	$S_1, C_1, S_2$	$C_2$	$C_3$	$S_3, S_4$		$S_1, C_1, S_2$	$C_2$	$C_3$	$S_3, S_4$	
				Proprietary	Models					
GPT-4.1	5/5, 4/5, 5/5	5 - 5/5	5/5	4/5, 5/5		5/5, 4/5, 5/5		5/5	4/5, 5/5	94.29
GPT-4o	5/5, 0/5, 5/5	5/5	2/5	1/5, 5/5	65.71	5/5, 0/5, 5/5	5/5	1/5	1/5, 5/5	62.86
GPT-4	5/5, 1/5, 5/5	0/5	1/5	2/5, 5/5	54.29	5/5, 1/5, 5/5	0/5	0/5	2/5, 5/5	51.43
GPT-4-Turbo	5/5, 1/5, 5/5	1/5	1/5	2/5, 5/5	57.14	5/5, 1/5, 5/5		1/5	2/5, 5/5	65.71
GPT-4.1-mini	$5/5, \frac{3}{5}, \frac{5}{5}$	5/5	4/5	3/5, 5/5	85.71	$5/5, \frac{3/5}{5}, 5/5$	3/5	4/5	3/5, 5/5	80.00
GPT-4o-mini	5/5, 1/5, 3/5	5/5	3/5	2/5, 4/5	65.71	5/5, 1/5, 3/5	5/5	1/5	2/5, 4/5	60.00
GPT-o4-mini	5/5, 2/5, 5/5		3/5	2/5, 5/5	77.14	5/5, 2/5, 5/5		2/5	2/5, 4/5	68.57
GPT-o3-mini	5/5, 5/5, 5/5	0/5	2/5	3/5, 5/5	71.43	$5/5, \frac{4/5}{5}, 5/5$		2/5	3/5, 5/5	68.57
GPT-3.5-Turbo	5/5, 0/5, 0/5	0/5	0/5	1/5, 3/5	25.71	5/5, 0/5, 0/5	2/5	0/5	1/5, 3/5	31.43
Claude-3-7-Sonnet	5/5, 2/5, 3/5		1/5	2/5, 3/5	51.43	5/5, 2/5, 3/5	2/5	1/5	2/5, 5/5	57.14
	, , , ,	,	,	Open-source	e Models	, ,		,	, ,	
				Huge M						
DeepSeek-V3	5/5, 3/5, 5/5	5/5	5/5	4/5, 5/5		5/5, 3/5, 5/5	4/5	5/5	4/5, 5/5	88.57
DeepSeek-R1	5/5, 0/5, 5/5		0/5	0/5, 5/5	42.86	5/5, 0/5, 5/5		0/5	0/5, 5/5	42.86
Owen3 235B	5/5, 0/5, 5/5		0/5	0/5, 5/5	42.86	5/5, 0/5, 5/5		0/5	0/5, 5/5	42.86
LLaMA-4 Maverick	5/5, 2/5, 4/5		3/5	3/5, 5/5	77.14	5/5, 2/5, 4/5		3/5	3/5, 5/5	71.43
LLaMA-4 Scout	5/5, 2/5, 5/5		4/5	2/5, 5/5	80.00	5/5, 2/5, 5/5		3/5	2/5, 5/5	77.14
	-1-1	- / -	7 -	Large M		-1-1	- / -	- / -	1 - 1 - 1 -	
DeepSeek-R1-70B	5/5, 0/5, 5/5	0/5	0/5	0/5, 5/5		5/5, 0/5, 5/5	0/5	0/5	0/5, 5/5	42.86
LLaMA-3-70B	5/5, 1/5, 5/5		1/5	1/5, 5/5		5/5, 1/5, 5/5		2/5	1/5, 5/5	60.00
	0/0, -/0, 0/	-/ -	-/ -	Medium 1		0/0, -/0, 0/0	-/ -	-/ ~	-/0,0/0	
Qwen QwQ 32B	5/5, 4/5, 5/5	3/5	4/5	4/5, 5/5		$5/5, \frac{4}{5}, 5/5$	2/5	4/5	4/5, 5/5	82.86
Owen3-30B	5/5, 0/5, 5/5		0/5	1/5, 5/5		$5/5, \frac{4}{5}, \frac{5}{5}$		0/5	1/5, 5/5	45.71
Gemma3-27B-instruc	t 5/5 0/5 0/5	0/5	0/5	0/5, 0/5	14.29	5/5, 0/5, 0/5		0/5	0/5, 0/5	14.29
Semmas 27B-mstruc	0,0,0,0,0,0	0/0	3/3	Small M		0,0,0,0,0	0/0	3/3	0,0,0,0	11.27
LLaMA-3-8B	5/5, $1/5$ , $5/5$		$\frac{1}{5}$	1/5, 5/5 $1/5, 5/5$		5/5, $1/5$ , $5/5$		$\frac{1}{5}$	1/5, 5/5	65.71
Owen-3-14B	5/5, 0/5, 5/5		0/5	0/5, 5/5	42.86	5/5, 0/5, 5/5		$\frac{2}{5}$	0/5, 4/5	40.00
Gemma3-12B-instruc			0/5		14.29			0/5		14.29
Gennias-12B-mstruc	(0/0, 0/0, 0/8)	0/0	0/5	0/5, 0/5	14.29	5/5, 0/5, 0/5	0/5	0/5	0/5, 0/5	14.29

in Appendix D). Eg: When asked to select skin cancer dataset, Gemma-3 27B Instruct repeatedly returned philosophical or sarcastic monologues in foreign languages, tutorials on freelancing, etc. (v) Task-Type and Modality Mismatch Due to Prior Assumptions: Agents can sometimes confuse tasks or ignore modality constraints due to frequency biases and shallow keyword matching instead of hierarchical task understanding. Eg: Recommending a malignant lesion segmentation dataset for a classification task or ultrasound datasets for histopathology-based breast cancer detection task. (vi) Procedural Overthinking and Paralysis by Analysis: The reasoning/thinking agents often delay execution by speculating about dataset structure or missing dependencies without being asked, potentially due to excessive internal reasoning without grounding in file system or available information (see Fig. 15 in Appendix D). Eg: DeepSeek R1 repeatedly debates whether a client dataset

## 4 CONCLUSION AND LIMITATION

should be selected without reading the dataset description file.

In this paper, we introduced **the first agent-driven FL framework** and an associated benchmark, **FedAgentBench**, for evaluating LLM agents across diverse tasks constituting typical FL workflows. The evaluation covers 24 LLMs with varying sizes and a wide range of FL sub-tasks with varying difficulty levels in six modality-specific FL settings that closely simulate real-world clinical FL environments. Our framework is privacy preserving, comprehensive and modular. It includes 201 medical datasets and 40 FL algorithms and can be easily extended to incorporate more functionalities, settings, and algorithms specific to the user requirement.

We investigated the impact of various factors like FL task complexity and granularity of guidance on the agent performance and analyzed the common failure modes of different agents. Our experiments reveal that across all environments, GPT-4.1 achieves almost perfect scores, under both fine-grained and goal-oriented prompting, whereas GPT-3.5-Turbo, Gemma3 series, and some Qwen variants consistently underperform across all stages and environments. DeepSeek-V3, Qwen QwQ 32B, and LLaMA-4 Maverick are the most reliable open-source agents across tasks. Unsurprisingly, fine-grained guidance consistently outperforms goal-oriented prompting, especially for less capable models. Our findings highlight that the order of complexity of the FL sub-tasks for most agents is: Label Harmonization > Data Pre-processing > Federated Training > Client Orchestration. Our experiments also show that larger model size does not necessarily correlate with better performance, *i.e.*, some mid-sized models (30–40B) outperform larger ones (70B+). E.g., Qwen QwQ 32B consistently outperforms Qwen3-235B and DeepSeek-R1-70B.

This work is a first step toward agent-driven FL and comes with some limitations: (a) We currently assume stable network conditions and do not model dynamic communication bandwidth. (b) We do not incorporate real-time monitoring or interruption mechanisms. (c) We do not simulate safety check or regulatory compliance assessment but it can be seamlessly integrated into the system.

```
486
      REFERENCES
487
      3d-mri-ultrasound-brain dataset. https://www.kaggle.com/datasets/shubhamcodez/
488
         3d-mri-ultrasound-brain-images. Accessed: 2025-05-22.
489
490
      Buet
                             https://www.kaggle.com/datasets/jarintasnim090/
491
        buet-breast-ultrasound-data. Accessed: 2025-05-22.
492
                                      https://www.kaggle.com/datasets/orvile/
493
494
        bus-uclm-breast-ultrasound-dataset, a. Accessed: 2025-05-22.
495
      Busi
                               https://www.kaggle.com/datasets/sabahesaraki/
496
        breast-ultrasound-images-dataset, b. Accessed: 2025-05-22.
497
498
      Breast-cancer
                   dataset.
                                 https://www.kaggle.com/datasets/aryashah2k/
499
        breast-ultrasound-images-dataset, a. Accessed: 2025-05-22.
500
      Breast-ultrasound-images dataset. https://www.kaggle.com/datasets/aryashah2k/
501
        breast-ultrasound-images-dataset, note = Accessed: 2025-05-22, b.
502
      Camus-human-heart-data.
                                 https://www.kaggle.com/datasets/shoybhasan/
504
         camus-human-heart-data. Accessed: 2025-05-22.
505
      Ct2us
                dataset.
                                     https://www.kaggle.com/datasets/siatsyx/
506
        ct2usforkidneyseg. Accessed: 2025-05-22.
507
508
      Carotid-ultrasound-image dataset.
                                      https://www.kaggle.com/datasets/orvile/
509
         carotid-ultrasound-images, a. Accessed: 2025-05-22.
510
      Carotid-artery dataset. https://www.kaggle.com/datasets/pahunichoudhary/
511
        carotid-artery-ultrasound-and-color-doppler, b. Accessed: 2025-05-22.
512
513
      Ddti
                                https://www.kaggle.com/datasets/dasmehdixtr/
514
        ddti-thyroid-ultrasound-images. Accessed: 2025-05-22.
515
      Dr dataset. https://www.kaggle.com/c/diabetic-retinopathy-detection/
516
         data. Accessed: 2025-05-22.
517
518
      Dermatology-tabular-dataset. https://www.kaggle.com/datasets/olcaybolat1/
519
         dermatology-dataset-classification, a. Accessed: 2025-05-22.
520
                            https://www.kaggle.com/datasets/farhatullah8398/
521
         skin-lesion-dermis-dataset, b. Accessed: 2025-05-22.
522
523
      Dermnet dataset. https://www.kaggle.com/datasets/shubhamgoel27/dermnet, c.
524
        Accessed: 2025-05-22.
525
526
      Dermquest dataset. http://www.dermquest.com, d. Accessed: 2025-05-22.
527
                                  https://www.kaggle.com/datasets/ankit8467/
528
         fetal-head-ultrasound-dataset-for-image-segment, a. Accessed: 2025-
529
        05-22.
530
531
      Fetal-health-classification dataset. https://www.kaggle.com/datasets/andrewmvd/
532
         fetal-health-classification, b. Accessed: 2025-05-22.
533
      Isic2020
                    dataset.
                                              https://kaggle.com/competitions/
534
        siim-isic-melanoma-classification, a. Accessed: 2025-05-22.
535
      Isic2024 dataset. https://kaggle.com/competitions/isic-2024-challenge, b.
        Accessed: 2025-05-22.
538
      Liver-histopathology dataset.
                                  http://kaggle.com/datasets/vibhingupta028/
539
```

liver-histopathology-fibrosis-ultrasound-images. Accessed: 2025-05-22.

```
540
      Mednode dataset. https://www.kaggle.com/datasets/prabhavsanga/med-node.
541
        Accessed: 2025-05-22.
542
      Monkeypox-skin-image-dataset.
                                      https://www.kaggle.com/datasets/nafin59/
543
        monkeypox-skin-lesion-dataset. Accessed: 2025-05-22.
544
545
      Pcos
                            https://www.kaggle.com/datasets/anaghachoudhari/
546
        pcos-detection-using-ultrasound-images. Accessed: 2025-05-22.
547
548
      Ph2dataset. https://www.kaggle.com/datasets/spacesurfer/ph2-dataset. Ac-
549
        cessed: 2025-05-22.
550
       Rsna-asnr dataset. https://www.kaggle.com/c/rsna-miccai-brain-tumor-radiogenomic-classifica
551
        Accessed: 2025-05-22.
552
553
                                               https://www.kaggle.com/datasets/
      Regensburg pediatric appendicitis dataset.
554
         joebeachcapital/regensburg-pediatric-appendicitis. Accessed: 2025-05-
556
    Robotic<sub>h</sub> and held_lumbar_spine_usdataset.. Accessed: 2025-05-22.
558
       Stare dataset.
                             https://www.kaggle.com/datasets/vidheeshnacode/
559
       stare-dataset. Accessed: 2025-05-22.
560
      Thyroid<sub>u</sub>ltrasounddataset.. Accessed: 2025-05-22.
561
562
       Us3m
                                      https://www.kaggle.com/datasets/timesxy/
563
      multimodal-breast-ultrasound-dataset-us3m. Accessed: 2025-05-22.
564
565
                              https://www.kaggle.com/datasets/imtkaggleteam/
       abdomen-mri dataset.
566
       abdomen-mri-object-detection. Accessed: 2025-05-22.
567
      alzheimer-mri dataset.
                                 ://www.kaggle.com/datasets/borhanitrash/alzheimer-mri-disease-
568
      classification-dataset.
569
570
      augmented-skin-conditions-image-dataset.
                                               https://www.kaggle.com/datasets/
571
       syedalinaqvi/augmented-skin-conditions-image-dataset. Accessed: 2025-
      05-22.
572
573
                                       https://www.kaggle.com/datasets/orvile/
       brain-cancer
                   dataset.
574
      brain-cancer-mri-dataset, a. Accessed: 2025-05-22.
575
576
      brain-mri-images dataset.
                                 https://www.kaggle.com/datasets/ashfakyeafi/
577
      brain-mri-images, b. Accessed: 2025-05-22.
578
      brain-tumor dataset.
                                 https://www.kaggle.com/datasets/ultralytics/
579
      brain-tumor, c. Accessed: 2025-05-22.
580
581
      brain-tumor-classification dataset.
                                               https://www.kaggle.com/datasets/
582
       sartajbhuvaji/brain-tumor-classification-mri, d. Accessed: 2025-05-22.
583
      brain-tumor-detection dataset.
                                    https://www.kaggle.com/datasets/abhranta/
584
      brain-tumor-detection-mri, e. Accessed: 2025-05-22.
585
586
      chasedb1 dataset. https://blogs.kingston.ac.uk/retinal/chasedb1/.
      dermatologic-ultrasound dataset.
                                    https://www.kaggle.com/datasets/alfageme/
588
       dermatologic-ultrasound-images. Accessed: 2025-05-22.
589
                              https://www.kaggle.com/datasets/humanaizedata/
591
       facial-mri-dataset-boost-your-ai-models. Accessed: 2025-05-22.
592
       fallmud dataset. https://www.kaggle.com/datasets/angeliqueloesch/fallmud.
593
      Accessed: 2025-05-22.
```

```
594
       fetalultrasoundbrain dataset.
                                 https://www.kaggle.com/datasets/rahimalargo/
595
       fetalultrasoundbrain. Accessed: 2025-05-22.
596
597
       fhms-ultrasound dataset.
                                       https://www.kaggle.com/datasets/jai132/
       fhms-ultrasound-dataset. Accessed: 2025-05-22.
598
       foraminal-stenosis dataset.
                                    https://www.kaggle.com/datasets/axondata/
600
       foraminal-stenosis-mri-dataset. Accessed: 2025-05-22.
601
602
       glioma dataset.
                                 https://www.kaggle.com/datasets/azharsaleem/
      mri-based-glioma-detection-dataset-with-masks. Accessed: 2025-05-22.
603
604
                                       https://www.kaggle.com/datasets/orvile/
605
       human-bone-fractures-image-dataset-hbfmid. Accessed: 2025-05-22.
606
      heart-mri
                dataset.
                                   https://www.kaggle.com/datasets/adarshsng/
607
      heart-mri-image-dataset-left-atrial-segmentation. Accessed: 2025-05-22.
608
609
       hippocampal dataset.
                                  https://www.kaggle.com/datasets/aryashah2k/
610
       hippocampal-sparing-dataset. Accessed: 2025-05-22.
611
       indian-diabetic-retinopathy dataset.
                                               https://www.kaggle.com/datasets/
612
       aaryapate198/indian-diabetic-retinopathy-image-dataset.
613
      2025-05-22.
614
615
                                       https://www.kaggle.com/datasets/orvile/
      isles
              dataset
616
       isles-2022-brain-stoke-dataset. Accessed: 2025-05-22.
617
       Fetal_health_classification dataset.. Accessed: 2025-05-22.
618
619
       mendeley<sub>b</sub>reast_u s datas et.. Accessed: 2025-05-22.
620
       mri-and-pet dataset.
                                 https://www.kaggle.com/datasets/grantmcnatt/
621
      mri-and-pet-dice-similarity-dataset. Accessed: 2025-05-22.
622
623
       multimodal-ultrasound-vascular-segmentation dataset.
                                                         https://www.kaggle.com/
624
       datasets/among22/multimodal-ultrasound-vascular-segmentation.
625
      cessed: 2025-05-22.
626
      oasis-1 dataset.
                             https://www.kaggle.com/datasets/mdfahimbinamin/
627
       oasis-1-fastsurfer-quickseq-segmentation-dataset. Accessed: 2025-05-22.
628
629
       phantom dataset.
                                 https://www.kaggle.com/datasets/ukeppendorf/
630
       frequently-traveling-human-phantom-fthp-dataset.
                                                                   Accessed: 2025-05-
631
       22.
632
       pmram
                dataset.
                                       https://www.kaggle.com/datasets/orvile/
633
      pmram-bangladeshi-brain-cancer-mri-dataset. Accessed: 2025-05-22.
634
635
      prostate-annotated dataset.
                                 https://www.kaggle.com/datasets/haithem1999/
636
      prostate-annotated-dataset-for-image-segmentation, a. Accessed: 2025-05-
637
638
                                       https://www.kaggle.com/datasets/dsptlp/
       prostate-mri-us dataset.
639
       prostate-mri-us-biopsy, b. Accessed: 2025-05-22.
640
641
       refuge2 dataset. https://www.kaggle.com/datasets/victorlemosml/refuge2.
642
      Accessed: 2025-05-22.
643
      skin_d isease_c ls_k apple dataset., a. Accessed: 2025-05-22.
644
645
      skin_i n fection dataset., b. Accessed: 2025-05-22.
646
                            https://www.kaggle.com/datasets/trainingdatapro/
       spinal-cord-dataset.
647
```

spinal-cord-dataset. Accessed: 2025-05-22.

```
stroke_head_mridataset.. Accessed: 2025-05-22.
```

668

673

685

694

- ultra-lr-hr-ultrasound dataset. https://www.kaggle.com/datasets/chirag2466/ ultra-lr-hr-ultrasound-image-dataset-for-research, a. Accessed: 2025-05-22.
- ultrasound-nerve-segmentation dataset. https://www.kaggle.com/competitions/ultrasound-nerve-segmentation, b. Accessed: 2025-05-22.
- wmh-dataset. https://www.kaggle.com/datasets/farahmo/wmh-dataset. Accessed: 2025-05-22.
- Mohammad Amin Abbasi, Farnaz Sadat Mirnezami, and Hassan Naderi. Hamraz: A culture-based persian conversation dataset for person-centered therapy using llm agents. *arXiv preprint arXiv:2502.05982*, 2025.
- Durmus Alp Emre Acar, Yue Zhao, Ramon Matas Navarro, Matthew Mattina, Paul N. Whatmough, and Venkatesh Saligrama. Federated learning based on dynamic regularization, 2021. URL https://arxiv.org/abs/2111.04263.
- Vaibhav Aggarwal, Ojasv Kamal, Abhinav Japesh, Zhijing Jin, and Bernhard Schölkopf. Dars:
  Dynamic action re-sampling to enhance coding agent performance by adaptive tree traversal. *arXiv* preprint arXiv:2503.14269, 2025.
- Walid Al-Dhabyani, Mohammed Gomaa, Hussien Khaled, and Aly Fahmy. Dataset of breast ultrasound images. *Data in Brief*, 28:1–5, 2020.
- Alibaba. Qwen 3: High-performance multilingual llms. https://qwenlm.github.io/, 2025. Accessed: 2025-05-22.
- Alibabaei78. Ebhi segmentation dataset. https://www.kaggle.com/datasets/alibabaei78/ebhi-seg, 2024. Accessed: 2024-05-22.
- Anthropic. Claude 3.7 sonnet. https://www.anthropic.com/index/claude, 2024. Accessed: 2025-05-22.
- Rodolfo Stoffel Antunes, Cristiano André da Costa, Arne Küderle, Imrana Abdullahi Yari, and Björn Eskofier. Federated learning for healthcare: Systematic review and architecture proposal. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 13(4):1–23, 2022.
- Manoj Ghuhan Arivazhagan, Vinay Aggarwal, Aaditya Kumar Singh, and Sunav Choudhary. Federated learning with personalization layers, 2019. URL https://arxiv.org/abs/1912.00818.
- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and Charles Sutton. Program Synthesis with Large Language Models, August 2021. URL http://arxiv.org/abs/2108.07732. arXiv:2108.07732 [cs].
- Reza Averly, Frazier N Baker, and Xia Ning. Liddia: Language-based intelligent drug discovery agent. *arXiv preprint arXiv:2502.13959*, 2025.
- Beosup. Lung segment dataset. https://www.kaggle.com/datasets/beosup/lung-segment, 2023. Accessed: 2025-05-22.
- Olivier Bernard et al. Deep learning techniques for automatic MRI cardiac multi-structures segmentation and diagnosis: is the problem solved? *IEEE Transactions on Medical Imaging*, 37(11): 2514–2525, 2018.
- Matt Berseth. Isic 2017-skin lesion analysis towards melanoma detection. *arXiv preprint* arXiv:1703.00523, 2017.
- Michał Byra, Grzegorz Styczynski, Cezary Szmigielski, Piotr Kalinowski, Łukasz Michałowski, Rafał Paluszkiewicz, Bogna Ziarkiewicz-Wróblewska, Krzysztof Zieniewicz, Piotr Sobieraj, and

Andrzej Nowicki. Transfer learning with deep convolutional neural network for liver steatosis assessment in ultrasound images. *International journal of computer assisted radiology and surgery*, 13(12):1895–1903, 2018.

Tianle Cai, Xuezhi Wang, Tengyu Ma, Xinyun Chen, and Denny Zhou. Large language models as tool makers. In *The Twelfth International Conference on Learning Representations*.

John Capocyan. Cellnet beta version. https://www.kaggle.com/datasets/johncapocyan/cellnet-beta-version, 2024. Accessed: 2024-05-22.

Dengsheng Chen, Jie Hu, Vince Junkai Tan, Xiaoming Wei, and Enhua Wu. Elastic aggregation for federated optimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 12187–12197, June 2023a.

Hao-Yuan Chen, Cheng-Pong Huang, and Jui-Ming Yao. Verbal process supervision elicits better coding agents. *arXiv preprint arXiv:2503.18494*, 2025a.

Hong-You Chen and Wei-Lun Chao. On bridging generic and personalized federated learning for image classification, 2022. URL https://arxiv.org/abs/2107.00778.

Junying Chen, Chi Gui, Anningzhe Gao, Ke Ji, Xidong Wang, Xiang Wan, and Benyou Wang. Cod, towards an interpretable medical agent using chain of diagnosis. *arXiv preprint arXiv:2407.13301*, 2024.

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating Large Language Models Trained on Code, July 2021. URL http://arxiv.org/abs/2107.03374. arXiv:2107.03374 [cs].

Yiqiang Chen, Wang Lu, Xin Qin, Jindong Wang, and Xing Xie. Metafed: Federated learning among federations with cyclic knowledge distillation for personalized healthcare, 2023b. URL https://arxiv.org/abs/2206.08516.

Zhaoling Chen, Xiangru Tang, Gangda Deng, Fang Wu, Jialong Wu, Zhiwei Jiang, Viktor Prasanna, Arman Cohan, and Xingyao Wang. Locagent: Graph-guided llm agents for code localization. *arXiv* preprint arXiv:2503.09089, 2025b.

Zhen Chen, Zhihao Peng, Xusheng Liang, Cheng Wang, Peigan Liang, Linsheng Zeng, Minjie Ju, and Yixuan Yuan. Map: Evaluation and multi-agent enhancement of large language models for inpatient pathways. *arXiv* preprint arXiv:2503.13205, 2025c.

Hojun Cho, Donghu Kim, Soyoung Yang, Chan Lee, Hunjoo Lee, and Jaegul Choo. Building resource-constrained language agents: A korean case study on chemical toxicity information. *arXiv* preprint arXiv:2503.17753, 2025.

Zhendong Chu, Shen Wang, Jian Xie, Tinghui Zhu, Yibo Yan, Jinheng Ye, Aoxiao Zhong, Xuming Hu, Jing Liang, Philip S Yu, et al. Llm agents for education: Advances and applications. *arXiv* preprint arXiv:2503.11733, 2025.

Noel Codella, Veronica Rotemberg, Philipp Tschandl, M Emre Celebi, Stephen Dusza, David Gutman, Brian Helba, Aadi Kalloo, Konstantinos Liopyris, Michael Marchetti, et al. Skin lesion analysis toward melanoma detection 2018: A challenge hosted by the international skin imaging collaboration (isic). *arXiv preprint arXiv:1902.03368*, 2019.

- Liam Collins, Hamed Hassani, Aryan Mokhtari, and Sanjay Shakkottai. Exploiting shared representations for personalized federated learning, 2023. URL https://arxiv.org/abs/2102.07078.
- Marc Combalia, Noel C. F. Codella, Veronica Rotemberg, Brian Helba, Veronica Vilaplana, Ofer Reiter, Cristina Carrera, Alicia Barreiro, Allan C. Halpern, Susana Puig, and Josep Malvehy. Bcn20000: Dermoscopic lesions in the wild, 2019. URL https://arxiv.org/abs/1908.02288.
- Roxana Daneshjou, Kailas Vodrahalli, Roberto A Novoa, Melissa Jenkins, Weixin Liang, Veronica Rotemberg, Justin Ko, Susan M Swetter, Elizabeth E Bailey, Olivier Gevaert, et al. Disparities in dermatology ai performance on a diverse, curated clinical image set. *Science advances*, 8(31): eabq6147, 2022.
- P. K. Darabi. Bone break classification image dataset. https://www.kaggle.com/datasets/pkdarabi/bone-break-classification-image-dataset, 2023. Accessed: 2025-05-22.
- Sergio MM de Faria, Jose N Filipe, Pedro MM Pereira, Luis MN Tavora, Pedro AA Assuncao, Miguel O Santos, Rui Fonseca-Pinto, Felicidade Santiago, Victoria Dominguez, and Martinha Henrique. Light field image dataset of skin lesions. In 2019 41st Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC), pp. 3905–3908. IEEE, 2019.
- Coen De Vente, Koenraad A Vermeer, Nicolas Jaccard, He Wang, Hongyi Sun, Firas Khader, Daniel Truhn, Temirgali Aimyshev, Yerkebulan Zhanibekuly, Tien-Dung Le, et al. Airogs: Artificial intelligence for robust glaucoma screening challenge. *IEEE transactions on medical imaging*, 43(1): 542–557, 2023.
- Etienne Decenciere, Guy Cazuguel, Xiwei Zhang, Guillaume Thibault, J-C Klein, Fernand Meyer, Beatriz Marcotegui, Gwénolé Quellec, Mathieu Lamard, Ronan Danno, et al. Teleophta: Machine learning and image processing methods for teleophthalmology. *Irbm*, 34(2):196–203, 2013.
- Deep Infra. Deepinfra models documentation, 2025. URL https://deepinfra.com/docs/models. Accessed: 2025-05-16.
- Google DeepMind. Gemma: Lightweight open models by google deepmind. https://ai.google.dev/gemma, 2024. Accessed: 2025-05-22.
- DeepSeek. Deepseek-rl and rl-distill model release. https://api-docs.deepseek.com, 2024a. Accessed: 2025-05-22.
- DeepSeek. Deepseek-v3: Advanced open llm. https://deepseek.com, 2024b. Accessed: 2025-05-22.
- DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025. URL https://arxiv.org/abs/2501.12948.
- Qixin Deng, Qikai Yang, Ruibin Yuan, Yipeng Huang, Yi Wang, Xubo Liu, Zeyue Tian, Jiahao Pan, Ge Zhang, Hanfeng Lin, et al. Composerx: Multi-agent symbolic music composition with llms. *arXiv preprint arXiv:2404.18081*, 2024.
- Yihe Deng and Paul Mineiro. Flow-dpo: Improving llm mathematical reasoning through online multi-agent learning. *arXiv preprint arXiv:2410.22304*, 2024.
- Yuyang Deng, Mohammad Mahdi Kamani, and Mehrdad Mahdavi. Adaptive personalized federated learning, 2020. URL https://arxiv.org/abs/2003.13461.
- Canh T. Dinh, Nguyen H. Tran, and Tuan Dung Nguyen. Personalized federated learning with moreau envelopes, 2022. URL https://arxiv.org/abs/2006.08848.
- Tuan Le Dinh. Monuseg 2018. https://www.kaggle.com/datasets/tuanledinh/monuseg2018, 2024. Accessed: 2024-05-22.
- Wei Dong. The ann arbor architecture for agent-oriented programming. *arXiv preprint* arXiv:2502.09903, 2025.

817

819

820

823

824

825

826

827

828

829

830

831 832

833

834

835 836

837

838

839

840 841

842

843

844 845

846 847

848

849

850 851

852

853

854 855

856

857

858

862

863

```
810
        Zhuoyun Du, Lujie Zheng, Renjun Hu, Yuyang Xu, Xiawei Li, Ying Sun, Wei Chen, Jian Wu, Haolei
811
        Cai, and Haohao Ying. Llms can simulate standardized patients via agent coevolution. arXiv preprint
812
        arXiv:2412.11716, 2024.
```

Gaurav Dutta. Fracturefusion: A symphony of bone breaks. https://www.kaggle.com/ 814 datasets/gauravduttakiit/fracturefusion-a-symphony-of-bone-breaks, 815 2023. Accessed: 2025-05-22. 816

Sachi Dwivedi. Kmc kidney histopathology dataset. https://www.kaggle.com/datasets/ 818 sachidwivedi1234/kmc-kidney-histopathology-dataset2, 2024. 2024-05-22.

821 factory.ai. Code Droid Technical Report, June 2024. URL https://www.factory.ai/news/ 822 code-droid-technical-report.

Alireza Fallah, Aryan Mokhtari, and Asuman Ozdaglar. Personalized federated learning with theoretical guarantees: A model-agnostic meta-learning approach. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), Advances in Neural Information Processing Systems, volume 33, pp. 3557–3568. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper\_files/paper/2020/ file/24389bfe4fe2eba8bf9aa9203a44cdad-Paper.pdf.

Adibvafa Fallahpour, Jun Ma, Alif Munim, Hongwei Lyu, and BO WANG. Medrax: Medical reasoning agent for chest x-ray. In Forty-second International Conference on Machine Learning.

Sorouralsadat Fatemi and Yuheng Hu. Enhancing financial question answering with a multi-agent reflection framework. In Proceedings of the 5th ACM International Conference on AI in Finance, pp. 530–537, 2024.

George Fatouros, Kostas Metaxas, John Soldatos, and Manos Karathanassis. Marketsenseai 2.0: Enhancing stock analysis through llm agents. arXiv preprint arXiv:2502.00415, 2025.

Jared Feng. Histopath-sn. https://www.kaggle.com/datasets/jaredfeng/ histopathsn, 2024. Accessed: 2024-05-22.

Jinghao Feng, Qiaoyu Zheng, Chaoyi Wu, Ziheng Zhao, Ya Zhang, Yanfeng Wang, and Weidi Xie. M<sup>^</sup> 3builder: A multi-agent system for automated machine learning in medical imaging. arXiv preprint arXiv:2502.20301, 2025.

Forderation. Breakhis 400x. https://www.kaggle.com/datasets/forderation/ breakhis-400x, 2024. Accessed: 2024-05-22.

Francisco Fumero, Silvia Alayón, José L Sanchez, Jose Sigut, and M Gonzalez-Hernandez. Rim-one: An open retinal image database for optic nerve evaluation. In 24th International Symposium on Computer-Based Medical Systems (CBMS), pp. 1–6. IEEE, 2011.

Fatemeh Ghezloo, Mehmet Saygin Seyfioglu, Rustin Soraki, Wisdom O Ikezogwo, Beibin Li, Tejoram Vivekanandan, Joann G Elmore, Ranjay Krishna, and Linda Shapiro. Pathfinder: A multi-modal multi-agent system for medical diagnostic decision-making applied to histopathology. arXiv preprint arXiv:2502.08916, 2025.

Arsham Gholamzadeh Khoee, Shuai Wang, Yinan Yu, Robert Feldt, and Dhasarathy Parthasarathy. Gatelens: A reasoning-enhanced llm agent for automotive software release analytics. arXiv e-prints, pp. arXiv-2503, 2025.

859 Mohamed Gobara. Multi-class knee osteoporosis x-rav 860 dataset. https://www.kaggle.com/datasets/mohamedgobara/ 861 multi-class-knee-osteoporosis-x-ray-dataset, 2023a. Accessed: 2025-05-22.

Mohamed Gobara. Osteoporosis database. https://www.kaggle.com/datasets/ mohamedgobara/osteoporosis-database, 2023b. Accessed: 2025-05-22.

868

869

870 871

872

873

885

892

897

898

899

900

901

905

906

907

908

909

Juraj Gottweis, Wei-Hung Weng, Alexander Daryin, Tao Tu, Anil Palepu, Petar Sirkovic, Artiom Myaskovsky, Felix Weissenberger, Keran Rong, Ryutaro Tanno, et al. Towards an ai co-scientist. *arXiv preprint arXiv:2502.18864*, 2025.

Zhibin Gou, Zhihong Shao, Yeyun Gong, Yujiu Yang, Minlie Huang, Nan Duan, Weizhu Chen, et al. Tora: A tool-integrated reasoning agent for mathematical problem solving. In *The Twelfth International Conference on Learning Representations*.

- Dennis Grinwald, Philipp Wiesner, and Shinichi Nakajima. Federated learning over connected modes, 2025. URL https://arxiv.org/abs/2403.03333.
- Fabian Gröger, Simone Lionetti, Philippe Gottfrois, Alvaro Gonzalez-Jimenez, Matthew Groh, Roxana Daneshjou, Alexander A Navarini, Marc Pouly, Labelling Consortium, et al. Towards reliable dermatology evaluation benchmarks. In *Machine Learning for Health (ML4H)*, pp. 101–128. PMLR, 2023.
- Fabian Gröger, Simone Lionetti, Philippe Gottfrois, Alvaro Gonzalez-Jimenez, Ludovic Amruthalingam, Matthew Groh, Alexander Navarini, and Marc Pouly. Intrinsic self-supervision for data quality audits. *Advances in Neural Information Processing Systems*, 37:92273–92316, 2024.
- Fabian Gröger, Simone Lionetti, Philippe Gottfrois, Alvaro Gonzalez-Jimenez, Ludovic Amruthalingam, Elisabeth Victoria Goessinger, Hanna Lindemann, Marie Bargiela, Marie Hofbauer, Omar Badri, et al. Cleanpatrick: A benchmark for image data cleaning. *arXiv preprint arXiv:2505.11034*, 2025.
- Matthew Groh, Caleb Harris, Luis Soenksen, Felix Lau, Rachel Han, Aerin Kim, Arash Koochek, and Omar Badri. Evaluating deep neural networks trained on clinical images in dermatology with the fitzpatrick 17k dataset. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 1820–1828, 2021.
- Groq. Groqcloud supported models, 2025. URL https://console.groq.com/docs/models. Accessed: 2025-05-16.
- Xuehang Guo, Xingyao Wang, Yangyi Chen, Sha Li, Chi Han, Manling Li, and Heng Ji. Syncmind: Measuring agent out-of-sync recovery in collaborative software engineering. *arXiv preprint arXiv:2502.06994*, 2025.
  - Yaming Guo, Kai Guo, Xiaofeng Cao, Tieru Wu, and Yi Chang. Out-of-distribution generalization of federated learning via implicit invariant relationships. In *International Conference on Machine Learning*, pp. 11905–11933. PMLR, 2023.
  - Shyam Gupta. Fracatlas. https://www.kaggle.com/datasets/shyamgupta196/fracatlas, 2023. Accessed: 2025-05-22.
- Izzeddin Gur, Hiroki Furuta, Austin V Huang, Mustafa Safdari, Yutaka Matsuo, Douglas Eck, and
   Aleksandra Faust. A real-world webagent with planning, long context understanding, and program
   synthesis. In *The Twelfth International Conference on Learning Representations*.
  - David Gutman, Noel CF Codella, Emre Celebi, Brian Helba, Michael Marchetti, Nabin Mishra, and Allan Halpern. Skin lesion analysis toward melanoma detection: A challenge at the international symposium on biomedical imaging (isbi) 2016, hosted by the international skin imaging collaboration (isic). arXiv preprint arXiv:1605.01397, 2016.
- Haashaatif. Melanoma histopathology dataset. https://www.kaggle.com/datasets/
   haashaatif/melanoma-histopathology-dataset, 2024. Accessed: 2024-05-22.
- Senyu Han, Lu Chen, Li-Min Lin, Zhengshan Xu, and Kai Yu. Ibsen: Director-actor agent collaboration for controllable and interactive drama script generation. *arXiv preprint arXiv:2407.01093*, 2024a.
- Shijie Han, Changhai Zhou, Yiqing Shen, Tianning Sun, Yuhua Zhou, Xiaoxia Wang, Zhixiao Yang, Jingshu Zhang, and Hongguang Li. Finsphere: A conversational stock analysis agent equipped with quantitative tools based on real-time database. *arXiv preprint arXiv:2501.12399*, 2025.

947

- Xuewen Han, Neng Wang, Shangkun Che, Hongyang Yang, Kunpeng Zhang, and Sean Xin Xu. Enhancing investment analysis: Optimizing ai-agent collaboration in financial research. In *Proceedings of the 5th ACM International Conference on AI in Finance*, pp. 538–546, 2024b.
- Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin Burns, Samir Puranik, Horace He, Dawn Song, and Jacob Steinhardt. Measuring Coding Challenge Competence With APPS, November 2021. URL http://arxiv.org/abs/2105.09938. arXiv:2105.09938 [cs].
- 926 HFUTYBX. Mhsi choledoch dataset preprocessed. https://www.kaggle.com/datasets/ 927 hfutybx/mhsi-choledoch-dataset-preprocessed-dataset, 2024. Accessed: 928 2024-05-22.
- Hmchuong. X-ray bone shadow suppression. https://www.kaggle.com/datasets/hmchuong/xray-bone-shadow-supression, 2023. Accessed: 2025-05-22.
- Huthayfa Hodeb. Nih chest x-rays (bbox version). https://www.kaggle.com/datasets/huthayfahodeb/nih-chest-x-rays-bbox-version, 2023. Accessed: 2025-05-22.
- Tzu-Ming Harry Hsu, Hang Qi, and Matthew Brown. Measuring the effects of non-identical data distribution for federated visual classification, 2019. URL https://arxiv.org/abs/1909.06335.
- Ruida Hu, Chao Peng, Xinchen Wang, and Cuiyun Gao. An Ilm-based agent for reliable docker environment configuration. *arXiv preprint arXiv:2502.13681*, 2025.
- Dong Huang, Jie M. Zhang, Michael Luck, Qingwen Bu, Yuhao Qing, and Heming Cui. AgentCoder:
  Multi-Agent-based Code Generation with Iterative Testing and Optimisation, May 2024a. URL
  http://arxiv.org/abs/2312.13010. arXiv:2312.13010 [cs].
- Qian Huang, Jian Vora, Percy Liang, and Jure Leskovec. MLAgentBench: Evaluating Language
  Agents on Machine Learning Experimentation. In Forty-first International Conference on Machine
  Learning, June 2024b. URL https://openreview.net/forum?id=1Fs1LvjYQW.
- IMT Kaggle Team. Dental radiography dataset. https://www.kaggle.com/datasets/imtkaggleteam/dental-radiography/data, 2023. Accessed: 2025-05-22.
- Yoshitaka Inoue, Tianci Song, and Tianfan Fu. Drugagent: Explainable drug repurposing agent with large language model-based reasoning. *arXiv preprint arXiv:2408.13378*, 2024.
- Md Ashraful Islam, Mohammed Eunus Ali, and Md Rizwan Parvez. Codesim: Multi-agent code generation and problem solving through simulation-driven planning and debugging. *arXiv preprint arXiv:2502.05664*, 2025.
- Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. LiveCodeBench: Holistic and Contamination Free Evaluation of Large Language Models for Code, June 2024. URL http://arxiv.org/abs/2403.07974. arXiv:2403.07974 [cs].
- Naman Jain, Jaskirat Singh, Manish Shetty, Liang Zheng, Koushik Sen, and Ion Stoica. R2e-gym: Procedural environments and hybrid verifiers for scaling open-weights swe agents. *arXiv preprint arXiv:2504.07164*, 2025.
- Yixing Jiang, Kameron C Black, Gloria Geng, Danny Park, James Zou, Andrew Y Ng, and Jonathan H Chen. Medagentbench: A realistic virtual ehr environment to benchmark medical llm agents. *arXiv* preprint arXiv:2501.14654, 2025.
- Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. SWE-bench: Can Language Models Resolve Real-World GitHub Issues?, April 2024. URL http://arxiv.org/abs/2310.06770. arXiv:2310.06770 [cs].
- Liqiang Jing, Zhehui Huang, Xiaoyang Wang, Wenlin Yao, Wenhao Yu, Kaixin Ma, Hongming Zhang, Xinya Du, and Dong Yu. DSBench: How Far Are Data Science Agents to Becoming Data Science

977

978

979

980

991

1013

```
Experts?, September 2024. URL http://arxiv.org/abs/2409.07703. arXiv:2409.07703
```

Lin Justin. Train val test tcga coad msi mss. https://www.kaggle.com/datasets/linjustin/train-val-test-tcga-coad-msi-mss, 2024. Accessed: 2024-05-22.

- Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In *International conference on machine learning*. PMLR, 2020.
- Jeremy Kawahara, Sara Daneshvar, Giuseppe Argenziano, and Ghassan Hamarneh. Seven-point checklist and skin lesion classification using multitask multimodal neural nets. *IEEE journal of biomedical and health informatics*, 23(2):538–546, 2018.
- Yubin Kim, Chanwoo Park, Hyewon Jeong, Yik Siu Chan, Xuhai Xu, Daniel McDuff, Hyeonhoon Lee, Marzyeh Ghassemi, Cynthia Breazeal, Hae Park, et al. Mdagents: An adaptive collaboration of llms for medical decision-making. *Advances in Neural Information Processing Systems*, 37: 79410–79452, 2024.
- Felipe Kitamura. Spr x-ray age and gender dataset. https://www.kaggle.com/datasets/felipekitamura/spr-x-ray-age-and-gender-dataset, 2022a. Accessed: 2025-05-22.
- Felipe Kitamura. Unifesp x-ray bodypart classification. https://www.kaggle.com/datasets/felipekitamura/unifesp-xray-bodypart-classification, 2022b. Accessed: 2025-05-22.
- Sourabh Kumar. Breast cancer histopathology. https://www.kaggle.com/datasets/sourabhkumar29/breast-cancer-histopathology, 2024. Accessed: 2024-05-22.
- Shrinidhi Kumbhar, Venkatesh Mishra, Kevin Coutinho, Divij Handa, Ashif Iquebal, and Chitta Baral. Hypothesis generation for materials discovery and design using goal-driven and constraint-guided llm agents. *arXiv preprint arXiv:2501.13299*, 2025.
- The Red Lad. Pannuke dataset (experimental). https://www.kaggle.com/datasets/theredlad/pannuke-dataset-experimental-data, 2024. Accessed: 2024-05-22.
- Pamela J LaMontagne, Tammie LS Benzinger, John C Morris, Sarah Keefe, Russ Hornbeck, Chengjie Xiong, Elizabeth Grant, Jason Hassenstab, Krista Moulder, Andrei G Vlassenko, et al. OASIS-3: longitudinal neuroimaging, clinical, and cognitive dataset for normal aging and Alzheimer disease. *MedRxiv*, pp. 1–37, 2019.
- Langgraph. Langgraph documentation, 2025. URL https://www.langchain.com/langgraph. Accessed: 2025-05-16.
- Gisang Lee, Sangwoo Park, Junyoung Park, Andrew Chung, Sieun Park, Yoonah Park, Byungju Kim, and Min-gyu Cho. Expanding search space with diverse prompting agents: An efficient sampling approach for llm mathematical reasoning. *arXiv preprint arXiv:2410.09780*, 2024.
- Jingoo Lee, Kyungho Lim, Young-Chul Jung, and Byung-Hoon Kim. Psyche: A multi-faceted patient simulation framework for evaluation of psychiatric assessment conversational agents. *arXiv* preprint arXiv:2501.01594, 2025.
- Tahir Lee. Uterine leiomyosarcoma histopathology. https://www.kaggle.com/datasets/tahirlee/uterine-leiomyosarcoma-histopathology, 2024. Accessed: 2024-05-22.
- Bin Lei, Yi Zhang, Shan Zuo, Ali Payani, and Caiwen Ding. Macm: Utilizing a multi-agent system for condition mining in solving complex mathematical problems. *arXiv* preprint arXiv:2404.04735, 2024.
- Binxu Li, Tiankai Yan, Yuanting Pan, Jie Luo, Ruiyang Ji, Jiayuan Ding, Zhe Xu, Shilong Liu, Haoyu Dong, Zihao Lin, et al. Mmedagent: Learning to use medical tools with multi-modal agent. In *EMNLP (Findings)*, 2024a.

1046

- Daliang Li and Junpu Wang. Fedmd: Heterogenous federated learning via model distillation, 2019. URL https://arxiv.org/abs/1910.03581.
- Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. Camel: Communicative agents for" mind" exploration of large language model society. *Advances in Neural Information Processing Systems*, 36:51991–52008, 2023a.
- Long Li, Weiwen Xu, Jiayan Guo, Ruochen Zhao, Xingxuan Li, Yuqian Yuan, Boqiang Zhang, Yuming Jiang, Yifei Xin, Ronghao Dang, et al. Chain of ideas: Revolutionizing research via novel idea development with llm agents. *arXiv preprint arXiv:2410.13185*, 2024b.
- Qinbin Li, Bingsheng He, and Dawn Song. Model-contrastive federated learning, 2021a. URL https://arxiv.org/abs/2103.16257.
- Qinbin Li, Zeyi Wen, Zhaomin Wu, Sixu Hu, Naibo Wang, Yuan Li, Xu Liu, and Bingsheng He. A survey on federated learning systems: Vision, hype and reality for data privacy and protection. *IEEE Transactions on Knowledge and Data Engineering*, 35(4):3347–3366, 2021b.
- Qinbin Li, Bingsheng He, and Dawn Song. Adversarial collaborative learning on non-IID features. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 19504–19526. PMLR, 23–29 Jul 2023b. URL https://proceedings.mlr.press/v202/li23j.html.
- 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, 2020a.
- Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. *Proceedings of Machine learning and systems*, 2: 429–450, 2020b.
- Tian Li, Shengyuan Hu, Ahmad Beirami, and Virginia Smith. Ditto: Fair and robust federated learning through personalization, 2021c. URL https://arxiv.org/abs/2012.04221.
- Vincent Li, Yule Fu, Tim Knappe, Kevin Han, and Kevin Zhu. Automating mathematical proof generation using large language model agents and knowledge graphs. *arXiv preprint arXiv:2503.11657*, 2025.
- Xiaoxiao Li, Meirui Jiang, Xiaofei Zhang, Michael Kamp, and Qi Dou. Fedbn: Federated learning on non-iid features via local batch normalization, 2021d. URL https://arxiv.org/abs/2102.07623.
- Paul Pu Liang, Terrance Liu, Liu Ziyin, Nicholas B. Allen, Randy P. Auerbach, David Brent, Ruslan Salakhutdinov, and Louis-Philippe Morency. Think locally, act globally: Federated learning with local and global representations, 2020. URL https://arxiv.org/abs/2001.01523.
- Xun Liang, Jiawei Yang, Yezhaohui Wang, Chen Tang, Zifan Zheng, Simin Niu, Shichao Song, Hanyu Wang, Bo Tang, Feiyu Xiong, et al. Surveyx: Academic survey automation via large language models. *arXiv preprint arXiv:2502.14776*, 2025.
- Sook-Lei Liew, Bethany P Lo, Miranda R Donnelly, Artemis Zavaliangos-Petropulu, Jessica N Jeong, Giuseppe Barisano, Alexandre Hutton, Julia P Simon, Julia M Juliano, Anisha Suri, et al. A large, curated, open-source stroke neuroimaging dataset to improve lesion segmentation algorithms. *Scientific data*, 9(1):320, 2022.
- Ryan Y Lin, Siddhartha Ojha, Kevin Cai, and Maxwell F Chen. Strategic collusion of llm agents: Market division in multi-commodity competitions. *arXiv preprint arXiv:2410.00031*, 2024.
- Ben Liu, Jihan Zhang, Fangquan Lin, Xu Jia, and Min Peng. One size doesn't fit all: A personalized conversational tutoring agent for mathematics instruction. 2025. URL https://arxiv.org/abs/2502.12633.
- Hong-Hsiang Liu and Yi-Wen Liu. Agent-driven large language models for mandarin lyric generation. In 2024 27th Conference of the Oriental COCOSDA International Committee for the Co-ordination

1104

1114

- 1080 and Standardisation of Speech Databases and Assessment Techniques (O-COCOSDA), pp. 1–6. IEEE, 1081
- Renpu Liu, Cong Shen, and Jing Yang. Federated representation learning in the under-parameterized 1083 regime, 2024a. URL https://arxiv.org/abs/2406.04596. 1084
- Siyi Liu, Chen Gao, and Yong Li. Large Language Model Agent for Hyper-Parameter Optimization, 1086 February 2024b. URL http://arxiv.org/abs/2402.01881. arXiv:2402.01881 [cs].
- 1087 Dental disease panoramic detection dataset. https://www.kaggle.com/ 1088 Lokisilvres. datasets/lokisilvres/dental-disease-panoramic-detection-dataset, 1089 2023. Accessed: 2025-05-22. 1090
- 1091 Wang Lu, Jindong Wang, Yiqiang Chen, Xin Qin, Renjun Xu, Dimitrios Dimitriadis, and Tao 1092 Qin. Personalized federated learning with adaptive batchnorm for healthcare, 2022. URL https: 1093 //arxiv.org/abs/2112.00734.
- 1094 Yuxuan Lu, Bingsheng Yao, Hansu Gu, Jing Huang, Jessie Wang, Laurence Li, Jiri Gesi, Qi He, Toby 1095 Jia-Jun Li, and Dakuo Wang. Uxagent: An llm agent-based usability testing framework for web 1096 design. arXiv preprint arXiv:2502.12561, 2025. 1097
- 1098 Junyu Luo, Weizhi Zhang, Ye Yuan, Yusheng Zhao, Junwei Yang, Yiyang Gu, Bohan Wu, Binqi Chen, Ziyue Qiao, Qingqing Long, et al. Large language model agent: A survey on methodology, 1099 applications and challenges. arXiv preprint arXiv:2503.21460, 2025. 1100
- 1101 Mi Luo, Fei Chen, Dapeng Hu, Yifan Zhang, Jian Liang, and Jiashi Feng. No fear of heterogeneity: 1102 Classifier calibration for federated learning with non-iid data, 2021. URL https://arxiv.org/ 1103 abs/2106.05001.
- Tianmi Ma, Jiawei Du, Wenxin Huang, Wenjie Wang, Liang Xie, Xian Zhong, and Joey Tianyi Zhou. 1105 Llm knows geometry better than algebra: Numerical understanding of llm-based agents in a trading 1106 arena. arXiv preprint arXiv:2502.17967, 2025. 1107
- 1108 Xiaosong Ma, Jie Zhang, Song Guo, and Wenchao Xu. Layer-wised model aggregation for person-1109 alized federated learning. In Proceedings of the IEEE/CVF Conference on Computer Vision and 1110 Pattern Recognition (CVPR), pp. 10092–10101, June 2022.
- 1111 Nikita Manaenkov. Annotated x-ray angiography dataset. https://www.kagqle.com/ datasets/nikitamanaenkov/annotated-x-ray-angiography-dataset, 2023. 1113 Accessed: 2025-05-22.
- Othmane Marfoq, Giovanni Neglia, Aurélien Bellet, Laetitia Kameni, and Richard Vidal. Federated 1115 multi-task learning under a mixture of distributions, 2022. URL https://arxiv.org/abs/ 1116 2108.10252. 1117
- 1118 Antonis Maronikolakis, Ana Peleteiro Ramallo, Weiwei Cheng, and Thomas Kober. What should i 1119 wear to a party in a greek taverna? evaluation for conversational agents in the fashion domain. arXiv 1120 preprint arXiv:2408.08907, 2024.
- 1121 Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. 1122 Communication-efficient learning of deep networks from decentralized data. In Artificial intelligence 1123 and statistics. PMLR, 2017. 1124
- 1125 Kai Mei, Xi Zhu, Wujiang Xu, Wenyue Hua, Mingyu Jin, Zelong Li, Shuyuan Xu, Ruosong Ye, Yingqiang Ge, and Yongfeng Zhang. Aios: Llm agent operating system. arXiv preprint 1126 arXiv:2403.16971, 2024. 1127
- 1128 Bjoern H Menze, Andras Jakab, Stefan Bauer, Jayashree Kalpathy-Cramer, Keyvan Farahani, Justin 1129 Kirby, Yuliya Burren, Nicole Porz, Johannes Slotboom, Roland Wiest, et al. The multimodal brain 1130 tumor image segmentation benchmark (BRATS). IEEE Transactions on Medical Imaging, 34(10): 1131
- 1993-2024, 2014.
- Meta-LLaMA. Llama 3 technical overview. https://ai.meta.com/llama, 2024. Accessed: 1133 2025-05-22.

- 1134 Meta-LLaMA. Llama 4 model card. https://ai.meta.com/llama, 2025. Accessed:
- 1135 2025-05-22. 1136
- Arindam Mitra, Luciano Del Corro, Guoqing Zheng, Shweti Mahajan, Dany Rouhana, Andres Codas, 1137
- Yadong Lu, Wei-ge Chen, Olga Vrousgos, Corby Rosset, et al. Agentinstruct: Toward generative 1138
- teaching with agentic flows. arXiv preprint arXiv:2407.03502, 2024. 1139
- 1140 Paul Timothy Mooney. Chest x-ray pneumonia. https://www.kaggle.com/datasets/
- 1141 paultimothymooney/chest-xray-pneumonia, 2018. Accessed: 2025-05-22. 1142
- 1143 Paul Timothy Mooney. Breast histopathology images. https://www.kaggle.com/
- datasets/paultimothymooney/breast-histopathology-images, 2024. 1144
- cessed: 2024-05-22. 1145
- 1146 WSH Munirah. Npc-88k-public dataset. https://www.kaggle.com/datasets/ 1147
- wshmunirah/npc-88k-public, 2024. Accessed: 2024-05-22. 1148
- 1149 Andrew MVD. Lung and colon cancer histopathological images. https://www.kagqle.com/
- 1150 datasets/andrewmvd/lung-and-colon-cancer-histopathological-images,
- 2024a. Accessed: 2024-05-22.
- 1152 Andrew MVD. Breast cancer cell segmentation. https://www.kaggle.com/datasets/ 1153
- andrewmvd/breast-cancer-cell-segmentation, 2024b. Accessed: 2024-05-22. 1154
- 1155 Mouheb Ben Nasr and Yassine Hachaïchi. Reinforcement learning agent for client selection in
- 1156 federated llms.

1167

1172

1178

1183

- A. Tuan Nguyen, Philip Torr, and Ser Nam Lim. Fedsr: A simple and effec-1158
- tive domain generalization method for federated learning. In S. Koyejo, S. Mo-1159
- hamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), Advances in Neural In-1160 formation Processing Systems, volume 35, pp. 38831–38843. Curran Associates, Inc.,
- 1161 URL https://proceedings.neurips.cc/paper files/paper/2022/
- 1162 file/fd946a6c99541fddc3d64a3ea39a1bc2-Paper-Conference.pdf. 1163
- 1164 Dinh C Nguyen, Quoc-Viet Pham, Pubudu N Pathirana, Ming Ding, Aruna Seneviratne, Zihuai Lin,
- 1165 Octavia Dobre, and Won-Joo Hwang. Federated learning for smart healthcare: A survey. ACM
- Computing Surveys (CSUR), 55(3):1–37, 2022b. 1166
- Huan Ning, Zhenlong Li, Temitope Akinboyewa, and M Naser Lessani. An autonomous gis agent 1168
- framework for geospatial data retrieval. *International Journal of Digital Earth*, 18(1):2458688, 2025. 1169
- 1170 Jaehoon Oh, Sangmook Kim, and Se-Young Yun. Fedbabu: Towards enhanced representation for
- 1171 federated image classification, 2022. URL https://arxiv.org/abs/2106.06042.
- Izunna Okpala, Ashkan Golgoon, and Arjun Ravi Kannan. Agentic ai systems applied to tasks in 1173
- financial services: Modeling and model risk management crews. arXiv preprint arXiv:2502.05439, 1174
- 2025. 1175
- 1176 OpenAI. Gpt-3.5 turbo overview. https://openai.com/gpt-3-5, 2023a. Accessed: 2025-
- 1177
- OpenAI. Gpt-4 technical report. https://openai.com/research/gpt-4, 2023b. Ac-1179
- cessed: 2025-05-22. 1180
- 1181 OpenAI. Gpt-4.1 technical report. https://openai.com/gpt-4, 2024a. Accessed: 2025-05-
- 1182
- 1184 OpenAI. Gpt-4-turbo performance details. https://platform.openai.com/docs/
- 1185 quides/gpt, 2024b. Accessed: 2025-05-22.
- OpenAI. Gpt-4o: Multimodal reasoning model. https://openai.com/index/gpt-4o, 1187 2024c. Accessed: 2025-05-22.

- OpenAI. Gpt-40 mini: Advancing cost-efficient intelligence. https://openai.com/index/
- gpt-4o-mini-advancing-cost-efficient-intelligence/, July 2024d. Accessed:
- 1190 2025-05-22.
- OpenAI. Openai api documentation, 2025. URL https://openai.com/api/. Accessed:
- 1193 2025-05-16.
- OpenAI. Introducing gpt-4.1 in the api. https://openai.com/index/gpt-4-1/, April
- 1195 2025a. Accessed: 2025-05-22.
- 1196

- OpenAI. Openai o3-mini. https://openai.com/index/openai-o3-mini/, January
- 1198 2025b. Accessed: 2025-05-22.
- OpenAI. Introducing openai o3 and o4-mini. https://openai.com/index/
- introducing-o3-and-o4-mini/, April 2025c. Accessed: 2025-05-22.
- 1201
- 1202 Orvile. Bone fracture dataset. https://www.kaggle.com/datasets/orvile/
- bone-fracture-dataset, 2023a. Accessed: 2025-05-22.
- Orvile. Simple vs comminuted fractures x-ray data. https://www.kaggle.com/datasets/
- orvile/simple-vs-comminuted-fractures-x-ray-data, 2023b. Accessed: 2025-
- 1206 05-22.
- 1208 Orvile. Human bone fractures image dataset (hbfmid). https://www.kaggle.com/
- 1209 datasets/orvile/human-bone-fractures-image-dataset-hbfmid, 2023c. Ac-
- 1210 cessed: 2025-05-22.
- Orvile. Digital knee x-ray images. https://www.kaggle.com/datasets/orvile/
- 1212 digital-knee-x-ray-images, 2023d. Accessed: 2025-05-22.
- 1214 Orvile. Gastric cancer histopathology tissue image
- 1215 dataset. https://www.kaggle.com/datasets/orvile/
- gastric-cancer-histopathology-tissue-image-dataset, 2024. Accessed:
- 1217 2024-05-22.
- Andre GC Pacheco, Gustavo R Lima, Amanda S Salomão, Breno Krohling, Igor P Biral, Gabriel G
- de Angelo, Fábio CR Alves Jr, José GM Esgario, Alana C Simora, Pedro BC Castro, et al. PAD-UFES-
- 1220 20: A skin lesion dataset composed of patient data and clinical images collected from smartphones.
- 1221 Data in Brief, 32:1–10, 2020.
- Jiayi Pan, Xingyao Wang, Graham Neubig, Navdeep Jaitly, Heng Ji, Alane Suhr, and Yizhe Zhang.
- Training software engineering agents and verifiers with swe-gym. arXiv preprint arXiv:2412.21139,
- 1224 1225 2024.
- 1226 Emilio Paspuel. Deep learning driven diagnosis of humerus frac-
- 1227 tures. https://www.kaggle.com/datasets/emiliopaspuel/
- 1228 deeplearning-driven-diagnosis-of-humerus-fractures, 2024. Accessed:
- 1229 2025-05-22.
- Bjarne Pfitzner, Nico Steckhan, and Bert Arnrich. Federated learning in a medical context: A
- systematic literature review. *ACM Transactions on Internet Technology (TOIT)*, 21(2):1–31, 2021.
- 1233 Bits N Pieces. Ovarian cancer and subtypes dataset histopathol-
- 1234 ogy. https://www.kaggle.com/datasets/bitsnpieces/
- ovarian-cancer-and-subtypes-dataset-histopathology, 2024. Accessed:
- 1236 2024-05-22.

- Jianing Qiu, Kyle Lam, Guohao Li, Amish Acharya, Tien Yin Wong, Ara Darzi, Wu Yuan, and Eric J
- Topol. Llm-based agentic systems in medicine and healthcare. *Nature Machine Intelligence*, 6(12):
- 1239 1418–1420, 2024.
- Tawsifur Rahman. Covid-19 radiography database. https://www.kaggle.com/datasets/tawsifurrahman/covid19-radiography-database, 2020. Accessed: 2025-05-22.

- Tawsifur Rahman. Aseptic loose hip implant x-ray database. https://www.kaggle.com/
- datasets/tawsifurrahman/aseptic-loose-hip-implant-xray-database,
- 1244 2022. Accessed: 2025-05-22.

1254

1272

1277

1283

1286

- Suraj Rajendran, Jihad S Obeid, Hamidullah Binol, Kristie Foley, Wei Zhang, Philip Austin, Joey
- Brakefield, Metin N Gurcan, and Umit Topaloglu. Cloud-based federated learning implementation
- across medical centers. JCO Clinical Cancer Informatics, 5:1–11, 2021.
- RANZCR. Ranzer clip catheter and line classification. https://www.kaggle.com/c/
- ranzcr-clip-catheter-line-classification, 2021. Accessed: 2025-05-22.
- Reasat. Histo image-text. https://www.kaggle.com/datasets/reasat/histo-img-text, 2024. Accessed: 2024-05-22.
- Sashank Reddi, Zachary Charles, Manzil Zaheer, Zachary Garrett, Keith Rush, Jakub Konečný,
- Sanjiv Kumar, and H. Brendan McMahan. Adaptive federated optimization, 2021. URL https:
- //arxiv.org/abs/2003.00295.
- Nicola Rieke, Jonny Hancox, Wenqi Li, Fausto Milletari, Holger R Roth, Shadi Albarqouni, Spyridon
- Bakas, Mathieu N Galtier, Bennett A Landman, Klaus Maier-Hein, et al. The future of digital health
- with federated learning. *NPJ digital medicine*, 3(1):119, 2020.
- B Madushani Rodrigo. Fracture multi-region x-ray data. https://www.kaggle.com/
- datasets/bmadushanirodrigo/fracture-multi-region-x-ray-data, 2022.
- 1264 Accessed: 2025-05-22.
- Daniel Rose, Chia-Chien Hung, Marco Lepri, Israa Alqassem, Kiril Gashteovski, and Carolin
- Lawrence. Meddxagent: A unified modular agent framework for explainable automatic differential
- diagnosis. arXiv preprint arXiv:2502.19175, 2025.
- Felix Sattler, Klaus-Robert Müller, and Wojciech Samek. Clustered federated learning: Model-
- agnostic distributed multi-task optimization under privacy constraints, 2019. URL https://
- 1271 arxiv.org/abs/1910.01991.
- Samuel Schmidgall and Michael Moor. Agentrxiv: Towards collaborative autonomous research.
- arXiv preprint arXiv:2503.18102, 2025.
- Dominik Schmidt, Zhengyao Jiang, and Yuxiang Wu. Introducing Weco AIDE, April 2024. URL
- https://www.weco.ai/blog/technical-report.
- Jonathan Scott, Hossein Zakerinia, and Christoph H. Lampert. Peffl: Personalized federated learning
- by learning to learn, 2025. URL https://arxiv.org/abs/2306.05515.
- 1280
- Constantin Seibold. Anatomy in chest x-rays (pax-ray++). https://www.kaggle.com/datasets/constantinseibold/anatomy-in-chest-x-rays-pax-ray, 2023. Ac-
- cessed: 2025-05-22.
- Aviv Shamsian, Aviv Navon, Ethan Fetaya, and Gal Chechik. Personalized federated learning using
- hypernetworks, 2021. URL https://arxiv.org/abs/2103.04628.
- Zitong Shi, Guancheng Wan, Wenke Huang, Guibin Zhang, Jiawei Shao, Mang Ye, and Carl Yang.
- Privacy-enhancing paradigms within federated multi-agent systems. In ICML 2025 Workshop on
- Collaborative and Federated Agentic Workflows.
- Jayanthi Siyaswamy, SR Krishnadas, Gopal Datt Joshi, Madhulika Jain, and A Ujjwaft Syed Tabish.
- Drishti-gs: Retinal image dataset for optic nerve head (onh) segmentation. In 2014 IEEE 11th
- international symposium on biomedical imaging (ISBI), pp. 53–56. IEEE, 2014.
- 1294 Spritan1. Yolo annotated chestxray 8 object detection. https://www.kaggle.com/
- datasets/spritan1/yolo-annotated-chestxray-8-object-detection, 2023. Accessed: 2025-05-22.

1316

- Ian Steenstra, Farnaz Nouraei, and Timothy W Bickmore. Scaffolding empathy: Training counselors with simulated patients and utterance-level performance visualizations. *arXiv preprint arXiv:2502.18673*, 2025.
- Osama H Taher. Heel dataset. https://www.kaggle.com/datasets/osamahtaher/heel-dataset, 2023. Accessed: 2025-05-22.
- Jiahao Tan and Xinpeng Wang. FL-bench: A federated learning benchmark for solving image classification tasks. URL https://github.com/KarhouTam/FL-bench.
- Jiahao Tan, Yipeng Zhou, Gang Liu, Jessie Hui Wang, and Shui Yu. pfedsim: Similarity-aware model aggregation towards personalized federated learning. *arXiv preprint arXiv:2305.15706*, 2023.
- Yue Tan, Guodong Long, Lu Liu, Tianyi Zhou, Qinghua Lu, Jing Jiang, and Chengqi Zhang. Fedproto: Federated prototype learning across heterogeneous clients, 2022. URL https://arxiv.org/abs/2105.00243.
- Xiangru Tang, Yuliang Liu, Zefan Cai, Yanjun Shao, Junjie Lu, Yichi Zhang, Zexuan Deng, Helan Hu, Kaikai An, Ruijun Huang, Shuzheng Si, Sheng Chen, Haozhe Zhao, Liang Chen, Yan Wang, Tianyu Liu, Zhiwei Jiang, Baobao Chang, Yin Fang, Yujia Qin, Wangchunshu Zhou, Yilun Zhao, Arman Cohan, and Mark Gerstein. ML-Bench: Evaluating Large Language Models and Agents for Machine Learning Tasks on Repository-Level Code, August 2024. URL http://arxiv.org/abs/2311.09835.
- Xiangru Tang, Tianyu Hu, Muyang Ye, Yanjun Shao, Xunjian Yin, Siru Ouyang, Wangchunshu Zhou, Pan Lu, Zhuosheng Zhang, Yilun Zhao, et al. Chemagent: Self-updating library in large language models improves chemical reasoning. *arXiv preprint arXiv:2501.06590*, 2025.
- Tapendu. Chest x-ray dataset for tuberculosis segmentation. https://www.kaggle.com/datasets/iamtapendu/chest-x-ray-lungs-segmentation, 2023a. Accessed: 2025-05-22.
- Tapendu. Rsna pneumonia processed dataset. https://www.kaggle.com/datasets/iamtapendu/rsna-pneumonia-processed-dataset, 2023b. Accessed: 2025-05-22.
- Mr. Tejas. Covid-19 and normal x-ray dataset (balanced). https://www.kaggle.com/datasets/mrtejas/covid-19-and-normal-x-ray-dataset-balanced, 2022. Accessed: 2025-05-22.
- Aryamaan Thakur. Rsna breast cancer detection roi 1024. https://www.kaggle.com/datasets/aryamaanthakur/rsna-breast-cancer-detection-roi-1024, 2024. Accessed: 2025-05-22.
- TrainingDataPro. Chest x-ray 17 diseases. https://www.kaggle.com/datasets/trainingdatapro/chest-xray-17-diseases, 2023. Accessed: 2025-05-22.
- TruthIsNeverLinear. Brecahad. https://www.kaggle.com/datasets/truthisneverlinear/brecahad, 2024. Accessed: 2024-05-22.
- Preet Viradiya. Covid-19 radiography dataset extended. https://www.kaggle.com/datasets/preetviradiya/covid19-radiography-dataset, 2023. Accessed: 2025-05-22.
- Xiangpeng Wan, Haicheng Deng, Kai Zou, and Shiqi Xu. Enhancing the efficiency and accuracy of underlying asset reviews in structured finance: The application of multi-agent framework. *arXiv* preprint arXiv:2405.04294, 2024.
- Cunshi Wang, Xinjie Hu, Yu Zhang, Xunhao Chen, Pengliang Du, Yiming Mao, Rui Wang, Yuyang Li, Ying Wu, Hang Yang, et al. Starwhisper telescope: Agent-based observation assistant system to approach ai astrophysicist. *arXiv preprint arXiv:2412.06412*, 2024a.
- Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. *arXiv* preprint arXiv:2305.16291, 2023.

- 1350 Jian Wang, Yinpei Dai, Yichi Zhang, Ziqiao Ma, Wenjie Li, and Joyce Chai. Training turn-by-turn 1351 verifiers for dialogue tutoring agents: The curious case of llms as your coding tutors. arXiv preprint
- 1352 arXiv:2502.13311, 2025a.

1383

- Jiuniu Wang, Zehua Du, Yuyuan Zhao, Bo Yuan, Kexiang Wang, Jian Liang, Yaxi Zhao, Yihen Lu, 1354 Gengliang Li, Junlong Gao, et al. Aesopagent: Agent-driven evolutionary system on story-to-video 1355
- production. arXiv preprint arXiv:2403.07952, 2024b. 1356
- 1357 Ruida Wang, Rui Pan, Yuxin Li, Jipeng Zhang, Yizhen Jia, Shizhe Diao, Renjie Pi, Junjie Hu, and 1358 Tong Zhang. Ma-lot: Multi-agent lean-based long chain-of-thought reasoning enhances formal
- 1359 theorem proving. arXiv preprint arXiv:2503.03205, 2025b. 1360
- Xin Wang, Yifan Zhang, Xiaojing Zhang, Longhui Yu, Xinna Lin, Jindong Jiang, Bin Ma, and 1361 Kaicheng Yu. Patentagent: Intelligent agent for automated pharmaceutical patent analysis. arXiv 1362 preprint arXiv:2410.21312, 2024c. 1363
- 1364 Zhizheng Wang, Qiao Jin, Chih-Hsuan Wei, Shubo Tian, Po-Ting Lai, Qingqing Zhu, Chi-Ping Day, 1365 Christina Ross, and Zhiyong Lu. Geneagent: self-verification language agent for gene set knowledge 1366 discovery using domain databases. arXiv preprint arXiv:2405.16205, 2024d.
- 1367 Ziyue Wang, Junde Wu, Chang Han Low, and Yueming Jin. Medagent-pro: Towards multi-modal 1368 evidence-based medical diagnosis via reasoning agentic workflow. arXiv preprint arXiv:2503.18968, 1369 2025c. 1370
- 1371 Abbi Ward, Jimmy Li, Julie Wang, Sriram Lakshminarasimhan, Ashley Carrick, Bilson Campana, 1372 Jay Hartford, Pradeep K. Sreenivasaiah, Tiya Tiyasirisokchai, Sunny Virmani, Renee Wong, Yossi 1373 Matias, Greg S. Corrado, Dale R. Webster, Margaret Ann Smith, Dawn Siegel, Steven Lin, Justin Ko, 1374 Alan Karthikesalingam, Christopher Semturs, and Pooja Rao. Creating an empirical dermatology dataset through crowdsourcing with web search advertisements. JAMA Network Open, 7(11): 1375 e2446615-e2446615, 11 2024. ISSN 2574-3805. 10.1001/jamanetworkopen.2024.46615. URL 1376 https://doi.org/10.1001/jamanetworkopen.2024.46615. 1377
- 1378 Robert Wasenmüller, Kevin Hilbert, and Christoph Benzmüller. Script-based dialog policy planning 1379 for llm-powered conversational agents: A basic architecture for an" ai therapist". arXiv preprint 1380 arXiv:2412.15242, 2024.
- 1381 Panlong Wu, Kangshuo Li, Junbao Nan, and Fangxin Wang. Federated in-context llm agent learning. arXiv preprint arXiv:2412.08054, 2024.
- 1384 Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun 1385 Zhang, Shaokun Zhang, Jiale Liu, et al. Autogen: Enabling next-gen llm applications via multi-agent 1386 conversation. In ICLR 2024 Workshop on Large Language Model (LLM) Agents.
- Wenbei Xie, Donglin Liu, Haoran Yan, Wenjie Wu, and Zongyang Liu. Mathlearner: A large 1388 language model agent framework for learning to solve mathematical problems. arXiv preprint 1389 arXiv:2408.01779, 2024. 1390
- 1391 Ancheng Xu, Di Yang, Renhao Li, Jingwei Zhu, Minghuan Tan, Min Yang, Wanxin Qiu, Mingchen 1392 Ma, Haihong Wu, Bingyu Li, et al. Autocbt: An autonomous multi-agent framework for cognitive 1393 behavioral therapy in psychological counseling. arXiv preprint arXiv:2501.09426, 2025a. 1394
- Zhenran Xu, Longyue Wang, Jifang Wang, Zhouyi Li, Senbao Shi, Xue Yang, Yiyu Wang, Baotian 1395 Hu, Jun Yu, and Min Zhang. Filmagent: A multi-agent framework for end-to-end film automation in 1396 virtual 3d spaces. arXiv preprint arXiv:2501.12909, 2025b. 1397
- 1398 Jian Yang, Wei Zhang, Jiaxi Yang, Yibo Miao, Shanghaoran Quan, Zhenhe Wu, Qiyao Peng, Liqun 1399 Yang, Tianyu Liu, Zeyu Cui, et al. Multi-agent collaboration for multilingual code instruction tuning. 1400 arXiv preprint arXiv:2502.07487, 2025a. 1401
- Xiyuan Yang, Wenke Huang, and Mang Ye. Fedas: Bridging inconsistency in personalized federated 1402 learning. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition 1403 (CVPR), pp. 11986–11995, June 2024.

- Yizhe Yang, Palakorn Achananuparp, Heyan Huang, Jing Jiang, Kit Phey Leng, Nicholas Gabriel Lim, Cameron Tan Shi Ern, and Ee-peng Lim. Cami: A counselor agent supporting motivational
- interviewing through state inference and topic exploration. *arXiv preprint arXiv:2502.02807*, 2025b.
- 1408 Yuzhe Yang, Yifei Zhang, Minghao Wu, Kaidi Zhang, Yunmiao Zhang, Honghai Yu, Yan Hu, and
- Benyou Wang. Twinmarket: A scalable behavioral and socialsimulation for financial markets. *arXiv*
- 1410 preprint arXiv:2502.01506, 2025c.
- Zhiqin Yang, Yonggang Zhang, Yu Zheng, Xinmei Tian, Hao Peng, Tongliang Liu, and Bo Han.
- Fedfed: Feature distillation against data heterogeneity in federated learning, 2023. URL https:
- //arxiv.org/abs/2310.05077.
- Beibei Yu, Tao Shen, Hongbin Na, Ling Chen, and Denqi Li. Mineagent: Towards remote-sensing
- mineral exploration with multimodal large language models. arXiv preprint arXiv:2412.17339,
- 1417 2024a.

1430

1445

- Dingyao Yu, Kaitao Song, Peiling Lu, Tianyu He, Xu Tan, Wei Ye, Shikun Zhang, and Jiang Bian.
- Musicagent: An ai agent for music understanding and generation with large language models. arXiv
- 1421 preprint arXiv:2310.11954, 2023.
- Yangyang Yu, Zhiyuan Yao, Haohang Li, Zhiyang Deng, Yuechen Jiang, Yupeng Cao, Zhi Chen,
- Jordan Suchow, Zhenyu Cui, Rong Liu, et al. Fincon: A synthesized llm multi-agent system with
- conceptual verbal reinforcement for enhanced financial decision making. Advances in Neural
- 1425 *Information Processing Systems*, 37:137010–137045, 2024b.
- Murong Yue, Wenhan Lyu, Wijdane Mifdal, Jennifer Suh, Yixuan Zhang, and Ziyu Yao. Mathvc:
- An Ilm-simulated multi-character virtual classroom for mathematics education. arXiv preprint
- 1429 arXiv:2404.06711, 2024.
- Taedong Yun, Eric Yang, Mustafa Safdari, Jong Ha Lee, Vaishnavi Vinod Kumar, S Sara Mahdavi,
- Jonathan Amar, Derek Peyton, Reut Aharony, Andreas Michaelides, et al. Sleepless nights, sugary
- days: Creating synthetic users with health conditions for realistic coaching agent interactions. arXiv
- 1433 preprint arXiv:2502.13135, 2025.
- Jingying Zeng, Hui Liu, Zhenwei Dai, Xianfeng Tang, Chen Luo, Samarth Varshney, Zhen Li, and
- Qi He. Cite before you speak: Enhancing context-response grounding in e-commerce conversational
- llm-agents. *arXiv preprint arXiv:2503.04830*, 2025.
- Huan Zhang, Yu Song, Ziyu Hou, Santiago Miret, and Bang Liu. Honeycomb: A flexible llm-based
- agent system for materials science. arXiv preprint arXiv:2409.00135, 2024a.
- Jianqing Zhang, Yang Hua, Hao Wang, Tao Song, Zhengui Xue, Ruhui Ma, and Haibing Guan.
- Fedala: Adaptive local aggregation for personalized federated learning. *Proceedings of the*
- AAAI Conference on Artificial Intelligence, 37(9):11237–11244, June 2023. ISSN 2159-5399.
- 10.1609/aaai.v37i9.26330. URL http://dx.doi.org/10.1609/aaai.v37i9.26330.
- Jie Zhang, Zhiqi Li, Bo Li, Jianghe Xu, Shuang Wu, Shouhong Ding, and Chao Wu. Federated
- learning with label distribution skew via logits calibration, 2022. URL https://arxiv.org/
- abs/2209.00189.
- Michael Zhang, Karan Sapra, Sanja Fidler, Serena Yeung, and Jose M. Alvarez. Personalized
- federated learning with first order model optimization, 2021. URL https://arxiv.org/abs/
- 1451 2012.08565.
- Ran Zhang and Steffen Eger. Llm-based multi-agent poetry generation in non-cooperative environ-
- ments. arXiv preprint arXiv:2409.03659, 2024.
- 1456 Yiqun Zhang, Xiaocui Yang, Xiaobai Li, Siyuan Yu, Yi Luan, Shi Feng, Daling Wang, and Yifei
- Zhang. Psydraw: A multi-agent multimodal system for mental health screening in left-behind children. *arXiv preprint arXiv:2412.14769*, 2024b.

Yuntong Zhang, Haifeng Ruan, Zhiyu Fan, and Abhik Roychoudhury. AutoCodeRover: Au-tonomous Program Improvement, July 2024c. URL http://arxiv.org/abs/2404.05427. arXiv:2404.05427 [cs]. Mingkai Zheng, Xiu Su, Shan You, Fei Wang, Chen Qian, Chang Xu, and Samuel Albanie. Can GPT-4 Perform Neural Architecture Search?, August 2023. URL http://arxiv.org/abs/ 2304.10970. arXiv:2304.10970 [cs]. Yuan Zhou, Peng Zhang, Mengya Song, Alice Zheng, Yiwen Lu, Zhiheng Liu, Yong Chen, and Zhaohan Xi. Zodiac: A cardiologist-level llm framework for multi-agent diagnostics. arXiv preprint arXiv:2410.02026, 2024. Zhuangdi Zhu, Junyuan Hong, and Jiayu Zhou. Data-free knowledge distillation for heterogeneous federated learning, 2021. URL https://arxiv.org/abs/2105.10056. 

A	PPEN	IDIX	
C	ONTI	ENTS	
	<u>.</u> .		٠
1	Intr	oduction and Background	1
2	Fed	AgentBench Framework	3
	2.1	Problem Formulation and Overview	3
	2.2	Client Dataset Curation and FL Algorithm Integration	3
	2.3	Federated Agentic Framework Construction	Δ
	2.4	Privacy Preserving and Modular Design	$\epsilon$
3	Exp	eriments and Results	6
	3.1	Implementation and Evaluation Details	$\epsilon$
	3.2	Main Results and Key Insights	7
	3.3	Agent Failure Analysis:	8
4	Con	clusion and Limitation	9
C	ONTI	ENTS OF APPENDIX	29
A	Rela	ated Works	31
	A.1	Federated Learning for Medical Image Analysis	31
	A.2	LLM Agent Applications	31
	A.3	LLM Agents for Machine Learning, Software Engineering, and Federated Learning	32
В	Tool	s and Agents in FedAgentBench Framework	33
	B.1	Collection of Tools accessed by the LLM Agents	33
	B.2	Role-specialized Agents	42
C	Task	ss and Algorithms in FedAgentBench Framework	49
	C.1	Dataset Details	49
	C.2	Sample dataset description files:	66
	C.3	Detecting and Addressing Data Quality Issues for Data Pre-Processing Agent	78
	C.4	Collection of Federated Learning algorithms	79
	C.5	LLMs as the agent core components	88
D	Resi	ults and Discussions	92
	D.1	Discussion on agentic performance in individual healthcare environment	93
	D.2	Discussion on time-efficiency	97
	D.3	Discussion on client selection, reasoning vs non-reasoning models and failure modes:	98
T.	Duis	agy Analysis of Harmonized Labels and Matadata	06

	E.1	Mutual Information Analysis	98
	E.2	Differential Privacy (DP) Proof	99
	E.3	k-Anonymity Analysis	99
	E.4	Privacy-Utility Trade-off	100
F	Broa	ader Social Impact	100
G	LLN	M Usage:	100

# A RELATED WORKS

# A.1 FEDERATED LEARNING FOR MEDICAL IMAGE ANALYSIS

Existing research on federated learning (FL) in medical image analysis has primarily focused on the development of machine learning algorithms to address technical challenges, such as data distribution shift, statistical and system heterogeneity, and communication efficiency (Antunes et al., 2022; Rajendran et al., 2021; Nguyen et al., 2022b; Pfitzner et al., 2021; Rieke et al., 2020). These efforts have produced a wide range of methods tailored for robust and scalable training under diverse and decentralized medical data environments. However, despite these advances, a significant barrier to real-world deployment persists: the complex set of operational and human-in-the-loop challenges encountered in practice.

Notably, existing FL benchmarks and studies rarely account for the intricacies of human factors—such as institutional workflows, task specification, annotation and curation requirements, and the expertise needed to orchestrate the entire FL pipeline across multiple healthcare institutions. These operational hurdles, including coordination among stakeholders, error handling, and workflow reproducibility, often constitute the most substantial obstacles to routine FL adoption in clinical settings.

This paper distinguishes itself from prior work by explicitly modeling and integrating these real-world operational challenges into the benchmarking process. By capturing both the algorithmic and human-centered aspects of FL deployment, our benchmark provides a more comprehensive and realistic evaluation platform. This enables the research community to move beyond algorithm-centric benchmarks and address the "elephant in the room", *i.e.*, the operational bottlenecks that ultimately determine the success or failure of federated learning in medical imaging practice.

#### A.2 LLM AGENT APPLICATIONS

AI agents, powered by large language models (LLMs), autonomous tool use, and decision-making workflows, are rapidly transforming a diverse range of application domains. In **healthcare**, LLM-based agents drive advances in clinical diagnosis Chen et al. (2024); Zhou et al. (2024); Wang et al. (2025c); Rose et al. (2025); Ghezloo et al. (2025); Li et al. (2024a); Jiang et al. (2025); Kim et al. (2024); Fallahpour et al., mental health and therapy Wasenmüller et al. (2024); Du et al. (2024); Zhang et al. (2024b); Lee et al. (2025); Xu et al. (2025a); Yang et al. (2025b); Steenstra et al. (2025); Abbasi et al. (2025), workflow optimization Feng et al. (2025); Yun et al. (2025); Chen et al. (2025c), and pharmaceutical research Wang et al. (2024c); Averly et al. (2025); Inoue et al. (2024). These agents support professionals through transparent reasoning, multi-modal data integration, and interactive, explainable decision support, as well as automated data processing and clinical research acceleration.

In **biomedical and materials science**, agents enhance literature analysis and hypothesis generation Liang et al. (2025); Li et al. (2024b); Schmidgall & Moor (2025); Gottweis et al. (2025), automate gene set knowledge discovery Wang et al. (2024d), and orchestrate complex scientific workflows, including astronomical observation Wang et al. (2024a) and materials design Zhang et al. (2024a); Kumbhar et al. (2025).

The field of **software engineering** benefits from LLM agents for code generation, repair, verification, and environment setup Dong (2025); Jain et al. (2025); Wang et al. (2025a); Chen et al. (2025a); Aggarwal et al. (2025); Chen et al. (2025b); Gholamzadeh Khoee et al. (2025); Hu et al. (2025); Lu et al. (2025); Pan et al. (2024); Yang et al. (2025a); Guo et al. (2025); Islam et al. (2025). These agents leverage specialized architectures, collaborative multi-agent strategies, and benchmarking frameworks for automated programming, debugging, and user experience testing.

In **finance**, AI agents automate structured finance workflows, simulate markets, optimize investment decisions, and manage risk Wan et al. (2024); Yang et al. (2025c); Yu et al. (2024b); Lin et al. (2024); Fatemi & Hu (2024); Han et al. (2024b; 2025); Fatouros et al. (2025); Okpala et al. (2025); Zeng et al. (2025). Multi-agent frameworks enable complex reasoning, robust QA, and the generation of explainable financial reports.

**Synthetic data generation** is advanced through multi-agent orchestration frameworks Mitra et al. (2024), improving post-training data quality and scalability for large language models.

In **chemistry and materials**, agents automate chemical reasoning Cho et al. (2025); Tang et al. (2025), accelerate drug and materials discovery, and enable hypothesis-driven research Zhang et al. (2024a); Kumbhar et al. (2025).

Mathematics education and scientific reasoning have seen the development of multi-agent reasoning and tutoring systems to tackle complex mathematical proofs, theorem proving, and adaptive instruction Lei et al. (2024); Xie et al. (2024); Lee et al. (2024); Deng & Mineiro (2024); Li et al. (2025); Wang et al. (2025b); Yue et al. (2024); Liu et al. (2025); Ma et al. (2025).

In **geospatial science**, agents facilitate autonomous GIS analysis and data retrieval Yu et al. (2024a); Ning et al. (2025), addressing the challenge of spatial reasoning and multi-source data fusion.

The domain of **multimedia and creative industries** is being transformed by AI agents capable of automating film production, music and lyric generation, story-to-video creation, fashion assistance, and poetry composition Xu et al. (2025b); Wang et al. (2024b); Han et al. (2024a); Maronikolakis et al. (2024); Deng et al. (2024); Yu et al. (2023); Zhang & Eger (2024); Liu & Liu (2024). These systems support multi-modal content creation and human-AI co-creation.

Overall, the emergence of LLM-powered agents marks a shift toward highly automated, context-aware, and collaborative AI systems with applications spanning healthcare, science, engineering, finance, education, and the creative arts.

# A.3 LLM AGENTS FOR MACHINE LEARNING, SOFTWARE ENGINEERING, AND FEDERATED LEARNING

The intersection of large language models (LLMs) and autonomous agents has made rapid advancements in machine learning and software engineering. Several works Chen et al. (2021); Hendrycks et al. (2021); Austin et al. (2021); Jain et al. (2024) assess model performance on code generation from natural language instructions. For example: AgentCoder (Huang et al., 2024a) reports 96.3% and 91.8% accuracy on HumanEval and MBPP, respectively. SWE-bench (Jimenez et al., 2024) advances the field by requiring models to resolve real-world pull requests from open-source repositories. Notably, model performance on SWE-bench continues to improve steadily (Zhang et al., 2024c; factory.ai, 2024).

Prior work has also leveraged LLMs for tasks such as hyperparameter optimization (Liu et al., 2024b) and neural architecture design (Zheng et al., 2023). MLAgentBench (Huang et al., 2024b) evaluates agents on 13 Kaggle and custom ML tasks, providing a baseline solution for each and measuring whether agents can achieve at least a 10% improvement. Similarly, ML-Bench (Tang et al., 2024) evaluates an agent's ability to generate code and interact with established ML repositories. AIDE, as reported by Weco AI (Schmidt et al., 2024), surpasses more than 50% of human competitors in Kaggle-style data science contests. DSBench (Jing et al., 2024) also introduces a Kaggle competition benchmark, but, like Weco AI, focuses primarily on data science tasks.

While benchmarking LLM agents for automated machine learning and data science has gained momentum across both academia and industry, all of these operate under the assumption of a centralized, single-site environment, limiting their applicability to the federated learning paradigm, which introduces unique challenges such as distributed data silos, partial observability, and multiparty coordination. Recent works on agentic FL frameworks include in-context learning in FL of LLM agents Wu et al. (2024), reinforcement learning agent for client selection Nasr & Hachaïchi, and privacy enhancing techniques in federated mult-agent systems Shi et al..

In contrast to these works, **FedAgentBench** is designed to address the real-world operational complexities in federated learning workflows by evaluating the agentic capabilities — particularly in high-stakes healthcare settings. **Rather than being "yet another" benchmark, FedAgentBench is motivated by a concrete and pressing need to reduce the human coordination bottlenecks that currently hinder scalable deployment of FL in practice. It provides a realistic testbed for assessing agent autonomy, adaptability, and reasoning in decentralized, privacy-preserving environments.** 

# B TOOLS AND AGENTS IN FEDAGENTBENCH FRAMEWORK

#### B.1 COLLECTION OF TOOLS ACCESSED BY THE LLM AGENTS

The following tools form the operational backbone of the LLM-based agents, enabling tasks such as file inspection, dataset organization, data cleaning, folder manipulation, and federated training orchestration. Corresponding code snippets for all 16 tools can be found in Listing 1.

- read\_files: Reads the content of one or more specified files and returns a dictionary mapping
  file paths to their contents. It supports UTF-8 text files and handles file access errors
  gracefully.
- move\_directory: Moves a source directory (including all files and subfolders) to a new destination.
- 3. **copy\_files**: Copies multiple individual files to specified destination paths. Accepts a mapping of source to destination file paths and ensures target directories are created as needed.
- 4. **write\_file**: Writes a given text string to a specified file path. It creates any missing directories in the path before writing.
- 5. **edit\_file**: Overwrites the contents of a specified file with new content. Used for completely replacing existing file content.
- 6. **run\_script**: Executes a given shell command (typically a Python script) using a secure subprocess or shell tool backend. Returns the result of the command execution.
- 7. **list\_files\_in\_second\_level**: Traverses the second-level entries of a root directory. For each subdirectory or file, it collects and returns metadata including the total number of files and a preview list of file paths (up to 10).
- 8. **preview\_file\_content**: Previews the contents of a CSV, JSON, or TXT file. Returns first 5 rows or entries and summary statistics such as total rows or elements.
- 9. run\_selfclean\_on\_dataset: Runs the data cleaning framework on an image folder to detect and optionally clean near duplicates, off-topic or irrelevant samples, and label errors. It generates internal diagnostic data in CSV format for inspection and removes samples based on a threshold. Within this process, we also achieve normalization and standardization.
- organize\_into\_subfolder: Reads a CSV containing image paths and labels, and organizes the corresponding images into class-specific subfolders within a specified destination directory.
- 11. **copy\_folder**: Copies all contents (files and subfolders) from a source directory to a destination directory. Ensures destination exists and performs a recursive copy.
- 12. **remove\_other\_files**: Recursively removes all non-image files from a directory structure. Keeps standard image formats (e.g., .jpg, .png, .bmp) and deletes all others.
- 13. **list\_folders**: Returns the names of all first-level subdirectories under a specified root directory. Useful for summarizing dataset structure.
- 14. **make\_folder**: Creates a new directory at a specified path. Used to set up target folders during label harmonization or preprocessing.
- 15. **copy\_images**: Copies all image files from a source folder to a specified target folder. Typically used during label harmonization to reorganize class-wise images.
- 16. **run\_federated\_method**: Launches federated learning using a specified algorithm and project directory. Executes a Python script with algorithm-specific parameters and returns algorithm performance.

#### Listing 1: Repository of tools used by LLM Agents

```
1. def read_files(file_paths: list) -> dict:
    """
    Read file contents and return as dictionary.

Args:
    file_paths: List of file paths to read
```

```
1782
1783
           Returns:
1784
               dict: Dictionary with {file_path: file_content} format
1785
           file_contents = {}
1786
1787
           for file_path in file_paths:
1788
               try:
1789
                   with open(file_path, 'r', encoding='utf-8') as file:
1790
                       content = file.read()
                       file_contents[file_path] = content
1791
               except (UnicodeDecodeError, PermissionError, FileNotFoundError)
1792
                   as e:
1793
                   print(f"Cannot read file {file_path}: {e}")
1794
                   file_contents[file_path] = None
1795
           return file_contents
1796
1797
       2. def move_directory(src_dir: str, dest_dir: str) -> str:
1798
1799
           Move source directory and its contents to destination directory,
1800
              creating a new subdirectory
           with the same name as the source directory.
1801
1802
           Aras:
1803
               src_dir: Source directory path (e.g., '/path/to/source/
1804
                   folder_name')
               dest_dir: Parent destination directory path (e.g., '/path/to/dest
1805
1806
                        A new subdirectory named 'folder_name' will be created
1807
1808
1809
           Returns:
               str: Operation result message
1810
1811
           Example:
1812
               If src_dir is '/path/to/source/folder_name' and dest_dir is '/
1813
                   path/to/dest'
1814
               the directory will be moved to '/path/to/dest/folder_name'
1815
           print(f"Running move_directory tool to move from {src_dir} to {
1816
              dest_dir}...")
1817
           try:
1818
               if not os.path.exists(src_dir):
1819
                   return f"Source directory {src_dir} does not exist"
1820
               # Get the source directory name
1821
               src_name = os.path.basename(src_dir.rstrip('/'))
1822
               target_dir = os.path.join(dest_dir, src_name)
1823
1824
               # If destination directory already exists, remove it first
1825
               if os.path.exists(target_dir):
                   shutil.rmtree(target_dir)
1826
1827
               # Move the directory
1828
               shutil.move(src_dir, target_dir)
1829
               return f"Directory {src_dir} has been successfully moved to {
                   target_dir}"
1830
1831
           except Exception as e:
1832
               return f"Error moving directory: {str(e)}"
1833
1834
      3. def copy_files(file_mapping: dict) -> str:
1835
```

```
1836
           Copy multiple files from source paths to destination paths.
1837
1838
          Args:
1839
               file_mapping (dict): A dictionary where keys are source file
                   paths and values are destination file paths.
1840
                   Example: {
1841
                       "/path/to/sourcel.txt": "/path/to/destination1.txt",
1842
                       "/path/to/source2.txt": "/path/to/destination2.txt"
1843
1844
           Returns:
1845
               str: A message indicating the result of the operation.
1846
1847
           print(f"Running copy_files tool to copy {file_mapping}...")
1848
           results = []
           for src, dest in file_mapping.items():
1849
               try:
1850
                   # Check if source file exists
1851
                   if not os.path.exists(src):
1852
                       results.append(f"Source file {src} does not exist.")
1853
                       continue
1854
                   # Create destination directory if it doesn't exist
1855
                   dest_directory = os.path.dirname(dest)
1856
                   if not os.path.exists(dest_directory):
1857
                       os.makedirs(dest_directory)
1858
1859
                   # Copy file
                   shutil.copy2(src, dest)
1860
                   results.append(f"File {src} successfully copied to {dest}")
1861
1862
               except Exception as e:
1863
                   results.append(f"Error copying file {src}: {e}")
1864
           # Return summary of all operations
1865
           return "\n".join(results)
1866
1867
      4. def write_file(content: str, file_path: str) -> None:
1868
          Write a given string of code to a specified file.
1869
1870
           This function creates the necessary directories for the file (if they
1871
               don't exist),
1872
           writes the content to the file, and handles any errors that may occur
1873
               during the process.
1874
          Aras:
1875
               content (str): The code or text you want to write into the file.
1876
               file_path (str): The full path (including filename) where the
1877
                   content will be saved.
1878
1879
          Example:
               write_file('print("Hello World")', 'scripts/hello.py')
1880
1881
          print(f"Running write_file tool to write {file_path}...")
1882
          try:
1883
               os.makedirs(os.path.dirname(file_path), exist_ok=True)
1884
               with open(file_path, 'w', encoding='utf-8') as file:
1885
                   file.write(content)
1886
1887
               print(f"File successfully written to: {file_path}")
1888
          except Exception as e:
               print(f"Error writing file: {e}")
```

```
1890
1891
       5. def edit_file(new_content: str, file_path: str) -> None:
1892
1893
           Completely overwrite a file with new content. The original file
              content will be replaced entirely.
1894
1895
           Aras:
1896
               new_content: Complete content to replace the existing file
1897
                   content. This should be the entire
1898
                            desired content of the file after editing, not just
                                the changes.
1899
               file_path: Path of the file to edit
1900
1901
           Note:
1902
               This function performs a complete overwrite operation. The
                   original content will be lost.
1903
               You must provide the complete desired final content, including
1904
                   both modified and unmodified parts.
1905
1906
           print(f"Running edit_file tool to edit {file_path}...")
1907
           try:
               with open(file_path, 'w', encoding='utf-8') as file:
1908
                   file.write(new_content)
1909
1910
               print(f"File {file_path} successfully edited.")
1911
           except Exception as e:
1912
               print(f"Error editing file: {e}")
1913
       6. def run_script(command: str) -> str:
1914
1915
           Execute shell command
1916
1917
           Args:
               command: Shell command to execute
1918
1919
           Returns:
1920
               str: Command execution result
1921
1922
           cmd_base, script_path = command.strip().split(maxsplit=1)
1923
           # Blindly quote the path
1924
           script_path = f'"{script_path}"'
1925
1926
           # Rebuild the final command
1927
           fixed_command = f"{cmd_base} {script_path}"
1928
           print(f"Executing fixed command: {fixed_command}")
1929
           print("Running run_script tool...")
1930
           shell_tool = ShellTool()
1931
           result = shell_tool.run({
1932
               "commands": [fixed_command]
1933
           })
           return result
1934
1935
       def natural_sort_key(s):
1936
1937
           Generate a key for natural sorting.
1938
           This function splits the string into numeric and non-numeric parts so
1939
               that,
1940
           for example, "file2" is sorted before "file10".
1941
1942
           return [int(text) if text.isdigit() else text.lower() for text in re.
              split(r'(\d+)', s)]
1943
```

```
1944
      def get_second_level_entries(root_dir):
1945
1946
           Retrieve all second-level entries (files and directories) under the
1947
              specified root directory,
           and sort them so that directories come first, then files. Both are
1948
              sorted naturally.
1949
1950
           trv:
1951
               entries = list(os.scandir(root_dir))
1952
           except Exception as e:
               print(f"Error scanning {root_dir}: {e}")
1953
               return []
1954
1955
           entries.sort(key=lambda e: (not e.is_dir(), natural_sort_key(e.name))
1956
           return entries
1957
1958
       def collect_all_files_from_directory(directory):
1959
1960
           Recursively collect all file paths from the given directory,
1961
           sorted naturally by their relative paths.
1962
           collected = []
1963
           for root, dirs, files in os.walk(directory):
1964
               dirs.sort(key=natural_sort_key)
1965
               files.sort(key=natural_sort_key)
1966
               for file in files:
1967
                   full_file_path = os.path.join(root, file)
                   relative_path = os.path.relpath(full_file_path, start=
1968
                       directory)
1969
                   collected.append((relative_path, full_file_path))
1970
           collected.sort(key=lambda tup: natural_sort_key(tup[0]))
1971
           return collected
1972
       7. def list_files_in_second_level(root_directory: str) -> dict:
1973
1974
           Traverse all second-level entries under the root directory and return
1975
                a summary dictionary.
1976
           print(f"Running list_files_in_second_level tool under {root_directory
1977
              } . . . ")
1978
           max_files = 10
1979
           results = []
1980
           second_level_entries = get_second_level_entries(root_directory)
1981
           for entry in second_level_entries:
1982
               if entry.is_file():
1983
                   result_dict = {
1984
                        "entry_name": entry.name,
1985
                        "entry_path": entry.path,
1986
                        "total_files": 1,
                        "files": [entry.path]
1987
1988
                   results.append(result_dict)
1989
               elif entry.is_dir():
1990
                   collected_files = collect_all_files_from_directory(entry.path
1991
                   total_file_count = len(collected_files)
1992
                   top_files = [full_path for _, full_path in collected_files[:
1993
                       max_files]]
1994
                   result_dict = {
1995
                        "entry_name": entry.name,
1996
                        "entry_path": entry.path,
                        "total_files": total_file_count,
1997
                        "files": top_files
```

```
1998
1999
                   results.append(result_dict)
2000
2001
           final_result = {"entries": results}
           print(final_result)
2002
           return final_result
2003
2004
       8. def preview_file_content(file_path: str) -> str:
2005
2006
           Preview the contents of CSV, JSON, or TXT files.
2007
           print(f"Running preview_file_content tool for {file_path}...")
2008
           if file_path.lower().endswith('.csv'):
2009
               rows = []
2010
               total\_rows = 0
2011
               try:
                   with open(file_path, 'r', encoding='utf-8') as f:
2012
                       reader = csv.reader(f)
2013
                        for row in reader:
2014
                            total_rows += 1
2015
                           if total_rows <= 5:</pre>
2016
                                rows.append(row)
               except Exception as e:
2017
                   return f"Error reading CSV file: {e}"
2018
2019
               preview_str = "CSV File Preview:\n"
2020
               for row in rows:
                   preview_str += ", ".join(row) + "\n"
2021
               preview_str += f"Total rows: {total_rows}"
2022
               return preview_str
2023
2024
           elif file_path.lower().endswith('.json'):
2025
               try:
                   with open(file_path, 'r', encoding='utf-8') as f:
2026
                       data = json.load(f)
2027
               except Exception as e:
2028
                   return f"Error reading JSON file: {e}"
2029
2030
               if isinstance(data, dict):
                   items = list(data.items())
2031
                   preview_items = items[:5]
2032
                   preview_str = "JSON File Preview (first 5 key-value pairs):\n
2033
2034
                   for key, value in preview_items:
2035
                       preview_str += f"{key}: {value}\n"
                   preview_str += f"Total key-value pairs: {len(items)}"
2036
               elif isinstance(data, list):
2037
                   preview_items = data[:5]
2038
                   preview_str = "JSON File Preview (first 5 elements):\n"
2039
                   for item in preview_items:
2040
                       preview_str += f"{item}\n"
                   preview_str += f"Total elements: {len(data)}"
2041
               else:
2042
                   preview_str = f"Unsupported JSON type: {type(data)}"
2043
               return preview_str
2044
2045
           elif file_path.lower().endswith('.txt'):
2046
               try:
                   with open(file_path, 'r', encoding='utf-8') as f:
2047
                       content = f.read()
2048
               except Exception as e:
2049
                   return f"Error reading TXT file: {e}"
2050
2051
               words = content.split()
               total_words = len(words)
```

```
2052
               preview_words = words[:10000]
2053
               preview_str = "TXT File Preview (first 10000 words):\n"
               preview_str += " ".join(preview_words)
2054
               preview_str += f"\nTotal words: {total_words}"
2055
               return "=== CSV Preview === \n" + preview_str
2056
2057
          else:
2058
               return "Unsupported file type. Only CSV, JSON, and TXT files are
2059
                   supported.'
2060
2061
       9. def run_selfclean_on_dataset(image_folder_path: str) -> None:
2062
2063
          Run SelfClean on an image folder and generate CSVs for near
2064
              duplicates, off-topic samples, and label errors.
2065
          Aras:
2066
               image_folder_path (str): Path to the root folder containing the
2067
                   images organized by class folders.
2068
2069
           sc_utils.init_distributed_mode = dummy_init_distributed_mode
2070
           # Patch torch.load for compatibility
2071
          original_torch_load = torch.load
2072
           def patched_torch_load(*args, **kwargs):
2073
               kwargs["weights_only"] = False
2074
               return original_torch_load(*args, **kwargs)
2075
           torch.load = patched_torch_load
2076
           resize_images_in_folder(image_folder_path)
2077
2078
          print("Loading dataset with ImageFolder...")
2079
          dataset = ImageFolder(root=image_folder_path)
2080
           parameters = copy.deepcopy(DINO_STANDARD_HYPERPARAMETERS)
2081
          parameters['model']['base_model'] = 'pretrained_imagenet_vit_tiny'
2082
2083
           print("Running SelfClean...")
2084
          selfclean = SelfClean(auto_cleaning=True)
          print("Selfclean loaded")
2085
2086
       def patched_load_pretrained(model_name=None, work_dir=None, **kwargs):
2087
               print("Using locally downloaded DINO checkpoint")
               local_model_path = "path/to/model"
2089
               model = sc_utils.Embedder.load_dino(ckp_path=local_model_path)
2090
               dummy_config = SimpleNamespace(model_type="ViT")
               dummy_augment_fn = lambda x: x
2091
               return model, dummy_config, dummy_augment_fn
2092
           sc_utils.Embedder.load_pretrained = patched_load_pretrained
2093
2094
           work_folder_path = {"..."}.get(image_folder_path, None)
2095
           issues = selfclean.run_on_dataset(
2096
               dataset=copy.copy(dataset),
2097
               pretraining_type=PretrainingType.DINO,
2098
               epochs=10,
2099
               batch_size=16,
2100
               save_every_n_epochs=1,
               dataset_name="...",
2101
               work_dir=work_folder_path,
2102
2103
2104
           df_near_duplicates = issues.get_issues("near_duplicates",
2105
              return_as_df=True)
```

```
2106
           df_off_topic_samples = issues.get_issues("off_topic_samples",
2107
              return_as_df=True)
2108
           df_label_errors = issues.get_issues("label_errors", return_as_df=True
2109
              )
2110
2111
       10. def organize_into_subfolder(root_directory: str,
2112
          destination_directory: str) -> dict:
2113
2114
          Organize images into class-wise subfolders using labels from a CSV
              file.
2115
2116
          trv:
2117
               csv_files = [f for f in os.listdir(root_directory) if f.endswith(
2118
                   ".csv")]
               if len(csv_files) != 1:
2119
                   return {"status": "error", "message": "Expected exactly one
2120
                       CSV file."}
2121
2122
               csv_path = os.path.join(root_directory, csv_files[0])
2123
               df = pd.read_csv(csv_path)
2124
               label_col = [col for col in df.columns if "label" in col.lower()
2125
                  [0]
2126
               file_col = [col for col in df.columns if "file" in col.lower() or
2127
                    "image" in col.lower() or "path" in col.lower()][0]
2128
2129
               moved_count = {}
               for _, row in df.iterrows():
2130
                   label = str(row[label_col]).strip()
2131
                   filename = str(row[file_col]).strip()
2132
                   src_path = filename
2133
                   if not os.path.exists(src_path):
                       continue
2134
2135
                   label_folder = os.path.join(destination_directory, label)
2136
                   os.makedirs(label_folder, exist_ok=True)
2137
                   dst_path = os.path.join(label_folder, os.path.basename(
2138
                       filename))
2139
                   shutil.copy2(src_path, dst_path)
                   moved_count[label] = moved_count.get(label, 0) + 1
2140
2141
               return {"status": "success", "moved": moved_count}
2142
          except Exception as e:
2143
               return {"status": "error", "message": str(e)}
2144
2145
       11. def copy folder(source directory: str, destination directory: str) ->
2146
           dict:
2147
2148
          Copies all files and subdirectories from source to destination.
2149
           try:
2150
               if not os.path.exists(source_directory):
2151
                   return {"status": "error", "message": f"Source folder does
2152
                       not exist: {source_directory}"}
2153
               os.makedirs(destination_directory, exist_ok=True)
2154
               for item in os.listdir(source_directory):
2155
                   src = os.path.join(source_directory, item)
2156
                   dst = os.path.join(destination_directory, item)
2157
                   if os.path.isdir(src):
2158
                       shutil.copytree(src, dst, dirs_exist_ok=True)
2159
                   else:
                       shutil.copy2(src, dst)
```

```
2160
2161
               return {"status": "success", "message": f"Copied from {
2162
                   source_directory} to {destination_directory}"}
2163
           except Exception as e:
               return {"status": "error", "message": str(e)}
2164
2165
2166
       12. def remove_other_files(root_directory: str) -> dict:
2167
2168
           Remove all non-image files from a directory and its subdirectories.
2169
           allowed_extensions = {'.jpg', '.jpeg', '.png', '.bmp', '.tiff', '.tif
2170
               ', '.gif', '.dcm', '.nii', '.nii.gz', '.mha', '.mhd', '.hdr', '.
2171
               img', '.nrrd'}
2172
           removed_files = []
2173
           for dirpath, _, filenames in os.walk(root_directory):
2174
               for filename in filenames:
2175
                   ext = os.path.splitext(filename)[1].lower()
2176
                   if ext not in allowed_extensions:
2177
                       file_path = os.path.join(dirpath, filename)
2178
                       try:
                           os.remove(file_path)
2179
                           removed_files.append(file_path)
2180
                       except Exception as e:
2181
                           print(f"Error removing {file_path}: {e}")
2182
           return {"status": "success", "removed_file_count": len(removed_files)
2183
               , "removed_files": removed_files}
2184
2185
2186
       13. def list_folders(root_directory: str) -> dict:
2187
           List subfolders in the given directory.
2188
2189
           folders = [f for f in os.listdir(root_directory) if os.path.isdir(os.
2190
               path.join(root_directory, f))]
2191
           return {"folders": folders}
2192
2193
       14. def make_folder(root_directory: str) -> dict:
2194
2195
           Create a new folder at the given path.
2196
2197
           try:
               os.makedirs(root_directory, exist_ok=True)
2198
               return {"status": "success", "message": f"Created folder: {
2199
                   root_directory}"}
2200
           except Exception as e:
2201
               return {"status": "error", "message": str(e)}
2202
2203
       15. def copy_images(src_folder: str, dst_folder: str) -> dict:
2204
2205
           Copies all image files from the source folder (including subfolders)
2206
               to the destination folder.
2207
2208
           Args:
               src_folder (str): Path to the source folder containing image
2209
                   files.
2210
               dst_folder (str): Path to the destination folder where images
2211
                   will be copied.
2212
2213
           Returns:
```

```
2214
               dict: Summary of copied images including total copied count and
2215
                   failed files.
2216
           allowed_extensions = {'.jpg', '.jpeg', '.png', '.bmp', '.tiff', '.tif
2217
               ', '.gif', '.dcm'}
2218
           copied_files = []
2219
           failed_files = []
2220
2221
           os.makedirs(dst_folder, exist_ok=True)
2222
           for root,
                      _, files in os.walk(src_folder):
2223
               for file in files:
2224
                   ext = os.path.splitext(file)[1].lower()
2225
                    if ext in allowed_extensions:
2226
                        src_path = os.path.join(root, file)
                        dst_path = os.path.join(dst_folder, file)
2227
2228
                        trv:
2229
                            shutil.copy2(src_path, dst_path)
2230
                            copied_files.append(file)
2231
                        except Exception as e:
2232
                            failed_files.append((file, str(e)))
2233
           return {
2234
               "status": "success",
2235
               "copied_count": len(copied_files),
2236
               "failed_count": len(failed_files),
               "failed_files": failed_files
2237
           }
2238
2239
2240
       16. def run_federated_method(project_directory: str, method_name: str) ->
2241
           Dict:
2242
           Run federated training using a specified method inside a given
2243
               project directory.
2244
2245
           try:
2246
               result = subprocess.run(
                    ["python", "/path/to/FL-bench/main.py", f"method={method_name
2247
                        }"],
2248
                   cwd=project_directory,
2249
                   stdout=subprocess.PIPE,
2250
                   stderr=subprocess.PIPE,
2251
                    text=True
               )
2252
2253
               return {
2254
                    "status": "success" if result.returncode == 0 else "failed",
2255
                   "stdout": result.stdout,
2256
                    "stderr": result.stderr,
                    "exit_code": result.returncode
2257
2258
           except Exception as e:
2259
               return {
2260
                    "status": "error",
2261
                    "message": str(e)
2262
               }
```

## B.2 ROLE-SPECIALIZED AGENTS

226322642265

2266

2267

To enable automated, modular, and scalable orchestration of federated learning workflows, we introduce a suite of seven specialized LLM agents within the FedAgentBench framework. Each

agent is assigned a distinct responsibility aligned with a specific stage of the FL pipeline, spanning from task interpretation and dataset selection to data preparation, label harmonization, algorithm selection, and training. These agents collectively simulate the collaborative behavior typically required from domain experts, data engineers, and FL researchers, while interacting through well-defined prompts and toolchains. Code snippets of all 7 role-specialized agents can be found in Listings 2-5 with each discussing agents of individual phases.

## RESPONSIBILITIES OF FEDAGENTBENCH AGENTS:

As a part of FedAgentBench, we design a modular and collaborative framework composed of seven specialized LLM agents, each responsible for a distinct role in the federated learning pipeline and operating via specific toolsets (if necessary) that allow them to automate key stages of client-server coordination, data preparation, and model training. Table 3 summarizes the roles of the seven specialized agents. Below, we describe the function of each agent in the context of the four major phases of the workflow.

- 1. Server Agent for Task Interpretation  $(S_1)$ : This agent parses the user-defined instruction to identify the intended task and required data modality. It then broadcasts this extracted requirement to all client agents to begin the dataset discovery process.
- 2. Client Selector Agent  $(C_1)$ : After receiving the task description from the server, this agent inspects the metadata of available datasets and determines which of them are relevant to the given task. The selection is based on textual descriptions stored in a structured JSON file. This task is facilitated using the read\_files function to analyze the dataset content. The agent responds with matching dataset names or returns "no dataset" if none are suitable.
- 3. Server Agent for Client Approval  $(S_2)$ : This agent is responsible for validating the responses returned by the client agents. If a client proposes one or more datasets, the server responds with "Approved. Prepare for training". If the client has no relevant data, the server sends "Client not needed for the task" to exclude them from training.
- 4. Data Pre-processor Agent  $(C_2)$ : This agent ensures the dataset is well-organized and free from noisy or irrelevant samples. It first checks whether the dataset is structured in class-specific subfolders. If not, it reorganizes the data accordingly. It then eliminates all non-image files and performs content-based cleaning to flag duplicates, off-topic, or mislabeled samples. These operations can be carried out using tools such as organize\_into\_subfolder, remove\_other\_files, and run\_selfclean\_on\_dataset discussed earlier. The agent concludes by signaling completion with "Data Cleaning Complete <end>".
- 5. Task conditioned Label Harmonizer Agent  $(C_3)$ : This agent unifies the class label space across multiple clients by remapping existing class folders into a shared label schema (e.g., from fine-grained categories to binary classes like malignant or benign). It first lists the current folder names, defines a harmonization mapping, and creates new folders to reflect the harmonized schema. This can be accomplished using list\_folders, make\_folder, and copy\_images functions mentioned earlier.
- 6. **FL Algorithm Selector Agent**  $(S_3)$ : This agent chooses the most appropriate federated learning algorithm for training based on the user's task requirement. It examines a JSON file describing available algorithms and selects one based on the alignment of its key idea and name with the user's intent. This process can be supported by the read\_files tool and results in a response such as "Algorithm Name: ... <end>".
- 7. Trainer Agent  $(S_4)$ : Once the data and algorithm are finalized, this agent launches federated training using the selected method. It delegates execution to the appropriate script that implements the algorithm. This can be done by calling the run\_federated\_method tool.

**Justification of Agent Design.** The decomposition into seven specialized agents is grounded in the need to modularize a complex and multi-phase federated learning pipeline that must accommodate the broad diversity of FL algorithms (as evidenced in FL-Bench, spanning aggregation-based,

2323

2324

2325

2326

2327

2328

2329

2330

233223332334

personalization-based, and representation-based strategies) and ensure automation across heterogeneous datasets and institutional constraints. The separation of concerns allows each agent to handle a distinct phase of the workflow: high-level task parsing  $(S_1)$ , distributed dataset discovery  $(C_1)$ , client validation  $(S_2)$ , data reorganization and quality control  $(C_2)$ , cross-client label harmonization  $(C_3)$ , FL algorithm selection conditioned on user intent  $(S_3)$ , and training orchestration  $(S_4)$ . This division aligns with the key bottlenecks in real-world FL deployment. The agent specialization ensures scalability, adaptability, and plug-and-play extensibility of the framework, enabling future integration of additional FL capabilities (e.g., fairness, security, cross-silo adaptation) without architectural redesign. The code snippets of the individual specialized agents are provided below:

CODE SNIPPETS OF SPECIALIZED AGENTS:

### Listing 2: Prompt definition for Client Orchestrator Agents

```
2335
      def create_server_to_client_communication_prompt_round_1():
2336
           system_prompt =
           You are a server agent in a Federated Learning setup, responsible for
2337
               communicating with the client agents.
2338
          From the user requirement, only extract the task and modality
2339
              information.
2340
           State this information and instruct the clients to respond with:
2341
           - The name of the selected dataset (that matches the user requirement
2342
           . . .
2343
          return system_prompt
2344
2345
      # Goal-oriented guidance
2346
      def create_selector_prompt(description_path, server_instruction):
          system_prompt = f"""
2347
           You are acting as a client agent in Federated Learning responsible
2348
              for selecting the datasets in your client based on the server
2349
              instructions: {server_instruction}.
2350
           I provide you with a list of dataset descriptions: {description_path
2351
               }, which is a json file that contains a list of dictionaries.
              Plan your workflow and solve the task:
2352
2353
           You have access to the tool:
2354
           read_files: This function reads a script file (such as a Python file)
2355
               so you can understand its content.
2356
           Return the chosen dataset names following {server_instruction}, so a
2357
              downstream peer agent can know the information accurately.
2358
           IMPORTANT: Give it only in this template for each dataset: **Dataset
2359
              Name **: .... If no suitable dataset for the given task exists,
2360
              the client should return "no dataset" and clearly explain why
2361
              before ending the conversation.
           Include <end> to end the conversation.
2362
2363
           return system_prompt
2364
       # Fine-grained guidance
2366
      def create_selector_prompt(description_path, server_instruction):
           system_prompt = f"""
2367
           You are acting as a client agent in Federated Learning responsible
2368
              for selecting the datasets in your client based on the server
2369
              instructions: {server_instruction}.
2370
           I provide you with a list of dataset descriptions: {description_path
2371
               }, which is a json file that contains a list of dictionaries.
           Every dictionary contains following entries: ["Dataset
2372
              Name", "Dataset Description", "dataset_path"].
2373
2374
           You have access to the tools:
2375
           read_files: This function reads a script file (such as a Python file)
               so you can understand its content.
```

```
2376
2377
          Here is the typical workflow you should follow:
2378
          1. Use read_files to read {description_path}, understand its content.
2379
           2. Choose all the datasets that match the server instructions.
              Remember, your choice should be mainly based on "dataset
2380
              descriptions" entry.
2381
           3. Return the chosen dataset names following {server_instruction}, so
2382
               a downstream peer agent can know the information accurately.
2383
           IMPORTANT: Give it only in this template for each dataset: **Dataset
2384
              Name **: .... If no suitable dataset for the given task exists,
              the client should return "no dataset" and clearly explain why
2385
              before ending the conversation.
2386
           4. Include <end> to end the conversation.
2387
2388
           return system_prompt
2389
      def create_server_to_client_communication_prompt_round_2(client_response)
2390
2391
           system_prompt = f"""
2392
           You are acting as a server agent for communicating with the client
2393
              agents in Federated Learning. Read the client response: {
2394
              client_response}
           If the client has returned one or more datasets, return the message:
2395
              "Approved. Prepare for training".
2396
           If the client has returned no dataset, return the message: "Client
2397
              not needed for the task".
2398
2399
           return system_prompt
2400
```

Listing 3: Prompt definition for Data Pre-processor Agent

```
2402
      # Goal-oriented guidance
2403
      def create_datacleaner_prompt(input_data_path, output_data_path,
2404
          server_response_round_2, description_path):
          system_prompt = f"""
2405
           You are a highly skilled data preparation and data cleaning agent
2406
              specializing in the medical domain. You MUST do your tasks ONLY
2407
              using the tools provided to you.
2408
           You MUST plan the workflow based on the instruction given below
2409
              sincerely and not bypass it.
           I provide you with server instruction {server_response_round_2}.
2410
           If the server mentions that the client is not needed, end the
2411
              conversation and do NOT do anything else. Instead, if it
2412
              instructs to prepare for training, you have three tasks:
2413
           1. Check if the dataset in {input_data_path} is already organized in
2414
              sub-folder format from dataset descriptions: {description_path}.
               If not, organize the data by grouping images of each class into
2415
                  their respective subfolders in your destination path: {
2416
                  output_data_path } .
2417
          2. Remove all non-image files from each sub-folder.
2418
           3. Clean client data by removing (a) near duplicate samples, (b) off
2419
              topic samples, (c) noisy label samples
2420
          You have access to the following tools. Plan and reason how to use
2421
              the following tools properly:
2422
           read_files: This function reads a script file (such as a Python file)
2423
               so you can understand its content.
          organize into subfolder: This function reads csv file, goes through
2424
              the labels column, creates subfolders and groups images inside
2425
              them based on labels column.
2426
          copy_folder: This function copies folder from source location to
2427
              destination location.
2428
           remove_other_files: This function checks the file extension of all
2429
              files in a given folder and deletes the files with non-image
              extensions.
```

```
2430
           run_selfclean_on_dataset: This function flags (a) near duplicate
2431
              samples, (b) off topic samples, (c) noisy label samples. Use this
2432
               to clean the dataset
2433
           Important rules you must follow:
2434

    You MUST use the run_selfclean_on_dataset tool to clean data!

2435
           - You MUST NOT modify the raw images manually.
2436
           - You MUST conclude your work by writing: "Data Cleaning Complete" <
2437
              end>.
2438
          return system_prompt
2439
2440
       # Fine-grained guidance
2441
      def create_datacleaner_prompt(input_data_path, output_data_path,
2442
          server_response_round_2, description_path):
          system_prompt = f"""
2443
           You are a highly skilled data preparation and data cleaning agent
2444
              specializing in the medical domain. I provide you with server
2445
              instruction {server_response_round_2}.
2446
           If the server mentions that the client is not needed, end the
2447
              conversation. If it instructs to prepare for training, you have
2448
              three tasks:
           1. Check if the dataset in {input_data_path} is already organized in
2449
              sub-folder format from dataset descriptions: {description_path}.
2450
               If not, organize the data by grouping images of each class into
2451
                  their respective subfolders in your destination path: {
2452
                  output_data_path } .
           2. Remove all non-image files from each sub-folder.
2453
           3. Clean client data by removing (a) near duplicate samples, (b) off
2454
              topic samples, (c) noisy label samples
2455
2456
          You have access to the tools:
2457
           read_files: This function reads a script file (such as a Python file)
               so you can understand its content.
2458
          organize_into_subfolder: This function reads csv file, goes through
2459
              the labels column, creates subfolders and groups images inside
2460
              them based on labels column.
2461
          copy_folder: This function copies folder from source location to
2462
              destination location.
           remove_other_files: This function checks the file extension of all
2463
              files in a given folder and deletes the files with non-image
2464
              extensions.
2465
           run_selfclean_on_dataset: This function flags (a) near duplicate
2466
              samples, (b) off topic samples, (c) noisy label samples. Use this
2467
               to clean the dataset
          clean_data: This function checks flagged samples from csv file and
2468
              removes them.
2469
2470
          Here is the typical workflow you should follow:
2471
           1. If the server instruction: {server_response_round_2} mentions that
2472
               the client is not needed, print <end> and end the conversation.
              Do NOT do anything further.
2473
           2. Instead, if it instructs you to prepare for training, use "
2474
              read_files" function to read and understand the dataset
2475
              description file in {description_path}. Check from there, if the
2476
              dataset in {input_data_path} is already organized as sub-folders.
2477
              If yes, copy the folder to the destination folder {
                 output_data_path } using the function "copy_folder" and go to
2478
                 step 4 below, skipping step 3.
2479
           3. If dataset is not organized as sub-folders, organize the data by
2480
              grouping images of each class into their respective subfolders in
2481
               the destination data path: {output_data_path} by using the
2482
              organize_into_subfolder function.
```

```
2484
           4. Go to each subfolder in the destination data path: {
2485
              output_data_path} and remove all non-image files by using
2486
              remove_other_files function.
2487
           5. Flag (a) near duplicate samples, (b) off topic samples, (c) noisy
              label samples using run_selfclean_on_dataset function.
2488
           6. Remove the flagged samples using clean_data function.
2489
2490
           Important rules you must follow:
2491
           - You MUST use the run_selfclean_on_dataset tool to clean data!
2492
           - You MUST NOT modify the raw images manually.
           - You MUST clean using the CSV outputs only
2493
           - You MUST conclude your work by writing: "Data Cleaning Complete" <
2494
              end>.
2495
2496
           return system_prompt
2497
```

# Listing 4: Prompt definition for Label Harmonization Agent

```
2500
       # Goal-oriented guidance
2501
      def label_harmonizer_prompt(input_data_path, output_data_path):
          system_prompt = f"""
2502
           You are an intelligent agent tasked with harmonizing medical image
2503
              labels in a Federated Learning environment.
2504
2505
          Your objective is to reorganize the dataset located at {
2506
              input_data_path} by grouping existing class folders into
2507
              standardized, harmonized categories (e.g., 'malignant', 'benign')
               based on the task specification.
2508
          You should inspect the current folder structure, define appropriate
2510
              label mappings to target categories, and reorganize the data into
2511
               the {output_data_path} directory using the available tools.
2512
          You have access to the following tools:
2513
            list_folders(path): Lists existing class folders in a dataset.
2514
            make_folder(path): Creates a new folder for a target label.
2515
           - copy_images(src_folder, dst_folder): Copies all image files from
2516
              the original to the harmonized destination folder.
2517
          Use these tools to achieve the goal of producing a clean, consistent
2518
              label space for downstream federated training.
2519
          When harmonization is complete, end your process with "<end>".
2521
          return system_prompt
2522
       # Fine-grained guidance
2523
      def label_harmonizer_prompt(input_data_path, output_data_path):
2524
          system_prompt = f"""
2525
           You are an intelligent agent for medical image label harmonization in
2526
               a Federated Learning setup.
          Your goal is to group existing class folders into harmonized target
2527
              categories (e.g., 'malignant', 'benign') by reorganizing the
2528
              folder structure.
2529
          This involves identifying the current class folders, mapping them to
2530
              new target labels, and copying images accordingly.
2531
          You have access to the tools:
2532
           - list_folders(path): Returns a list of subfolder names in the given
2533
              path.
2534
           - make_folder(path): Creates a new directory at the specified path.
2535
           - copy_images(src_folder, dst_folder): Copies all image files from
2536
              the source to the destination folder.
2537
          Here is the typical workflow you should follow:
```

```
2538
           1. Inspect class structure: Use 'list_folders("{input_data_path}")'
2539
              to get all existing class folder names.
2540
           2. Define label mapping: Based on user requirements (e.g., binary
2541
              classification), decide how existing class names map to target
              classes (coarse labels like 'malignant' and 'benign').
2542
           3. Prepare new folders: For each target class, use 'make_folder("{
2543
              output_data_path}/<class_name>") ` to create destination folders.
2544
           4. Move data: For each source class, use 'copy_images' to move all
2545
              image files to their new harmonized folder.
2546
2547
           return system_prompt
2548
```

## Listing 5: Prompt definition for Federated Trainer Agents

```
2551
       # Goal-oriented guidance
2552
      def FL_algorithm_selector_prompt(algorithm_description_path):
2553
          system_prompt = f"""
2554
           You are a server agent in a Federated Learning setup responsible for
2555
              selecting the most appropriate federated learning algorithm based
2556
               on the human users task requirement.
2557
          You are provided with a list of algorithm descriptions in the file {
2558
              algorithm_description_path}, formatted as a JSON list of
2559
              dictionaries. Each dictionary contains information about an
2560
              algorithm, including its name, full name, and key idea.
2561
          Your objective is to analyze the algorithm descriptions and identify
2562
              the method that best aligns with the users intent. Focus
2563
              primarily on the "Full Name" and "Key idea" fields to determine
2564
              relevance.
2565
2566
          You have access to the following tool:
           - read_files: This function reads a script file (such as a Python
2567
              file) so you can understand its content.
2568
2569
          Once you have selected the most suitable algorithm, return it in the
2570
              format:
2571
          Algorithm Name: <selected_algorithm>
2572
          Conclude your response with "<end>".
2573
2574
          return system_prompt
2575
2576
       # Fine-grained guidance
2577
      def FL_algorithm_selector_prompt(algorithm_description_path):
2578
          system_prompt = f"""
2579
           You are acting as a server agent in Federated Learning responsible
2580
              for selecting the federated learning algorithm in your client
2581
              based on the human user requirement.
           I provide you with a list of algorithm descriptions: {
2582
              algorithm_description_path}, which is a json file that contains a
2583
               list of dictionaries.
2584
          Every dictionary contains following entries: ["algorithm", "Full Name
2585
               ", "Key idea"].
2586
          You have access to the tools:
2587
           read_files: This function reads a script file (such as a Python file)
2588
               so you can understand its content.
2589
2590
          Here is the typical workflow you should follow:
2591
          1. Use read_files to read {algorithm_description_path}, understand
              its content.
```

```
2592
           2. Choose the algorithm that best matches the server instructions.
2593
              Remember, your choice should be mainly based on "Full Name", "Key
2594
               idea" entries.
           3. Return the chosen algorithm as Algorithm Name: ....
2595
           4. Include <end> to end the conversation.
2596
2597
          return system_prompt
2598
2599
2600
       def FL_trainer_prompt(project_directory, selected_algorithm):
           system_prompt = f"""
2601
           You are a trainer agent that performs federated learning with
2602
              selected clients using the chosen algorithm: {selected_algorithm}
2603
           You have access to the tools:
2604
           run_federated_method: Runs the specified federated learning method
2605
          Use run_federated_method to run the specific federated learning
2606
              algorithm: {selected_algorithm} and report results.
2607
2608
           return system_prompt
2609
```

Table 3: Summary of Specialized Agents and Their Responsibilities in Federated Learning Workflow

Agent	Agent Name	Role Description	Phase
$S_1$	Server Agent for	Parses user instructions to extract task and	Phase 1: Client
	Task Interpreta-	modality requirements; broadcasts the re-	Selection
	tion	quirement to all client agents to begin	
		dataset selection.	
$C_1$	Client Selector	Evaluates dataset metadata to identify rele-	Phase 1: Client
	Agent	vant datasets for the task based on textual	Selection
		descriptions in a JSON file; responds with	
		matched datasets or "no dataset".	
$S_2$	Server Agent for	Reviews responses from clients; approves	Phase 1: Client
	Client Approval	those with relevant datasets for training or	Selection
		excludes irrelevant ones.	
$C_2$	Data Pre-	Organizes dataset into class-wise subfold-	Phase 2: Data
	processor Agent	ers, removes non-image files, and performs	Preparation
		data cleaning (e.g., de-duplication, noise fil-	
		tering, off-topic detection).	
$C_3$	Task-conditioned	Reorganizes client label spaces into har-	Phase 3: Label
	Label Harmonizer	monized schema by mapping fine-grained	Harmonization
	Agent	classes to broader target labels (e.g.,	
		malignant, benign).	
$S_3$	FL Algorithm Se-	Selects the most appropriate federated learn-	Phase 4: FL Algo-
	lector Agent	ing algorithm based on the user's task by	rithm Selection
		analyzing algorithm metadata.	
$S_4$	Trainer Agent	Executes the federated learning training us-	Phase 4: Feder-
		ing the chosen algorithm and the approved	ated Training
		client datasets.	

## C TASKS AND ALGORITHMS IN FEDAGENTBENCH FRAMEWORK

#### C.1 Dataset Details

To enable systematic benchmarking across a broad range of real-world clinical scenarios, FedA-gentBench includes 201 publicly available datasets spanning six major medical imaging modalities: Dermatology (25 datasets), Ultrasound (33), Fundus (63), X-Ray (32), MRI (28), and Histopathology (20). These datasets comprise both 2D and 3D imaging formats and cover a wide array of task types, including classification (e.g., tumor detection, cancer subtype identification), grading/staging (e.g.,

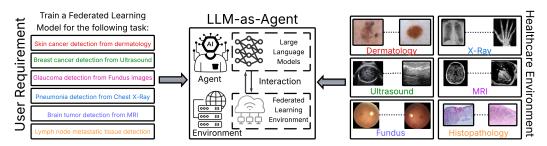


Figure 8: Sample tasks and datasets in FedAgentBench

diabetic retinopathy, cancer severity), segmentation (e.g., lesion, tumor, or stroke localization), object detection, regression, image reconstruction, and registration.

Each client in FedAgentBench is simulated by grouping one or more of these datasets, thereby reflecting the diversity and data heterogeneity found in real-world healthcare settings. For each client, a datacard is constructed, compiling metadata sourced from the original dataset publication, repository, or project website. This metadata includes information on imaging modality, data dimensionality, task type, class schema, and clinically relevant attributes, ensuring traceability and reproducibility.

In the following subsections, we provide a detailed breakdown of the dataset description for each imaging modality.

#### **DERMATOLOGY:**

The dermatology dataset collection curated for this benchmark represents one of the most comprehensive and heterogeneous sets assembled for machine learning research in skin disease analysis. Spanning over 25 datasets, the collection includes both photographic and dermoscopic images, structured tabular data, and multi-modal formats. The classification tasks range from binary cancer detection (e.g., benign vs. malignant in *ISIC2020*, *Mednode*) to fine-grained multi-class diagnosis involving over twenty conditions (e.g., *Dermnet*, *Derma7PT*, *skinL2\_dataset*). Several datasets such as *DDI\_skin\_dataset* and *fitzpatrick17k* are designed to ensure skin tone diversity, while others like *Monkeypox\_Skin\_Image\_Dataset* and *skin-infection-disease-dataset* address emerging and infectious conditions. Additionally, datasets like *PH2Dataset*, *ISIC2016–2024*, and *Dermis* support segmentation and localization, enabling both classification and pixel-wise lesion analysis. This diversity reflects a realistic, clinically relevant spectrum of dermatological challenges, and is particularly well-suited for benchmarking federated learning agents under varying input types, diagnostic complexity, and data distributions. The code snippets for dermatology dataset description file can be found in Listing 6. The description of each dataset is summed up below:

- **1.** augmented\_skin\_condition\_dataset\_kaggle. The augmented\_skin\_condition\_dataset\_kaggle dataset aug is designed for multi-class skin disease classification. It contains photographic images of six dermatological conditions: Acne, Carcinoma, Eczema, Keratosis, Milia, and Rosacea, supporting automated detection and differentiation of common skin ailments.
- **2. DDI\_skin\_dataset.** The *DDI\_skin\_dataset* Daneshjou et al. (2022) is a skin cancer classification resource with strong representation of diverse skin tones. Each image is annotated as benign or malignant, enabling the development of robust melanoma and non-melanoma skin cancer detection algorithms for varied populations.
- **3. Derma7PT.** Derma7PT Kawahara et al. (2018) is a multi-class skin disease classification dataset, annotated with ten distinct diagnostic categories: basal cell carcinoma, nevus, dermatofibroma, lentigo, melanoma, melanoma metastasis, melanosis, miscellaneous, seborrheic keratosis, and vascular lesion. It is suitable for fine-grained disease discrimination in clinical dermatology.
- **4. Dermatology\_tabular dataset.** The *Dermatology\_tabular* Der (a) dataset provides structured clinical features for diagnosing various skin diseases. It is intended for the development and bench-

marking of machine learning models using tabular (non-image) data for dermatological decision support.

- **5. Dermis.** *Dermis* Der (b) is a dual-purpose dataset supporting both skin lesion classification (benign vs malignant) and lesion segmentation. It is suitable for the development of algorithms targeting melanoma recognition and precise lesion boundary detection.
- **6. Dermnet.** Der (c) is a broad dermatology image dataset encompassing 23 disease categories, ranging from inflammatory conditions (e.g., eczema, psoriasis) to infectious (bacterial, viral, fungal), neoplastic (melanoma, carcinoma), and other rare skin diseases. It is valuable for comprehensive multi-class skin disease classification.
- **7. Dermquest.** Der (d) offers images for both classification (benign vs malignant) and segmentation of skin lesions, supporting research in melanoma detection as well as pixel-wise lesion analysis.
- **8. fitzpatrick17k.** The *fitzpatrick17k* Groh et al. (2021) dataset features a wide range of dermatological disease images, annotated with three high-level categories: non-neoplastic, benign, and malignant. Its diverse cases make it well suited for studying skin cancer classification across various skin tones.
- **9. ISIC2018\_HAM10000.** The *ISIC2018\_HAM10000* Codella et al. (2019) dataset is a standard benchmark for skin lesion diagnosis and segmentation, including cases such as melanocytic nevus, benign keratosis, melanoma, basal cell carcinoma, actinic keratosis, vascular lesions, and dermatofibroma. It is used for both classification and lesion segmentation.
- **10. ISIC\_2016.** *ISIC\_2016* Gutman et al. (2016) supports binary classification (benign vs malignant) and lesion segmentation for skin cancer detection, with a focus on melanoma diagnosis in clinical dermoscopic images.
- **11. ISIC\_2017.** *ISIC\_2017* Berseth (2017) targets the detection and segmentation of melanoma and seborrheic keratosis in dermoscopic images, supporting both binary and multi-label skin cancer classification tasks.
- **12. ISIC\_2019.** The *ISIC\_2019* Combalia et al. (2019) dataset offers an expanded benchmark for skin disease classification, with images labeled for nine conditions including melanoma, nevus, basal cell carcinoma, actinic keratosis, and others, facilitating studies in multi-class lesion recognition.
- **13. ISIC\_2020.** *ISIC\_2020* **ISI** (a) is a binary skin lesion classification dataset, primarily focused on discriminating benign from malignant lesions in dermoscopic images for melanoma screening.
- **14. ISIC\_2024.** The *ISIC\_2024* ISI (b) dataset continues the ISIC challenge series with an updated collection focused on binary melanoma (benign vs malignant) classification for automated skin cancer diagnosis.
- **15. Mednode.** *Mednode* MED is a binary classification dataset distinguishing between melanoma and nevus, intended for the development and validation of melanoma detection models.
- **16.** Monkeypox\_Skin\_Image\_Dataset. The *Monkeypox\_Skin\_Image\_Dataset* Mon supports image-based classification of viral skin diseases, including Monkeypox, Chickenpox, Measles, and Normal skin, for research on differential diagnosis of infectious exanthems.
- **17. PAD\_UFES\_20.** *PAD\_UFES\_20* Pacheco et al. (2020) provides images and diagnostic labels for six skin tumor types: melanoma, melanocytic nevus, basal cell carcinoma, actinic keratosis/Bowen's disease, seborrheic keratosis, and squamous cell carcinoma, supporting both single- and multi-class lesion classification.

- **18. PH2Dataset.** The *PH2Dataset* PH2 contains dermoscopic images and expert-annotated segmentation masks for three classes: common nevus, atypical nevus, and melanoma, making it suitable for both lesion segmentation and classification.
- **19. scin\_dataset.** *scin\_dataset* Ward et al. (2024) is a multi-class classification dataset including a range of common skin diseases, such as acne, pigmentary problems, nail disorders, hair loss, and others, for developing comprehensive skin disease classifiers.
- **20. skin\_disease\_3\_class.** The *skin\_disease\_3\_class* dataset comprises images for classifying three skin diseases: acne, atopic dermatitis, and basal cell carcinoma.
- **21. skin\_disease\_classification\_kaggle.** *skin\_disease\_classification\_kaggle* ski (a) is a small dataset for multi-class classification of acne, eye bags, and redness, designed for image-based diagnosis of common cosmetic and inflammatory skin conditions.
- **22. skin\_disease\_kaggle\_dataset.** The *skin\_disease\_kaggle\_dataset* supports multi-class skin disease classification for ten clinically relevant categories, including atopic dermatitis, basal cell carcinoma, eczema, melanoma, nevi, psoriasis, seborrheic keratosis, and infectious diseases.
- **23. Skin Disease\_Robo.** *Skin Disease\_Robo* is a skin disease dataset for both image classification and object detection. It provides bounding box annotations for ten skin disease classes, including acne, atopic dermatitis, eczema, leprosy, psoriasis, ringworm, and warts.
- **24. skin-infection-disease-dataset.** The *skin-infection-disease-dataset* ski (b) focuses on the classification of eight infectious skin diseases, covering bacterial, fungal, parasitic, and viral infections such as cellulitis, impetigo, athlete's foot, ringworm, cutaneous larva migrans, chickenpox, and shingles.
- **25. skinL2\_dataset.** The *skinL2\_dataset* de Faria et al. (2019) is a skin cancer classification resource annotated for eight disease classes, including basal cell carcinoma, dermatofibroma, hemangioma, melanoma, nevus, psoriasis, seborrheic keratosis, and others, facilitating both melanoma and non-melanoma skin lesion research.

#### **ULTRASOUND:**

The ultrasound dataset collection constitutes a diverse and representative corpus of ultrasound images. Spanning over 33 datasets, this collection captures the breadth of clinical applications across multiple anatomical regions (e.g., breast, fetal brain, liver, thyroid, heart, vascular system, musculoskeletal structures), imaging modalities (e.g., B-mode, Doppler, color flow), and task types (e.g., classification, segmentation, super-resolution, registration). Classification challenges range from binary diagnostic tasks such as benign vs. malignant lesion detection (e.g., BUSI, Mendeley, BUET BUSD) to multi-class pathological condition analysis (e.g., PCOS detection, fetal health classification). Several datasets, such as FALLMUD and fetal head US, are curated to support precise biometric measurements and fetal growth monitoring, while others such as CAMUS and leg segmentation datasets are tailored for structure delineation critical in cardiology and musculoskeletal rehabilitation, respectively. The inclusion of multimodal and cross-domain datasets—such as MUS-V (vascular segmentation from Doppler and B-mode), CT2US (CT-to-ultrasound adaptation), and Ultra LR-HR (super-resolution) further enhances the heterogeneity of input formats and computational tasks. In addition, the dataset collection includes rare or emerging clinical tasks such as dermatologic ultrasound, liver fibrosis staging, and hemangioma classification, reflecting real-world diagnostic diversity. This rich variation of organs, pathologies, modalities, and task complexities makes the benchmark exceptionally well-suited for evaluating federated learning agents under diverse diagnostic conditions, cross-institutional generalization scenarios, and clinically realistic constraints.

**1. Breast Ultrasound Images (BUSI):** This dataset BUS (b) is used for images of breast tumors annotated as benign, malignant, or normal. Specifically, it aims to detect and classify breast tumors into benign, malignant, or normal categories, and delineate the exact tumor boundaries in ultrasound images.

- 2808
  2809
  2809
  2810
  2. B-mode fatty liver US images: This dataset Byra et al. (2018) is used for ultrasound images used to classify liver steatosis severity. Specifically, it aims to assess and classify the degree of fatty liver disease (hepatic steatosis) using grayscale B-mode ultrasound scans.
  - **3. Fetal health classification:** This dataset Fet (b) is used for ultrasound data related to fetal health status. Specifically, it aims to evaluate fetal condition based on cardiotocographic or ultrasound signals to classify into normal, suspected, or pathological health status.
  - **4. Robotic handheld lumbar spine US:** This dataset Rob is used for ultrasound images of lumbar spine captured with robotic devices. Specifically, it aims to identify and segment vertebrae and surrounding spinal anatomy from ultrasound images acquired by a robotic handheld device for navigation.
  - **5. BUS-UCLM:** This dataset BUS (a) is used for breast ultrasound dataset from uclm annotated for tumors. Specifically, it aims to differentiate between benign and malignant breast lesions and segment the tumor region for further morphological analysis.
  - **6. Regensburg pediatric appendicitis:** This dataset Reg is used for ultrasound images of pediatric patients for appendicitis diagnosis. Specifically, it aims to distinguish between pediatric patients with and without appendicitis based on ultrasound scans of the abdomen.
  - **7. Breast Ultrasound Images:** This dataset Bre (b) aims to support breast cancer diagnosis by classifying tumors and extracting the region of interest (ROI) for clinical examination.
  - **8. BUS-UC:** This dataset Al-Dhabyani et al. (2020) is used for breast ultrasound dataset from university of california. Specifically, it aims to classify ultrasound-detected breast abnormalities and perform segmentation to assist in diagnostic workflows.
  - **9. Fetal head US dataset:** This dataset Fet (a) is used for images focused on fetal head for biometry (e.g., hc, bpd). Specifically, it aims to extract biometric measurements such as biparietal diameter (BPD) and head circumference (HC) through segmentation of the fetal head.
  - **10. Carotid Ultrasound Images:** This dataset Car (a) is used for ultrasound images of carotid arteries, with plaque annotations. Specifically, it aims to detect carotid artery plaques and measure intima-media thickness (IMT) to evaluate cardiovascular risk.
  - 11. Ultrasound breast images (for cancer): This dataset is used for breast cancer detection. Specifically, it aims to classify breast lesions as benign or malignant in 2D ultrasound scans for early cancer detection.
  - **12. 3D MRI Ultrasound brain images:** This dataset 3D is used for magnetic resonance elastography and ultrasound for brain imaging. Specifically, it aims to analyze brain stiffness and segment relevant anatomical regions in elastography-enhanced 3D ultrasound volumes.
  - **13. CAMUS Human Heart:** This dataset CAM is used for 2D echocardiographic sequences with lv, myocardium, and la labels. Specifically, it aims to segment key cardiac structures such as the left ventricle (LV), myocardium, and left atrium from 2D echocardiography sequences.
  - **14. CT2US for Kidney Seg:** This dataset CT2 is used for CT-derived kidney masks mapped to US domain. Specifically, it aims to leverage CT-derived kidney masks to train ultrasound-based models for accurate kidney segmentation under domain adaptation.
  - **15. Breast Cancer Image Dataset:** This dataset Bre (a) is used for breast cancer detection. Specifically, it aims to differentiate benign and malignant breast lesions to assist in non-invasive cancer diagnosis.

- **16. DDTI: Thyroid US Images:** This dataset DDT is used for digital database for thyroid imaging with nodule annotations. Specifically, it aims to detect and classify thyroid nodules and delineate their contours to support risk stratification and clinical reporting.
- **17. Thyroid Ultrasound:** This dataset Thy is used for thyroid nodule dataset. Specifically, it aims to perform classification and detailed boundary segmentation of thyroid nodules from grayscale ultrasound scans.
- **18.** Multimodal Breast US Dataset (US3M): This dataset US3 is used for multimodal dataset with us, mri, mammo for breast lesions. Specifically, it aims to fuse features from mammography, MRI, and ultrasound to enhance breast tumor classification using multimodal representations.
- **19. Liver histopathology (Fibrosis):** This dataset Liv is used for ultrasound images labeled with fibrosis grades based on biopsy. Specifically, it aims to grade liver fibrosis severity from ultrasound images based on corresponding histopathological findings from biopsy.
- **20. Prostate MRI and Ultrasound:** This dataset pro (b) is used for prostate cancer detection using mri and us fusion. Specifically, it aims to segment the prostate gland and align ultrasound scans with MRI images for guided prostate biopsy or treatment planning.
- **21.** Carotid artery US & Color Doppler This dataset Car (b) is used for detecting stenosis and plaque buildup in the carotid arteries. It typically includes segmentation of the vessel wall and atherosclerotic plaque, along with classification of stenosis severity using Doppler blood flow analysis.
- **22. PCOS Detection using Ultrasound Images** This dataset PCO involves classifying ovarian ultrasound images to detect Polycystic Ovary Syndrome (PCOS). Features such as ovarian volume, follicle count, and echogenicity are commonly used for diagnosis.
- **23. Ultra LR-HR Ultrasound Dataset** An ultrasound dataset ult (a) used for super-resolution tasks, where low-resolution ultrasound images are enhanced or reconstructed into high-resolution versions.
- **24.** MUS-V (Multimodal Ultrasound Vascular Segmentation) This dataset mul integrates multiple ultrasound modalities such as B-mode and Doppler to improve the accuracy of vascular segmentation tasks.
- **25. BUET BUSD** Developed by the Bangladesh University of Engineering and Technology BUE, this breast ultrasound dataset is used for both classification and segmentation of lesions.
- **26. Dermatologic Ultrasound Images** An emerging application of ultrasound for skin lesions der. This dataset is used for classifying dermatological conditions such as melanomas, cysts, or benign tumors.
- **27. FHMS Ultrasound Dataset** This is a fetal head ultrasound dataset fhm.
- **28. Mendeley Breast Ultrasound Dataset** A publicly available dataset men containing 780 images labeled as benign, malignant, or normal. It is frequently used for breast lesion classification.
- **29. FALLMUD** Fetal Abdomen and Longitudinal Liver Measurement in Ultrasound Dataset fal is used for segmentation of the fetal abdomen and liver, important for fetal growth monitoring.
- **30.** Leg Segmentation Ultrasound This dataset leg focuses on segmenting muscles, tendons, and fasciae in ultrasound images of the lower limbs. It has applications in physical therapy and sports medicine.

- **31. Fetal Ultrasound Brain** A dataset of fetal brain ultrasounds fet, commonly used for segmenting brain structures such as the lateral ventricles and midline. It supports fetal development tracking.
- **32. Ultrasound Image Set of Hemangiomas** This dataset includes ultrasound images of hemangiomas, which are benign vascular tumors. It is used for classifying these from other types of soft tissue lesions.
- **33. Ultrasound Nerve Segmentation** This dataset ult (b) comprises ultrasound images for identifying nerve structures of the neck. This would lead to improvement in catheter placement and contribute to reduction in post-surgical pain.

## X-RAY:

The X-ray dataset collection in FedAgentBench represents a highly diverse benchmark suite, encompassing 32 datasets across multiple diagnostic and anatomical categories. It includes chest, bone, knee, dental, and vascular imaging modalities, with tasks ranging from binary classification (e.g., pneumonia vs. normal in *pneumonia*, COVID-19 vs. normal in *cov\_19* and *cov19\_normal*) to complex multi-class and object detection tasks (e.g., *xray\_17\_diseases*, *8\_object\_detection*, and *RSNA-breast-cancer-detection*). Several datasets offer bounding box or pixel-wise segmentation annotations (*NIH\_bbox*, *lung\_segmentation*, *PAX-Ray++*), while others contain structured metadata (e.g., *spr\_age\_gender*, *knee*, *RANZCR*), enabling multi-modal reasoning and demographic prediction. This collection also includes modality-bridging datasets like *HBFMID* that pair X-ray and MRI scans, and datasets that focus on disease-specific localization such as *humerus\_fractures*, *HeelBone*, and *FracAtlas*. Collectively, the X-ray corpus provides a robust foundation for evaluating LLM agents on a wide range of radiological tasks—spanning classification, segmentation, detection, and clinical interpretation under realistic federated learning constraints. The exact dataset descriptions prepared for the client selection agents are provided in Listing 7 and summarized below:

- **1. cov\_19.** The *cov\_19* dataset Rahman (2020) comprises chest X-ray images collected by an international team of researchers, featuring COVID-19 positive cases alongside normal and viral pneumonia images. Initially released with 219 COVID-19, 1,341 normal, and 1,345 viral pneumonia images, the dataset has since expanded to include 3,616 COVID-19 cases, 10,192 normal cases, 6,012 lung opacity (non-COVID lung infection) cases, and 1,345 viral pneumonia cases. Each update has added more images and corresponding lung masks. Data sourcing and ongoing updates make this dataset a valuable resource for developing robust models for COVID-19 and other lung diseases.
- **2. bone\_frac.** The *bone\_frac* dataset Rodrigo (2022) includes X-ray images of fractured and non-fractured bones across various anatomical regions, such as the lower and upper limbs, lumbar spine, hips, and knees. The images are divided into train, test, and validation sets, each containing both classes, making the dataset suitable for training and evaluating bone fracture detection and classification algorithms.
- **3. chest\_tuberculosis\_segmentation.** The *chest\_tuberculosis\_segmentation* dataset Tapendu (2023a) consists of 704 chest X-ray images sourced from the Montgomery County Chest X-ray Database (USA) and the Shenzhen Chest X-ray Database (China). It includes tuberculosis-positive and normal images, accompanied by lung segmentation masks and clinical metadata (e.g., age, gender, county of origin). The combination of images and annotations makes it suitable for tuberculosis detection, segmentation, and broader deep learning tasks in medical imaging.
- **4. xray\_17\_diseases.** The *xray\_17\_diseases* dataset TrainingDataPro (2023) offers chest X-ray images in both .jpg and .dcm formats, labeled for a diverse set of thoracic diseases, including abscess, ARDS, atelectasis, atherosclerosis, cardiomegaly, emphysema, fractures, pneumonia, tuberculosis, and more. The dataset supports research in neurology, radiology, and oncology, enabling the development and evaluation of models for automated disease detection, diagnosis, and classification.
- **5. spr\_age\_gender.** The *SPR Age and Gender* dataset Kitamura (2022a) contains X-ray images in .png format with accompanying CSV files specifying patient age and gender. It is designed for research on patient demographic prediction from radiographic data.

- **6. unifesp.** The *UNIFESP X-Ray Body Part Classification* dataset Kitamura (2022b) comprises 2,481 DICOM-format X-ray images annotated by radiology residents. The dataset covers 20 anatomical body parts (plus an "other" category), with categorical labels assigned to each image, supporting multi-label classification tasks and body part recognition in medical imaging.
- **7. knee.** This dataset Orvile (2023d) contains 1,650 high-quality digital X-ray images of the knee, manually annotated by medical experts using the Kellgren and Lawrence grading system for osteoarthritis severity. The images are 8-bit grayscale and are accompanied by metadata and cartilage region annotations, facilitating research in automated knee osteoarthritis detection and grading.
- **8.** c19\_radiograph. The *c19\_radiograph* dataset Viradiya (2023) is a comprehensive chest X-ray collection curated by a team from Qatar University and the University of Dhaka, with COVID-19, normal, lung opacity, and viral pneumonia cases. The database is built from multiple public and hospital sources and contains extensive clinical labels and patient metadata, enabling detailed studies of COVID-19 pneumonia and related conditions.
- **9. simple\_vs\_community.** This bone fracture dataset Orvile (2023b) is structured to distinguish between simple and comminuted fractures, comprising over 7,500 images for simple fractures and more than 8,500 for comminuted fractures. It combines hospital records and web-sourced images, and includes extensive data augmentation, providing a challenging dataset for fracture classification and segmentation tasks.
- **10. nih\_bbox.** The *NIH Chest X-ray* dataset Hodeb (2023) consists of 112,120 images from 30,805 patients, each labeled for thoracic diseases using text-mined radiology reports. The dataset features bounding box annotations for localization, supports weakly-supervised learning, and includes metadata on disease classes, patient demographics, and imaging protocols.
- **11. bone\_break.** The *bone\_break* dataset Darabi (2023) focuses on the classification of various bone fracture types using X-ray images. It encompasses multiple fracture classes, such as avulsion, comminuted, fracture-dislocations, greenstick, hairline, impacted, longitudinal, oblique, pathological, and spiral fractures, supporting the development of automated fracture classification systems.
- **12. cov19\_normal.** This balanced dataset Tejas (2022) contains 800 high-quality chest X-ray images, equally divided between COVID-19 positive and normal cases (400 each). The curated and balanced nature makes it ideal for deep learning studies on COVID-19 detection.
- **13. dental.** The *dental* dataset IMT Kaggle Team (2023) consists of dental radiographs, enabling the evaluation of hard and soft tissue changes, jawbone development in children, and the detection of injuries in facial and oral structures. It is suitable for a range of dental diagnostic research tasks.
- **14. bone\_frac\_small.** A focused dataset Orvile (2023a) for bone fracture classification and localization in tibia and fibula bones, *bone\_frac\_small* features X-ray images in PNG format. Some images have been validated by medical experts at the University of Gondar, Ethiopia. The dataset includes enhanced and augmented images for robust model development.
- **15. knee\_osteoporosis.** Sourced from Mendeley Data, the *knee\_osteoporosis* dataset Gobara (2023b) contains X-rays categorized into three classes: normal, osteopenia, and osteoporosis. It is intended for studies on bone density assessment and osteoporosis detection.
- **16. RNSA\_pneumonia.** A pre-processed version of the RSNA Pneumonia Detection Challenge dataset, *RNSA\_pneumonia* Tapendu (2023b) includes PNG images and mask-based bounding box annotations. Associated metadata, such as patient information and bounding box coordinates, is provided in CSV format for easy integration into machine learning pipelines.
- **17. 8\_object\_detection.** The *Chest X-ray 8 Subset* Spritan1 (2023) is tailored for object detection in thoracic diseases, containing 790 images with 984 bounding boxes. Annotations are available in YOLO and Pascal VOC formats, and the dataset includes 14 thoracic disease classes, facilitating the development of object detection models in medical imaging.

- **18. HBFMID.** The *Human Bone Fractures Multi-modal Image Dataset* (HBFMID) Orvile (2023c) includes 1,539 annotated images (X-ray and MRI) covering fractures at multiple anatomical sites. The dataset is divided into training, validation, and testing sets and has undergone preprocessing (auto-orientation, resizing, contrast adjustments), supporting research in multi-modal fracture diagnosis.
- **19. FracAtlas.** FracAtlas Gupta (2023) comprises over 14,000 X-ray scans collected from three major hospitals in Bangladesh, with 4,083 images manually annotated for bone fracture classification, localization, and segmentation. Annotations were conducted by expert radiologists and validated by a medical officer, providing a high-quality benchmark for fracture analysis.
- **20. pneumonia.** The *pneumonia* dataset Mooney (2018) contains 5,863 chest X-ray images (anterior-posterior) of pediatric patients, labeled as either pneumonia or normal. Images underwent strict quality control and multi-expert grading, making the dataset reliable for training AI systems in pneumonia detection.
- **21.** pax\_ray. The *PAX-Ray++* dataset Seibold (2023) contains 7,377 chest radiographs (frontal and lateral views), with pseudo-labeled annotations for anatomical segmentation generated from projected thorax CT scans. The dataset is designed for segmentation tasks in chest X-ray analysis.
- **22. lung\_segmentation.** This dataset Beosup (2023) consists of over 500 X-ray scans labeled by radiologists, supporting machine learning research in lung region segmentation.
- **23. shadow.** The *shadow* dataset Hmchuong (2023) includes normal and bone-suppressed chest X-ray images, along with augmented samples. It is intended for research on bone shadow suppression to aid in lung disease diagnosis.
- **24. Angiography.** The *ARCADE* dataset Manaenkov (2023) features 3,000 X-ray coronary angiography frames with expert annotations for vessel segmentation, SYNTAX scoring, and stenosis detection. It is organized by task and includes cross-validated annotations, providing a rich resource for AI research in coronary artery disease diagnostics.
- **25. dental\_panoramic.** This panoramic dental radiograph dataset Lokisilvres (2023) includes segmentation masks for 31 dental disease classes, such as caries, crowns, implants, bone loss, fractures, and more. It is intended for comprehensive dental disease detection and segmentation research.
- **26. ALHI.** The *ALHI* dataset Rahman (2022) is a curated collection of 200 hip implant X-ray images from various medical sources, annotated and validated by orthopedic and clinical experts. The dataset includes images with diverse implant types and clinical conditions, supporting research on hip implant assessment.
- **27. humerus\_fractures.** The *humerus\_fractures* dataset Paspuel (2024) contains X-ray images depicting both fractured and non-fractured humeri, supporting automated diagnosis of humerus fractures through deep learning.
- **28.** multiclass\_knee\_osteoporosis. This dataset Gobara (2023a) offers X-ray images and patient records classified into normal, osteopenia, and osteoporosis categories, facilitating the automated diagnosis and classification of knee osteoporosis.
- **29. rsna-breast-cancer-detection.** The *RSNA Breast Cancer Detection* dataset Thakur (2024) provides breast X-ray image regions of interest (ROIs) in PNG format, without labels, for studies on automated detection in breast imaging.
- **30. RANZCR.** The *RANZCR* dataset RANZCR (2021) is intended for detecting the presence and position of catheters and lines on chest X-rays. It contains image IDs, binary labels for multiple types of catheters, and patient identifiers, along with associated CSV metadata.

- **31. FractureFusion.** FractureFusion Dutta (2023) is a diverse dataset capturing a wide variety of bone fracture cases, including avulsion, comminuted, greenstick, and spiral fractures, suitable for developing comprehensive fracture classification models.
- **32. HeelBone.** The *Heel Bone X-Ray* dataset Taher (2023) comprises 3,956 foot X-rays labeled for normal, heel spur, and severe heel spur complications. Images were sourced from Kirkuk General Hospital and cross-verified by orthopedic and radiology specialists, supporting disease classification in foot imaging.

### HISTOPATHOLOGY:

The histopathology dataset collection in FedAgentBench covers a wide range of diseases and task types, making it a comprehensive benchmark for evaluating LLM agents in digital pathology. It spans various cancer types, including breast (e.g., breast\_histo, BreaKHis\_400X, BreCaHAD), ovarian (ovarian\_cancer), gastric (gastric\_cancer), kidney (kmc\_kidney), melanoma, and nasopharyngeal carcinoma (NPC-88k-Public). The datasets support multiple learning paradigms such as binary and multi-class classification (lung\_and\_colon, EBHI), segmentation (MonuSeg, PanNuke), detection of mitotic figures (ULMS), and multimodal image-to-text learning (histo-img-text). Some datasets, like choledoch, incorporate hyperspectral imaging, while others like CellNet aggregate thousands of high-resolution images across organ types, facilitating generalization studies. Fine-grained annotations by expert pathologists (e.g., in BreCaHAD, NPC-88k-Public, MonuSeg) add clinical reliability. Together, these datasets reflect a realistic landscape of digital histopathology rich in diagnostic complexity, varied in modality and scale, and suitable for evaluating both general-purpose and specialized LLM agents in federated clinical settings. The exact dataset descriptions for each file are available in Listing 8 and summarized as follows:

- **1. breast\_histo.** The *Breast Histopathology Images* dataset Mooney (2024) focuses on Invasive Ductal Carcinoma (IDC), the most common breast cancer subtype. The original dataset comprises 162 whole mount slides scanned at 40x magnification, from which 277,524 patches of size  $50 \times 50$  were extracted (198,738 IDC negative and 78,786 IDC positive). Patch filenames encode patient ID, spatial coordinates, and IDC class (0 for non-IDC, 1 for IDC). Only images are provided, with no additional labels.
- **2. BreaKHis\_400X.** The *BreaKHis\_400X* dataset Forderation (2024) is derived from the BreaKHis database, which contains microscopic biopsy images of benign and malignant breast tumors. This subset includes images acquired at 400x optical zoom, with training and test data stored in separate folders. Images only are provided; no labels are included.
- **3. lung\_and\_colon.** The *Lung and Colon Cancer Histopathological Images* dataset MVD (2024a) contains 25,000 JPEG images of size 768 × 768 pixels, covering five classes: lung benign tissue, lung adenocarcinoma, lung squamous cell carcinoma, colon adenocarcinoma, and colon benign tissue. Images were generated from HIPAA-compliant and validated original samples (750 lung and 500 colon images) and augmented using the Augmentor package to create a balanced dataset of 5,000 images per class.
- **4. gastric\_cancer.** The *Gastric Cancer Histopathology Tissue Image Dataset* (GCHTID) Orvile (2024) comprises 31,096 non-overlapping images (224 × 224 pixels), extracted from H&E-stained pathological slides from Harbin Medical University Cancer Hospital. Images are categorized into eight tissue types, including adipose, background, debris, lymphocytes, mucus, smooth muscle, normal colon mucosa, cancer-associated stroma, and tumor, enabling research on the tumor microenvironment in gastric cancer.
- **5. gastro\_cancer\_msi\_vs\_mss.** The *Gastrointestinal Cancer MSI MSS Prediction* dataset Justin (2024) contains histological images for the classification of microsatellite instability (MSI) versus microsatellite stability (MSS) in gastrointestinal cancer, supporting research in histopathology image analysis with CNNs and transfer learning.

- **6. breast\_cancer\_segmentation.** The *Breast Cancer Cell Segmentation* dataset MVD (2024b) contains 58 H&E stained histopathology images with expert annotations for breast cancer cell detection and segmentation. The challenging task is cell segmentation for subsequent classification into benign and malignant cells, supported by ground truth data for algorithm development.
- **7. ovarian\_cancer.** The *Ovarian Cancer & Subtypes Dataset Histopathology* Pieces (2024) contains histopathology images representing four subtypes of ovarian cancer as well as non-cancerous tissue. The dataset is referenced as: Kasture, Kokila (2021), "Ovarian-Cancer&SubtypesDatasetHistopathology", Mendeley Data, V1, doi: 10.17632/kztymsrjx9.1.
- **8. breast\_cancer\_histo.** The *Breast Cancer Histopathology* dataset Kumar (2024) includes JPG images labeled as benign or malignant, supporting automated breast cancer classification from histopathological images.
- **9. BreCaHAD.** The *BreCaHAD* (Breast Cancer Histopathological Annotation and Diagnosis) dataset TruthIsNeverLinear (2024) comprises 162 annotated H&E-stained images, supporting automated classification of histological structures into six classes: mitosis, apoptosis, tumor nuclei, non-tumor nuclei, tubule, and non-tubule. See: https://bmcresnotes.biomedcentral.com/articles/10.1186/s13104-019-4121-7.
- **10. melanoma.** The *melanoma* dataset Haashaatif (2024) is designed for the development of deep learning models for nuclei and tissue segmentation in melanoma H&E-stained histopathology. It addresses challenges of melanocyte mimicry and includes nuclei and tissue annotations to facilitate studies on tumor-infiltrating lymphocytes and predictive/prognostic tasks.
- 11. choledoch. The *Choledoch* dataset HFUTYBX (2024) introduces both microscopy hyperspectral and color images for cholangiocarcinoma, including 880 scenes from 174 individuals (689 partial cancer, 49 full cancer, 142 non-cancer). All cancer areas are precisely labeled by experienced pathologists. More information is available in: https://ieeexplore.ieee.org/document/8869757. The dataset includes suggested train/val/test splits.
- **12. histopath-sn.** The *histopath-sn* Kaggle dataset Feng (2024) focuses on classifying patches and patients from bronchus and lung samples. Both images and labels are provided, with recommended train and test splits given in train\_labels.csv and test\_labels.csv.
- **13. ULMS.** The *Uterine Leiomyosarcoma (ULMS)* dataset Lee (2024) targets mitosis detection in ULMS, the most common uterine sarcoma. Images were collected in collaboration with pathologists and annotated for mitosis, aiding AI-based approaches for automatic mitosis detection and grading.
- **14. MonuSeg.** The *MonuSeg* dataset Dinh (2024) comprises 24 training images (originally 30,  $1000 \times 1000$  pixels) with 21,623 annotated nuclei from seven organs, and a test set of 58 images (8 from MonuSeg, 50 from the TNBC dataset). Annotations were made by one expert pathologist and two research fellows using consensus peer review.
- **15.** kmc\_kidney. The *KMC Kidney Histopathology* dataset Dwivedi (2024) includes non-cancerous (Grade-0) and cancerous (Grades 1-4) images of renal clear cell carcinoma, collected at Kasturba Medical College (KMC), India. Images were stained with H&E and labeled according to grade, supporting studies in kidney cancer histopathology.
- **16. histo-img-text.** The *histo-img-text* dataset Reasat (2024) comprises histopathology image-text pairs, including over 32k PNGs, 40k JPGs, and a CSV file with 217,052 captioned image entries. The dataset is designed for multimodal studies, such as image-to-text and vision-language modeling.
- **17. cellnet.** *CellNet* is a large, curated dataset Capocyan (2024) featuring over 120,000 high-quality medical images from more than 20 organ/cancer classes. Images were aggregated from diverse repositories and medical labs, supporting comprehensive research in computational pathology.

**18. PanNuke.** The *PanNuke* dataset Lad (2024) is a semi-automatically generated nuclei instance segmentation and classification dataset. It covers 481 visual fields across 19 tissue types, containing 205,343 labeled nuclei with segmentation masks, enabling tissue type segmentation and generalization to new tissue domains.

**19. NPC-88k-Public.** The *NPC-88k-Public* dataset Munirah (2024) includes 88,000 histopathology patches from 17 whole slide images across three institutions. Annotated regions include normal, lymphoid hyperplasia (LHP), nasopharyngeal inflammation (NPI), and nasopharyngeal carcinoma (NPC), with concordance among at least two pathologists.

**20. EBHI.** The *EBHI* dataset Alibabaei78 (2024) comprises 4,456 histopathology images and corresponding ground truth segmentations, including normal, polyp, low-grade and high-grade intraepithelial neoplasia, serrated adenoma, and adenocarcinoma. Images are paired with ground truth labels to support segmentation and classification research.

### MRI:

Our collection of 28 Magnetic Resonance Imaging (MRI) datasets supports a diverse array of machine learning tasks such as binary and multi-class classification, anatomical and pathological segmentation, anomaly detection, multi-modal image registration, and physiological parameter estimation. The included datasets range from unlabeled brain scans (*Brain MRI Images*) to richly annotated clinical benchmarks such as *BraTS*, *WMH*, and *ISLES 2015*, covering tumor segmentation, white matter lesion detection, and ischemic stroke assessment. Cardiac datasets like *ACDC* facilitate diagnosis of specific heart conditions, while spine-related datasets such as the *RSNA 2024 Lumbar Spine Challenge* and *Foraminal Stenosis MRI* target degenerative spinal diseases. Other specialized collections, including *Facial MRI*, *Prostate MRI*, and multi-modal datasets (e.g., *MRI-PET Brain Scans*), enable crossdomain generalization and analysis. Together, this curated set of MRI datasets provides a foundation for training and benchmarking AI systems across a broad range of anatomical regions and diagnostic challenges.

**1. Brain MRI Images** A Kaggle dataset bra (b) containing diverse brain MRI images sourced from multiple datasets, offering a range of anatomical variations and imaging contrasts.

**2. Alzheimer Classification** Brain MRI dataset alz labeled for Alzheimer's disease classification into four categories: Mild Demented, Moderate Demented, Non-Demented, and Very Mild Demented.

**3. Brain Cancer** Brain MRI images bra (a) collected from hospitals in Bangladesh for classification into Brain Glioma, Brain Meningioma, and Pituitary Tumor classes.

**4. Brain Tumour** A labeled brain tumor dataset bra (e) for binary classification (tumor vs. non-tumor) and unlabeled prediction samples for testing.

**5. 4 Class Brain Tumour** A brain MRI dataset bra (d) for classifying tumors into Benign, Malignant, and Pituitary types.

**6. Heat MRI Left Atrial Segmentation** A segmentation dataset hea of left atrial structures in cardiac MRI provided by King's College London.

**7. PMRAM** MRI brain cancer dataset pmr with four classes (Glioma, Meningioma, Pituitary, No Tumor), standardized to 512×512 resolution and augmented from 1600 base images.

**8. Hippocampal Sparing** Unlabeled DICOM-format MRI slices hip of 25 patients for hippocampal sparing studies, organized per patient.

**9. Spine** Spine MRI scans spi from a single patient with labeled dystrophic anomalies and accompanying radiology reports.

- 10. Brain Tumour CT MRI A brain tumor dataset bra (c) composed of both MRI and CT images, labeled for tumor detection and suitable for training diagnostic models.
  - **11. BraTS 2019** Multimodal brain MRI dataset Menze et al. (2014) (T1, T1Gd, T2, FLAIR) with expert segmentations for tumor subregions, formatted as NIfTI (.nii.gz) files.
  - **12. Bone Fractures MRI X-ray** Multi-modal dataset hbf including MRI and X-ray scans for bone fracture detection across different body regions.
  - **13. Alzheimer Detection** Preprocessed MRI scans LaMontagne et al. (2019) from the OASIS-1 dataset labeled for Alzheimer's detection tasks.
  - **14. Stroke Head MRI** MRI brain scans str with segmentations of stroke lesions from patients with cerebrovascular conditions.
  - **15. MRI PET Brain Scans** Paired MRI and PET DICOM scans mri for brain tumors, aimed at multi-modal registration and Dice score evaluation.
  - **16. OASIS-1** Processed MRI scans of 1688 subjects across Alzheimer's Disease (AD), Cognitively Normal (CN), and Mild Cognitive Impairment (MCI) groups oas.
  - **17. Abdomen MRI** Abdominal MRI dataset abd with object detection annotations and bounding boxes in CSV format.
  - **18. Facial MRI** Facial MRI scans fac including sagittal and axial slices for anomaly detection, segmentation, and 3D anatomical modeling.
  - **19. Prostate** Multi-parametric prostate MRI scans pro (a) with manual segmentations for clinical segmentation research.
  - **20. Glioma** TCGA-LGG-based MRI dataset gli for low-grade glioma detection with segmentation masks and associated genomics metadata.
  - **21. Phantom** Longitudinal MRI dataset pha of a single healthy subject scanned on 116 scanners over 2.5 years to analyze scanner variability.
  - **22.** ACDC: Automated Cardiac Diagnosis Challenge Dataset The ACDC Bernard et al. (2018) dataset consists of cine-MRI scans, categorized into five balanced cardiac pathology classes: Normal (NOR), Myocardial Infarction (MINF), Dilated Cardiomyopathy (DCM), Hypertrophic Cardiomyopathy (HCM), and Abnormal Right Ventricle (ARV). Each class is defined by specific clinical parameters such as ejection fraction, wall thickness, and ventricular volumes, supporting robust machine learning development for automated cardiac function assessment.
  - **23. Foraminal Stenosis MRI Dataset** This dataset for comprises high-resolution lumbar spine MRI scans with segmentation masks and foraminal measurements, aimed at detecting and analyzing foraminal stenosis. It supports tasks such as nerve channel size analysis, stenosis classification, and monitoring of spinal degenerative conditions, enabling precise anatomical assessment and aiding in early diagnosis and treatment planning.
  - **24. RSNA 2024 Lumbar Spine Degenerative Classification Challenge** This RSNA-ASNR RSN dataset includes five lumbar spine degenerative conditions—Left/Right Neural Foraminal Narrowing, Left/Right Subarticular Stenosis, and Spinal Canal Stenosis—using lumbar spine MRI. The dataset includes severity scores (Normal/Mild, Moderate, Severe) across five disc levels (L1/L2 to L5/S1), enabling automated classification to support diagnosis and treatment planning for lower back pain and related conditions.

- 25. ATLAS v2.0 The Anatomical Tracings of Lesions After Stroke (ATLAS) v2.0 Liew et al. (2022) dataset provides manually segmented T1-weighted MRI scans of individuals with stroke lesions. It includes lesion masks and anatomical metadata for over 600 subjects, with the aim of facilitating the development and benchmarking of automated stroke lesion segmentation methods.
  - **26. BraTS** The Brain Tumor Segmentation (BraTS) dataset provided through the Medical Segmentation Decathlon (MSD), comprises multi-modal MRI scans (T1, T1-Gd, T2, and FLAIR) of glioma patients with expert annotations of tumor sub-regions including the enhancing tumor, peritumoral edema, and necrotic core.
  - **27. WMH** The White Matter Hyperintensities (WMH) dataset wmh consists of T1 and FLAIR MRI scans from multiple institutions with voxel-wise annotations of WMH regions. Originally compiled for the WMH Segmentation Challenge at MICCAI 2017, the dataset captures variability across scanners and populations, making it a robust benchmark for automated WMH detection methods.
  - **28. ISLES 2015 (SISS)** The Ischemic Stroke Lesion Segmentation (ISLES) 2015 challenge dataset isl, specifically the Sub-Acute Ischemic Stroke Lesion Segmentation (SISS) subtask, offers multimodal MRI scans (including FLAIR, T1, DWI) with corresponding lesion masks for patients in the subacute phase post-stroke. It supports the development of methods for accurate ischemic stroke lesion segmentation and includes cases with diverse lesion locations and volumes.

### FUNDUS:

- Our Fundus image datasets span a broad range of tasks and clinical applications, reflecting the diagnostic richness of retinal imaging. These include segmentation datasets such as *Drishti-GS*, *RIMONE*, and *ONH Segmentation* for optic disc/cup analysis in glaucoma, and vessel segmentation benchmarks like *DRIVE* and *CHASE\_DB1* for vascular assessment. Classification datasets such as *APTOS*, *MESSIDOR*, and *ARIA* support diabetic retinopathy grading, while multi-label datasets like *RFMID* and *ODIR-5K* address a broader set of ocular diseases. Lesion-level annotations in datasets like *IDRiD* and *E-Ophtha* enable fine-grained detection of diabetic pathologies. Additionally, niche datasets such as *e-ROP*, *Ocular Toxoplasmosis*, and *AMDP* target rare or longitudinal conditions. Others focus on preprocessed imaging (*CLAHE* + *ESRGAN Split FD*) or multi-modal metadata (*SMDG*, *DrHagis*). This diversity supports robust benchmarking across segmentation, classification, enhancement, and multimodal learning, forming the backbone of data-driven ophthalmic model development.
- **1. Drishti-GS** This dataset Sivaswamy et al. (2014) is used for glaucoma detection, providing optic disc and cup segmentation masks. It supports both segmentation and glaucoma classification tasks.
- **2. STARE** The STARE dataset STA is used for retinal disease diagnosis and retinal vessel segmentation. Its main tasks include vessel segmentation and lesion detection.
- **3. IDRiD** The Indian Diabetic Retinopathy Image Dataset (IDRiD) ind provides pixel-level annotations for diabetic retinopathy (DR) lesions. It is used for lesion segmentation and DR grading.
- **4. DR** This dataset DR is used for classifying diabetic retinopathy across 5 severity levels.
- **5. RIMONE** A glaucoma dataset Fumero et al. (2011) providing optic disc and cup annotations, mainly used for segmentation and glaucoma classification.
- **6. REFUGE** A unified glaucoma evaluation dataset ref, widely used for optical disc/cup segmentation and glaucoma classification.
- **7. CHASE\_DB1** This dataset cha contains child retinal images with annotated vessels. It is primarily used for vessel segmentation tasks.
- **8. E-Ophtha** Designed for diabetic retinopathy research, this dataset Decenciere et al. (2013) includes images annotated for exudates and hemorrhages, supporting lesion detection.

- 9. ARIA A retinal image dataset used in diabetic retinopathy screening. It is mainly employed for DR classification.
- 10. IOSTAR A dataset of multi-modal retinal images, particularly used for optic disc segmentation tasks.
- 11. HRF The High-Resolution Fundus dataset is used for both vessel and optic disc segmentation, offering detailed structural annotations.
- 12. LES-AV This dataset supports artery and vein classification, distinguishing vessel types in fundus images.
  - **13. PRIME-FP20** It is a high-resolution dataset of fundus images used for optic disc segmentation.
- **14. RIGA+** This is a glaucoma dataset derived from multiple sources, used for optic disc and cup segmentation.
- **15. APTOS** It is part of the Kaggle Diabetic Retinopathy Challenge (2019), this dataset is used to grade DR severity from fundus images.
- **16. MESSIDOR** It is a classic and widely used diabetic retinopathy dataset, primarily for classification tasks.
- **17. DRIVE** It is one of the earliest vessel segmentation datasets, often used as a benchmark in fundus segmentation.
  - **18. ORIGA** The ORIGA dataset provides optic disc and cup annotations for segmentation task and glaucoma detection.
    - **19. ODIR-5K** The ODIR (Ocular Disease Intelligent Recognition) dataset contains over 5,000 retinal fundus images with multi-label annotations for eight ocular diseases, including diabetic retinopathy, glaucoma, cataract, AMD, hypertension, and others. It supports multi-label classification tasks.
    - **20. RFMID** The Retinal Fundus Multi-Disease Image Dataset (RFMID) includes 3,200+ images annotated for 19 different conditions. It is intended for multi-label classification tasks and supports the development of fundus-based diagnostic models for diverse ocular diseases.
    - **21. MESSIDOR-2 DF** MESSIDOR-2 is the second edition of the MESSIDOR diabetic retinopathy dataset. It includes fundus images with diabetic retinopathy severity labels.
    - **22. Glaucoma datasets** (**EYE-PACS**) EYE-PACS is a large-scale dataset used primarily for diabetic retinopathy grading in the Kaggle challenge.
  - **23. Retina blood vessel segmentation dataset** This fundus dataset is used for vessel segmentation.
  - **24. DR Diagnosis Dataset** This dataset is used for classifying diabetic retinopathy severity based on retinal fundus images.
- 25. DDR Dataset The Diabetic Retinopathy Detection from Retina Images (DDR) dataset includes
   fundus images annotated for DR severity and pixel-level lesion types (e.g., exudates, hemorrhages).
   It supports both classification and lesion segmentation tasks.
  - **26. Hypertensive Retinopathy** This dataset contains fundus images annotated for signs of hypertensive retinopathy. While rare and usually hospital-specific, it is used for classification and grading of HR severity.

- 3402
   3403
   3404
   3404
   3405
   27. SUSTECH + SYSU Dataset This entry combines data from SUSTech and Sun Yat-sen Glaucoma, diabetic retinopathy, and related diseases. It supports classification tasks across multiple disease categories.
  - **28. RITE** The Retinal Images vessel Tree Extraction (RITE) dataset, derived from DRIVE, includes ground truth for artery and vein segmentation. It is used to differentiate between arterial and venous vessels in retinal images.
  - **29. CLAHE** + **ESRGAN Split FD** This dataset represents a preprocessed variant of fundus images where contrast enhancement (CLAHE) and super-resolution techniques (ESRGAN) have been applied. It is used to improve image quality for downstream classification tasks.
  - **30. Myopia Image Dataset** This dataset consists of retinal fundus images labeled for myopia classification.
  - **31. ACRIMA** ACRIMA is fundus dataset used for glaucoma detection.
  - **32.** and **33.** Retina Fundus Dataset (CHASE\_DB1, DRIVE) CHASE\_DB1 and DRIVE are fundus datasets used for retinal vessel segmentation, i.e., for segmenting blood vessels in fundus images.
  - **34.** Cataract Classification Dataset This is used for binary classification of cataract presence in fundus images.
  - **35. MURED** The Multicenter Retinal Disease Dataset (MURED) aggregates retinal images across multiple institutions and includes annotations for diabetic retinopathy, glaucoma, age-related macular degeneration (AMD), and other conditions. It is primarily used for multi-class classification of retinal diseases.
  - **36. Optic Disc Cup Fundus Image** This dataset contains annotations for optic disc and cup structures. These datasets are used for segmentation tasks and for calculating cup-to-disc ratio, an important indicator in glaucoma diagnosis.
  - **37. ROFT** This is a retinal and ocular fundus image dataset with 8 disease labels for fundus images normal, diabetes, glaucoma, cataract, age-related macular degeneration, hypertension, pathological myopia and other diseases/abnormalities. It also has 7 labels for OCT: age-related macular degeneration, diabetic macular edema, epiretinal membrane, normal, retinal artery occlusion, retinal vein occlusion, vitreomacular interface diseases.
  - **38.** Eye Disease Image Dataset A fundus dataset for detection of eye-related 10 conditions central serous chorioretinopathy, diabetic retinopathy, disc edema, glaucoma, healthy, macular scar, myopia, pterygium, retinal detachment, and retinitis pigmentosa.
  - **39. FIVES** The FIVES dataset (Fundus Image Vessel Extraction and Segmentation) is used for vessel segmentation tasks. It provides pixel-level annotations for blood vessel structures.
  - **40. AMDP Dataset** This refers to the Age-related Macular Degeneration Prediction dataset which is longitudinal ophthalmic dataset.
  - **41. AGAR 300** A Microaneurysms Fundus Dataset that consists of color fundus images showing signs of microaneurysms for early DR detection.
  - **42. SMDG** It is a standardized fundus glaucoma dataset consisting of full-fundus glaucoma images with image metadata on optic disc/cup segmentation and blood vessel segmentation.
  - **44. Fundus segmentation dataset** It is a unified retinal image dataset for assessing glaucoma with reference segmentation labels of optic disc and cup.

3458 3459

3460 3461

3462

3463

3464 3465

3466

3467 3468

3469

3470 3471

3472 3473

3474 3475

3476

3477 3478

3479 3480

3481

3482 3483

3484

3485 3486

3487

3488 3489

3490

3491 3492

3493

3494

3495 3496

3497

3498 3499

3500

3501

3502 3503

3504

3505

3506

3507

3508

- **45.** Hypertensive retinopathy dataset It is a fundus dataset for binary classification regarding 3457 presence or absence of hypertensive retinopathy.
  - **46. DR grading dataset** It is a fundus dataset for grading the severity of diabetic retinopathy.
  - 47. G1020 dataset It is a fundus dataset for glaucoma classification and contain 1020 high resolution colour fundus images. It also provides annotations for glaucoma diagnosis, optic disc and cup segmentation, vertical cup to disc ratio, etc.
  - **48.Ocular Toxoplasmosis dataset** It is a fundus dataset used for detection of Toxoplasmosis chorioretinitis and has three classes - healthy eye, active and inactive chorioretinitis.
  - **49. ARIA dataset** It is a fundus dataset used for detection of any of three classes: healthy, AMD and Diabetes.
  - **50. Fundus 4 categories dataset** It is a fundus dataset used for detection of normal, cataract, glaucoma and diabetic retinopathy.
  - **51. ONH Segmentation dataset** It is an optic disc and cup mask segmentation fundus dataset
  - **52. DrHagis dataset** It is a fundus dataset for detection of diabetic retinopathy, hypertension, age-related macular degenration and glaucoma.
  - **53. Driona DB dataset** It is a fundus dataset for optic disc segmentation.
  - **54. Cattle Retinal Fundus Images** A unique dataset featuring retinal fundus images from cattle, useful for comparative studies and veterinary ophthalmology research.
  - **55. Preprocessed Eye Diseases Fundus Images** It offers preprocessed fundus images enhanced using techniques like CLAHE and ESRGAN, facilitating improved classification performance.
  - **56. Retina Fundus Image Registration Dataset (FIRE)** It comprises 129 retinal images forming 134 image pairs, designed for evaluating image registration algorithms.
  - **57. 1000 Fundus Images with 39 Categories** This dataset comprises 1,000 fundus images categorized into 39 distinct classes, offering a diverse set for multi-class classification tasks.
  - **58. PAPILA: Retinal Fundus Images Dataset** The PAPILA dataset includes fundus images and clinical data from both eyes of individual patients for glaucoma assessment. It provides optic disc and cup segmentations, along with patient-level glaucoma labels derived from clinical evaluations.
  - **59. Diabetic Retinopathy Diagnosis Dataset** A large-scale retinal image dataset designed for the diagnosis of diabetic retinopathy, supporting medical image analysis and automated disease grading.
  - **60. Vessel Tree Extraction Dataset** This dataset supports comparative research on artery and vein segmentation or classification in retinal fundus images, facilitating the development and benchmarking of vessel-type analysis methods.
  - 61. DiaRetDB1: Diabetic Retinopathy Benchmark Dataset DiaRetDB1 includes retinal fundus images with expert-annotated ground truth for key lesions such as hard and soft exudates, microaneurysms, and hemorrhages, along with both the original images and raw annotation data.
  - **62.** SynFundus The SynFundus is a synthetic fundus dataset includes annotations for eleven ocular diseases: diabetic retinopathy, age-related macular degeneration, anomalies of the optic nerve, choroidal retinal pathology, degenerative and pathological myopia, diabetic macular edema, epimacular membrane, glaucoma, hypertensive retinopathy, and retinal vein occlusion. These

3511

3512 3513

3514

3515

3516 3517

3518 3519

3520

3521

3522

conditions cover a broad range of structural and vascular retinal abnormalities, supporting diverse diagnostic research in ophthalmology.

**63. AIROGS** The AIROGS dataset De Vente et al. (2023) comprises fundus photographs from diverse ethnicities and imaging devices. It supports two main tasks: referable glaucoma classification and detection of ungradable images to simulate real-world screening conditions.

## C.2 SAMPLE DATASET DESCRIPTION FILES:

Sample dataset description files are shown in Listings 6-8. The datasets are then partitioned into different clients and utilized by the client selector agents to decide whether to choose the client for federated analysis.

Listing 6: Dataset Descriptions for Dermatology Modality

```
3524
       [
3525
3526
               "Dataset Name": "augmented_skin_condition_dataset_kaggle",
3527
               "Dataset Description": "augmented_skin_condition_dataset_kaggle
3528
                   is a skin disease classification dataset containing images of
                    six different dermatological conditions: 'Acne', 'Carcinoma
3529
                   ', 'Eczema', 'Keratosis', 'Milia', and 'Rosacea'. It contains
3530
                    six subfolders, with each subfolder containing images of the
3531
                    corresponding class (disease) specified in the name of the
3532
                   subfolder. ",
3533
               "Dataset
                   Path": "skin_dataset/augmented_skin_condition_dataset_kaggle"
3534
           },
3535
3536
               "Dataset Name": "DDI_skin_dataset",
3537
               "Dataset Description": "DDI_skin_dataset is a skin cancer
                   classification dataset with diverse skin tone representation
                   that contains 1 subfolder 'images' and 2 CSV files. Focus on
3539
                   the columns: 'DDI_file' (for the image path) and 'malignant'
                   (the class label) of the csv file 'ddi_metadata.csv'. 'True'
3541
                   in 'malignant' column means malignant whereas 'False' means
3542
                   benign. ",
               "Dataset Path": "skin_dataset/DDI_skin_dataset"
3543
           },
3544
3545
               "Dataset Name": "Derma7PT",
3546
               "Dataset Description": "Derma7PT is a skin disease classification
3547
                    dataset containing a subfolder 'images' and a csv file 'meta
                   .csv'. Focus on the columns 'clinic' and 'derm' for the image
3548
                    file path as well as the column 'diagnosis' of the csv file
3549
                   that has 10 disease types: 'basal cell carcinoma', 'nevus',
3550
                   dermatofibroma', 'lentigo', 'melanoma', 'melanoma metastasis
', 'melanosis', 'miscellaneous', 'seborrheic keratosis', '
3551
3552
                   vascular lesion'. ",
               "Dataset Path": "skin_dataset/Derma7P"
3553
           },
3554
3555
               "Dataset Name": "Dermatology_tabular dataset",
3556
               "Dataset Description": "Dermatology_tabular dataset is a tabular
                   (non-image) dataset containing clinical features for
3558
                   diagnosing skin diseases. ",
               "Dataset Path": "skin_dataset/Dermatology_tabular dataset"
3559
           },
3560
3561
               "Dataset Name": "Dermis",
3562
               "Dataset Description": "Dermis is a skin disease dataset with
                   benign and malignant cases, supporting both classification
                   and segmentation tasks. It has two sub-folders 'benign' and
```

```
3564
                  melanoma'. In each of these sub-folders, we have two sub-
3565
                   folders 'contour' (that has the segmentation masks) and '
3566
                   images' (that has the original images). ",
               "Dataset Path": "skin_dataset/Dermis"
3567
           },
3568
3569
               "Dataset Name": "Dermnet",
3570
               "Dataset Description": "Dermnet contains a very broad collection
3571
                   of skin disease images. It has 23 sub-folders covering 23
                   disease categories namely 'Acne and Rosacea', 'Actinic
3572
                   Keratosis Basal Cell Carcinoma and other Malignant Lesions',
3573
                   'Atopic Dermatitis Photos', 'Bullous Disease Photos', '
3574
                   Cellulitis Impetigo and other Bacterial Infections', 'Eczema
3575
                  Photos', 'Exanthems and Drug Eruptions', 'Hair Loss (Alopecia
3576
                   ) and other Hair Diseases', 'Herpes HPV and other STDs Photos
                   ', 'Light Diseases and Disorders of Pigmentation', 'Lupus and
3577
                   other Connective Tissue Diseases', 'Melanoma Skin Cancer
3578
                  Nevi and Moles', 'Nail Fungus and other Nail Disease', '
3579
                  Poison Ivy Photos and other Contact Dermatitis', 'Psoriasis
3580
                  pictures and Lichen Planus and related Diseases', 'Scabies
                  Lyme Disease and other Infestations and Bites', 'Seborrheic
3581
                  Keratoses and other Benign Tumors', 'Systemic Disease', '
3582
                  Tinea Ringworm Candidiasis and other Fungal Infections', '
3583
                  Urticaria Hives', 'Vascular Tumors', 'Vasculitis Photos', '
3584
                  Warts Molluscum and other Viral Infections'. ",
3585
               "Dataset Path": "skin dataset/Dermnet"
3586
           },
3587
               "Dataset Name": "Dermquest",
3588
               "Dataset Description": "Dermquest is a skin disease
3589
                  classification and segmentation dataset containing images of
3590
                  benign and malignant skin diseases. It has two sub-folders '
3591
                  benign' and 'melanoma'. In each of these sub-folders, we have
                   two sub-folders 'contour' (that has the segmentation masks)
3592
                   and 'images' (that has the original images). ",
3593
               "Dataset Path": "skin_dataset/Dermquest"
3594
           },
3595
3596
               "Dataset Name": "fitzpatrick17k",
               "Dataset Description": "fitzpatrick17k is a large skin lesion
3597
                  dataset with a wide range of dermatological diseases. It has
3598
                   a sub-folder 'finalfitz17k' which contains all images and two
3599
                   csv files 'fitzpatrick17k_disease.csv' and '
3600
                   Fitzpatrick17k_morphology.csv'. Focus on the column 'md5hash'
3601
                   for filename and the column 'three_partition_label' that
                   contains three disease labels: 'non-neoplastic', 'benign',
3602
                  malignant' in the file 'fitzpatrick17k_disease.csv'. ",
3603
               "Dataset Path": "skin_dataset/fitzpatrick17k"
3604
          },
3605
3606
               "Dataset Name": "ISIC2018_HAM10000",
               "Dataset Description": "ISIC2018 HAM10000 is a skin lesion
3607
                   classification and segmentation dataset. It has a sub-folder
3608
                   'HAM10000_images_combined_600x450' that contains original
3609
                   images as well as a sub-folder 'HAM10000_segmentations_mask'
3610
                   that contains the corresponding segmentation masks. The
3611
                   classification labels can be found in the 'dx' column of the
                   csv file 'ISIC2018_Task3_Test_GroundTruth.csv' including '
3612
                  Melanocytic Nevus (nv)', 'Benign Keratosis-like Lesions (bkl)
3613
                   ', 'Melanoma (mel)', 'Basal Cell Carcinoma (bcc)', 'Actinic
3614
                  Keratosis / Bowen's Disease (akiec)', 'Vascular Lesions (vasc
3615
                   )', 'Dermatofibroma (df)'. The corresponding image names can
3616
                  be found in the column 'image_id' of the same csv file. ",
               "Dataset Path": "skin_dataset/ISIC2018_HAM10000"
3617
          },
```

```
3618
3619
               "Dataset Name": "ISIC_2016",
3620
               "Dataset Description": "ISIC_2016 is a skin lesion dataset for
3621
                   classification and segmentation, focused on skin cancer
                   detection. It has two sub-folders 'ISBI2016_ISIC_images' that
3622
                   contain original images and 'ISBI2016_ISIC_segmentation_mask
3623
                   ' that has segmentation masks. The csv file '
3624
                   ISBI2016_ISIC_binary_classification_Training_GroundTruth.csv'
3625
                   has two columns - the first column being image names and
3626
                   second column being binary disease labels: 'benign' and '
                  malignant'. ",
3627
               "Dataset Path": "skin_dataset/ISIC_2016"
3628
           },
3629
3630
               "Dataset Name": "ISIC_2017",
               "Dataset Description": "ISIC_2017 is a skin lesion classification
3631
                   and segmentation dataset with a focus on melanoma and
3632
                   seborrheic keratosis diagnosis. It has two sub-folders: '
3633
                   images' that contain original images and 'Segmentation_masks'
3634
                    that has segmentation masks. There is a csv file 'ISIC-2017
3635
                   GroundTruth with the columns 'image id' that contains image
                   filenames, 'melanoma' that contains binary labels
3636
                   corresponding to presence (1) and absence (0) of melanoma,
3637
                   and 'seborrheic keratosis' that contains binary labels
3638
                   corresponding to presence (1) and absence (0) of seborrheic
3639
                   keratosis. ",
3640
               "Dataset Path": "skin_dataset/ISIC_2017"
3641
           },
3642
               "Dataset Name": "ISIC_2019",
3643
               "Dataset Description": "ISIC_2019 is an extended skin disease
3644
                   classification dataset. It has one sub-folder: 'images' that
3645
                   contain original images. In the CSV file '
                   ISIC_2019_Training_GroundTruth.csv', the 'image' column
3646
                   contains the image file names and 9 other columns represent
3647
                   the presence (1) or absence (0) of 9 classes namely Melanoma
3648
                   (MEL), Nevus (NV), Basal Cell Carcinoma (BCC), Actinic
3649
                   Keratosis / Bowen's Disease (AK), Benign Keratosis-like
3650
                   Lesions (BKL), Dermatofibroma (DF), Vascular Lesions (VASC),
                   Squamous Cell Carcinoma (SCC) and Unknown (UNK). ",
3651
               "Dataset Path": "skin_dataset/ISIC_2019"
3652
3653
3654
               "Dataset Name": "ISIC_2020",
3655
               "Dataset Description": "ISIC_2020 is a binary classification
                  dataset of skin lesions (benign vs malignant). It has one sub
3656
                   -folder: 'images' that contain original images. In the CSV
3657
                   file 'ISIC 2020 Training GroundTruth.csv', the 'image name'
3658
                   column contains the image file names and the '
3659
                  benign_malignant' column contains the corresponding disease
3660
                   labels on malignant or benign. ",
               "Dataset Path": "skin_dataset/ISIC_2020"
3661
           },
3662
3663
               "Dataset Name": "ISIC_2024",
3664
               "Dataset Description": "ISIC_2024 is an updated ISIC skin disease
3665
                    dataset primarily for melanoma classification (binary:
                   benign vs malignant). It has one sub-folder: 'images' that
3666
                   contain original images. In the CSV file '
3667
                   {\tt ISIC\_2024\_Training\_GroundTruth.csv',\ the\ 'isic\_id'\ column}
3668
                   contains the image file names and the 'malignant' column
3669
                   contains the corresponding disease labels on malignant or
3670
                   benign. '0' means benign and '1' means malignant. ",
               "Dataset Path": "skin_dataset/ISIC_2024"
3671
           },
```

```
3672
3673
               "Dataset Name": "Mednode",
3674
               "Dataset Description": "Mednode is a skin disease dataset for
3675
                   binary classification. It has 2 sub-folders covering 2
                   disease categories namely melanoma and nevus. ",
3676
               "Dataset Path": "skin_dataset/Mednode"
3677
           },
3678
3679
               "Dataset Name": "Monkeypox_Skin_Image_Dataset",
               "Dataset Description": "Monkeypox_Skin_Image_Dataset is a dataset
3680
                    for skin disease classification and has four sub-folders (
3681
                   with data belonging to the corresponding disease category)
3682
                   named: 'Chickenpox', 'Measles', 'Monkeypox', and 'Normal'. ",
3683
               "Dataset Path": "skin_dataset/Monkeypox_Skin_Image_Dataset"
3684
           },
3685
               "Dataset Name": "PAD_UFES_20",
3686
               "Dataset
3687
                   Description": "PAD_UFES_20 is a skin disease classification
3688
                   dataset. It contains a sub-folder 'images' containing the
                   original images and a csv file called 'metadata.csv' that
3689
                   contains the image ids under the column 'img_id' and disease
3690
                   labels under the column 'diagnostic' which contains 6 disease
3691
                    categories with corresponding abbreviations: Melanoma (MEL),
3692
                    Melanocytic Nevus (NEV), Basal Cell Carcinoma (BCC), Actinic
3693
                    Keratosis / Bowen's Disease (ACK), Seborrheic Keratosis (SEK
                   ), and Squamous Cell Carcinoma (SCC). ",
3694
               "Dataset Path": "skin_dataset/PAD_UFES_20"
3695
           },
3696
3697
               "Dataset Name": "PH2Dataset",
3698
               "Dataset Description": "PH2Dataset is a skin lesion
                   classification and segmentation dataset. It has a sub-folder
3699
                   'PH2 Dataset images' which in turn has two sub-folders'
3700
                   all_dermoscopic_images' that contain all the original images
3701
                   and 'segmentation_mask' that contain all the segmentation
3702
                   masks. The folder has an xlsx file called 'PH2_dataset.xlsx'
3703
                   with a column called 'Image Name' that contains the image ids
                    and a column 'Clinical Diagnosis' three disease classes : '
3704
                   Common Nevus', 'Atypical Nevus', and 'Melanoma' marked with '
3705
                   X' whenever that category is present in a given image. ",
3706
               "Dataset Path": "skin_dataset/PH2Dataset"
3707
           },
3708
               "Dataset Name": "scin_dataset",
3709
               "Dataset Description": "scin_dataset is a multi-class skin
3710
                   disease classification dataset. It has a sub-folder '
3711
                   scin_images' that contains all the original images and two
3712
                   csv files. Follow the 'scin_cases.csv' file which has the
                   image ids in the column 'case_id' and the disease classes
3713
                   under the 'related category' which should include the 9
3714
                   diseases: 'RASH', 'LOOKS_HEALTHY', 'OTHER_ISSUE_DESCRIPTION',
3715
                   'NAIL_PROBLEM', 'GROWTH_OR_MOLE', 'ACNE', 'PIGMENTARY_PROBLEM', 'HAIR_LOSS', 'OTHER_HAIR_PROBLEM'. ",
3716
3717
               "Dataset Path": "skin_dataset/scin_dataset"
3718
           },
3719
               "Dataset Name": "skin_disease_3_class",
3720
               "Dataset Description": "skin_disease_3_class is a skin disease
3721
                   classification dataset that consists of a sub-folder 'images'
3722
                   which in turn has three sub-folders each consisting of one
3723
                   of the three classes indicated by the sub-folder name: 'acne
                   ', 'atopic dermatitis', and 'basal cell carcinoma'. ",
3724
               "Dataset Path": "skin_dataset/skin_disease_3_class"
3725
           },
```

```
3726
3727
               "Dataset Name": "skin_disease_classification_kaggle",
3728
               "Dataset Description": "skin_disease_classification_kaggle is a
3729
                   skin disease classification dataset with a sub-folder 'files'
                   that again contains three sub-folders each containing one of
3730
                   the three classes: 'acne', 'eye bags', and 'redness'. ",
3731
               "Dataset Path": "skin_dataset/skin_disease_classification_kaggle"
3732
           },
3733
3734
               "Dataset Name": "skin_disease_kaggle_dataset",
               "Dataset
3735
                   Description": "skin_disease_kaggle_dataset is a skin cancer
3736
                   detection dataset that has 10 sub-folders for 10 disease
3737
                   classes with the corresponding sub-folder names: 'Atopic
                   Dermatitis', 'Basal Cell Carcinoma (BCC)', 'Benign Keratosis-
3738
                   like Lesions (BKL)', 'Eczema', 'Melanocytic Nevi (NV)', '
3739
                  Melanoma', 'Psoriasis pictures Lichen Planus and related
3740
                  diseases', 'Seborrheic Keratoses and other Benign Tumors', '
3741
                   Tinea Ringworm Candidiasis and other Fungal Infections', and
3742
                   'Warts Molluscum and other Viral infections'. ",
               "Dataset Path": "skin_dataset/skin_disease_kaggle_dataset"
3743
           },
3744
3745
               "Dataset Name": "Skin Disease_Robo",
3746
               "Dataset Description": "Skin Disease_Robo is a skin disease
3747
                   classification and object detection dataset. It has one sub-
                   folder 'image' that contains all the original images and a
3748
                   csv file 'bounding_box_annotations.csv' with a column called
3749
                   'filename' that has all the image names and column 'class'
3750
                   that has 10 disease class labels: 'Acne', 'Atopic Dermatitis
3751
                   ', 'Chicken Skin', 'Eczema', 'Hansen's Disease-Leprosy', '
                  Hansen's Disease-Leprosy-severe', 'Healthy skin', 'Psoriasis
3752
                   ', 'Ringworm', 'Warts'. It also contains coordinates for
3753
                  bounding box annotations for lesions in the columns 'xmin', '
3754
                  ymin', 'xmax', and 'ymax'. ",
3755
               "Dataset Path": "skin_dataset/Skin Disease_Robo"
3756
           },
3757
               "Dataset Name": "skin-infection-disease-dataset",
3758
               "Dataset Description": "skin-infection-disease-dataset is a skin
3759
                   disease classification dataset focusing on infectious skin
3760
                   diseases. It has 8 sub-folders consisting diseases of each
3761
                   category - BA-cellulitis, BA-impetigo, FU-athlete-foot, FU-
3762
                  nail-fungus, FU-ringworm, PA-cutaneous-larva-migrans, VI-
                   chickenpox, VI-shingles. ",
3763
               "Dataset Path": "skin_dataset/skin-infection-disease-dataset"
3764
          },
3765
3766
               "Dataset Name": "skinL2_dataset",
               "Dataset Description": "skinL2_dataset is a skin cancer
3767
                   classification dataset with 8 sub-folders containing 8
3768
                   classes: 'Basal-cell Carcinoma', 'Dermatofibroma',
3769
                   Hemangioma', 'Melanoma', 'Nevus', 'Psoriasis', 'Seborrheic
3770
                   Keratosis', and 'Others'. Optional metadata is available in '
3771
                  PlenoISLA_DatasetV1_info.xlsx'. ",
3772
               "Dataset Path": "skin_dataset/skinL2_dataset"
           }
3773
3774
```

Listing 7: Dataset Descriptions for X-Ray Modality

```
3778 [
3779 {
    "Dataset Name": "cov_19",
```

```
3780
               "Dataset Description": "This is a database of chest X-ray images
3781
                   for COVID-19 positive cases along with Normal and Viral
3782
                  Pneumonia images. It has 3616 COVID-19 positive cases along
3783
                  with 10,192 Normal, 6012 Lung Opacity (Non-COVID lung
                   infection), and 1345 Viral Pneumonia images and corresponding
3784
                    lung masks organized in different sub-folders.",
3785
               "Dataset Path": "xray/cov_19"
3786
          },
3787
3788
               "Dataset Name": "bone_frac",
               "Dataset Description": "This dataset comprises fractured and non-
3789
                   fractured X-ray images covering all anatomical body regions,
3790
                   including lower limb, upper limb, lumbar, hips, knees, etc.
3791
                   The dataset is categorized into two subfolders containing
3792
                   fractured and non-fractured radiographic images.",
               "Dataset Path": "xray/bone_frac"
3793
           },
3794
3795
               "Dataset Name": "chest_tuberculosis_segmentation",
3796
               "Dataset Description": "This dataset consists of 704 chest X-ray
3797
                   images for tuberculosis (TB) detection. The dataset contains
                  both tuberculosis-positive and normal chest X-rays and are
3798
                   accompanied by lung segmentation masks (in separate
3799
                   subfolders) and clinical metadata as csv files.",
3800
               "Dataset Path": "xray/chest_tuberculosis_segmentation"
3801
          },
3802
               "Dataset Name": "xray/17_diseases",
3803
               "Dataset Description": "The dataset consists of a collection of
3804
                   chest X-ray images in .jpg and .dcm formats. Types of
3805
                   diseases in the dataset: Abscess, Ards, Atelectasis,
3806
                  Atherosclerosis of the aorta, Cardiomegaly, Emphysema,
3807
                   Fracture, Hydropneumothorax, Hydrothorax, Pneumonia,
3808
                   Pneumosclerosis, Post inflammatory changes, Post traumatic
                   ribs deformation, Sarcoidosis, Scoliosis, Tuberculosis and
3809
                   Venous congestion arranged in different subfolders.",
3810
               "Dataset Path": "xray/17_diseases"
3811
3812
               "Dataset Name": "spr_age_gender",
3813
               "Dataset Description": "SPR X-Ray Age and Gender Dataset. Used to
3814
                   help predict the age and gender of the patient based on the
3815
                   X-Ray image. Contains .png x-ray images in image subfolder
3816
                   with csv file containing gender and age.",
3817
               "Dataset Path": "xray/spr_age_gender"
3818
           },
3819
               "Dataset Name": "unifesp",
3820
               "Dataset Description": "The UNIFESP X-Ray Body Part
3821
                   Classification Dataset. This is a dataset of 2481 X-rays of
3822
                   20 body parts + others, annotated in a multilabel fashion by
                   radiology residents. Images are in DICOM format and Labels
3823
                   are categorical in csv file: Abdomen = 0, Ankle = 1, Cervical
3824
                   Spine = 2, Chest = 3, Clavicles = 4, Elbow = 5, Feet = 6,
3825
                   Finger = 7, Forearm = 8, Hand = 9, Hip = 10, Knee = 11, Lower
3826
                   Leg = 12, Lumbar Spine = 13, Others = 14, Pelvis = 15,
3827
                   Shoulder = 16, Sinus = 17, Skull = 18, Thigh = 19, Thoracic
                   Spine = 20, Wrist = 21",
3828
               "Dataset Path": "xray/unifesp"
3829
           },
3830
3831
               "Dataset Name": "knee",
3832
               "Dataset Description": "It has 1,650 high-quality digital X-ray
3833
                   images of knee joints with a metadata file.",
               "Dataset Path": "xray/knee"
```

```
3834
           },
3835
3836
               "Dataset Name": "c19_radiograph",
               "Dataset Description": "COVID-19, lung opacity, normal and viral
3837
                  pneumonia chest X-ray (CXR) images are arranged in different
3838
                   sub-folders.",
3839
               "Dataset Path": "xray/c19_radiograph"
3840
           },
3841
3842
               "Dataset Name": "simple_vs_community",
               "Dataset Description": "Bone Fracture X-ray Dataset: Simple vs.
3843
                  Comminuted Fractures organized in different subfolders",
3844
               "Dataset Path": "xray/simple_vs_community"
3845
           },
3846
               "Dataset Name": "nih_bbox",
3847
               "Dataset Description": "This NIH Chest X-ray Dataset is comprised
3848
                    of 112,120 X-ray images with disease labels from 30,805
3849
                   unique patients. It has images in the image folder along with
3850
                    a label.csv with Class labels: 8 classes - Infiltrate,
3851
                  Atelectasis, Pneumonia, Cardiomegaly, Effusion, Pneumothorax,
3852
                   Mass, Nodule.",
               "Dataset Path": "xray/nih_bbox"
3853
           },
3854
3855
               "Dataset Name": "bone break",
3856
               "Dataset Description": "The dataset covers a range of bone
3857
                   fracture classes, such as avulsion fractures, comminuted
                   fractures, fracture-dislocations, greenstick fractures,
3858
                  hairline fractures, impacted fractures, longitudinal
3859
                   fractures, oblique fractures, pathological fractures, and
3860
                   spiral fractures organized in separate subfolders",
3861
               "Dataset Path": "xray/bone_break"
3862
           },
3863
               "Dataset Name": "cov19_normal",
3864
               "Dataset Description": "This dataset comprises a total of 800
3865
                  high-quality chest X-ray images, with 400 images depicting
3866
                   COVID-19 infected patients and the other 400 images
                   representing normal cases (i.e., non-COVID patients) arranged
3867
                    in separate sub-folders.",
3868
               "Dataset Path": "xray/cov19_normal"
3869
           },
3870
3871
               "Dataset Name": "dental",
               "Dataset Description": "Dental radiographs along with labels in
3872
                   csv files",
3873
               "Dataset Path": "xray/dental"
3874
           },
3875
3876
               "Dataset Name": "bone_frac_small",
               "Dataset Description": "This dataset is designed for developing
3877
                  machine learning models for bone fracture classification and
3878
                   localization in tibia and fibula bones. It contains X-ray
3879
                   images in .PNG format along with labels in csv file",
3880
               "Dataset Path": "xray/bone_frac_small"
3881
           },
3882
               "Dataset Name": "knee_osteoporosis",
3883
               "Dataset Description": "This knee XRay dataset contains 3 classes
3884
                   : normal, Osteopenia ,and Osteoporosis arranged in separate
3885
                   subfolders",
3886
               "Dataset Path": "xray/knee_osteoporosis"
3887
           },
```

```
"Dataset Name": "RNSA_pneumonia",
3889
               "Dataset Description": "This dataset is a pre-processed version
3890
                   of the RSNA Pneumonia Detection Challenge dataset in PNG
3891
                   format along with the associated bounding box annotations as
                  mask images. The metadata, including the patient information
3892
                   and bounding box coordinates, has been extracted and saved in
3893
                   CSV format.",
3894
               "Dataset Path": "xray/RNSA_pneumonia"
3895
           },
3896
               "Dataset Name": "8_object_detection",
3897
               "Dataset Description": "Overview: The Chest X-ray 8 Subset
3898
                  dataset is a curated collection of chest radiographs for
3899
                   object detection models on thoracic diseases, with 790 images
3900
                   and 984 annotated bounding boxes in YOLO and Pascal VOC
3901
                   formats for diverse ML frameworks. Classes and Labels
                   contained in associated csv file: 14 thoracic disease classes
3902
                    including Atelectasis, Cardiomegaly, Effusion, Infiltrate,
3903
                   Nodule, Mass, Pneumonia, Pneumothorax.",
3904
               "Dataset Path": "xray/8_object_detection"
3905
           },
3906
               "Dataset Name": "HBFMID",
3907
               "Dataset Description": "Human Bone Fractures Multi-modal Image
3908
                  Dataset (HBFMID) is a collection of 1539 annotated medical
3909
                   images (X-ray and MRI) covering bone fractures in various
3910
                   locations (elbow, finger, forearm, humerus, shoulder, femur,
                   shinbone, knee, hipbone, wrist, spinal cord, and some healthy
3911
                   bones) contained in the Image folder along with associated
3912
                   csv file",
3913
               "Dataset Path": "xray/HBFMID"
3914
           },
3915
               "Dataset Name": "FracAtlas",
3916
               "Dataset Description": "It is a dataset of more than 14,000 X-Ray
3917
                    scans for classification, localization and segmentation of
3918
                  bone fractures. All the scans are available in JPG format
3919
                   along with proper annotations in separate sub-folders",
3920
               "Dataset Path": "xray/FracAtlas"
3921
           },
3922
               "Dataset Name": "pneumonia",
3923
               "Dataset Description": "There are 5,863 X-Ray images (JPEG) and 2
3924
                   categories (Pneumonia/Normal) arranged in separate sub-
3925
                   folders",
               "Dataset Path": "xray/pneumonia"
3926
           },
3927
3928
               "Dataset Name": "pax_ray",
3929
               "Dataset Description": "The PAX-Ray++ Dataset is a high-quality
3930
                   dataset designed to facilitate segmentation tasks for
                   anatomical structures in chest radiographs available in Image
3931
                   subfolder and annotations in mask subfolder.",
3932
               "Dataset Path": "xray/pax_ray"
3933
           },
3934
3935
               "Dataset Name": "lung_segmentation",
               "Dataset Description": "This dataset contains over 500 x-ray
3936
                  scans with clinical labels collected by radiologists
3937
                  available in separate subfolders.",
3938
               "Dataset Path": "xray/lung_segmentation"
3939
           },
3940
           {
               "Dataset Name": "shadow",
3941
```

```
3942
               "Dataset Description": "Normal Chest X-ray images and Bone Shadow
3943
                    images along with csv file.",
3944
               "Dataset Path": "xray/shadow"
3945
           },
3946
               "Dataset Name": "Angiography",
3947
               "Dataset Description": "The ARCADE dataset (Automatic Region-
3948
                   based Coronary Artery Disease Diagnostics using X-ray
3949
                   Angiography) is organized into two task-specific directories
                   ('Task_Syntax_Segmentation' and 'Task_Stenosis_Segmentation'), each containing flattened 'Images/' and 'masks/' subfolders
3950
3951
3952
               "Dataset Path": "xray/Angiography"
3953
           },
3954
               "Dataset Name": "dental_panoramic",
3955
               "Dataset Description": "Dental Disease Panoramic Dataset with
3956
                   segmentations on 31 classes: Classes: 0: Caries, 1: Crown, 2:
3957
                    Filling, 3: Implant, 4: Malaligned, 5: Mandibular Canal, 6:
3958
                   Missing teeth, 7: Periapical lesion, 8: Retained root, 9:
3959
                   Root Canal Treatment, 10: Root Piece, 11: Impacted tooth, 12:
                    Maxillary sinus, 13: Bone Loss, 14: Fracture teeth, 15:
3960
                   Permanent Teeth, 16: Supra Eruption, 17: TAD, 18: Abutment,
3961
                   19: Attrition, 20: Bone defect, 21: Gingival former, 22:
3962
                   Metal band, 23: Orthodontic brackets, 24: Permanent retainer,
3963
                    25: Post-core, 26: Plating, 27: Wire, 28: Cyst, 29: Root
3964
                   resorption, 30: Primary teeth organized as different sub-
3965
                   folders",
               "Dataset Path": "xray/dental_panoramic"
3966
           },
3967
3968
               "Dataset Name": "ALHI",
3969
               "Dataset Description": "All images include a stem and a cup of
3970
                   the hip implant, and the images have to be X-ray images along
                    with csv file containing metadata.",
3971
               "Dataset Path": "xray/ALHI"
3972
3973
           },
3974
               "Dataset Name": "humerus_fractures",
3975
               "Dataset Description": "Deep Learning-Driven Diagnosis of Humerus
3976
                    Fractures from Radiographic Data. Images contain x-ray
3977
                   images of humerus fractures and non-fractures in separate
3978
                   subfolders.",
3979
               "Dataset Path": "xray/humerus_fractures"
3980
           },
3981
               "Dataset Name": "multiclass_knee_osteoporosis",
3982
               "Dataset Description": "The dataset is divided into three primary
3983
                    categories: (1) Normal: Images of knees with no signs of
3984
                   osteoporosis., (2) Osteopenia: Images showing early stages of
3985
                    bone density loss, and (3) Osteoporosis: Images indicating
                   advanced bone density degradation organized as different
3986
                   subfolders",
3987
               "Dataset Path": "xray/multiclass_knee_osteoporosis"
3988
           },
3989
               "Dataset Name": "rsna-breast-cancer-detection",
3990
               "Dataset Description": "Region of Interests extracted from breast
3991
                    X-ray images. There are no labels, just .png images.",
3992
               "Dataset Path": "xray/rsna-breast-cancer-detection"
3993
           },
3994
           {
               "Dataset Name": "RANZCR",
```

```
3996
               "Dataset Description": "For detecting the presence and position
3997
                   of catheters and lines on chest x-rays. The .csv file
3998
                   contains image IDs, binary labels, and patient IDs with
3999
                   columns: Columns: StudyInstanceUID (unique ID for each image)
                   , ETT - Abnormal (endotracheal tube placement abnormal), ETT
4000
                   - Borderline (borderline abnormal), ETT - Normal (normal),
4001
                  NGT - Abnormal (nasogastric tube placement abnormal), NGT -
4002
                   Borderline (borderline abnormal), NGT - Incompletely Imaged (
4003
                   inconclusive due to imaging), NGT - Normal (normal), CVC
4004
                   Abnormal (central venous catheter placement abnormal), CVC -
                   Borderline (borderline abnormal), CVC - Normal (normal), Swan
4005
                   Ganz Catheter Present, PatientID (unique ID for each patient
4006
                   ) . ",
4007
               "Dataset Path": "xray/RANZCR"
4008
           },
4009
               "Dataset Name": "FractureFusion",
4010
               "Dataset Description": "From avulsion fractures to spiral
4011
                   fractures, this dataset is a rich repository of diverse cases
4012
                   , including comminuted fractures, fracture-dislocations,
4013
                   greenstick fractures, hairline fractures, impacted fractures,
4014
                    longitudinal fractures, oblique fractures, pathological
                   fractures arranged as different subfolders",
4015
               "Dataset Path": "xray/FractureFusion"
4016
           },
4017
4018
               "Dataset Name": "HeelBone",
               "Dataset Description": "Heel Bone X-Ray Dataset consists of 3,956
4019
                    X-ray images of the foot, primarily focused on detecting and
4020
                    classifying heel bone diseases with annotations arranged in
4021
                   label.csv",
4022
               "Dataset Path": "xray/HeelBone"
4023
           }
4024
```

Listing 8: Dataset Descriptions for Histopathology Modality

```
4027
4028
4029
               "Dataset Name": "breast_histo",
4030
               "Dataset Description": "Breast Histopathology Images with
                   Invasive Ductal Carcinoma (IDC). There's no labels for this
4031
                   dataset, only images.",
4032
               "Dataset Path": "histopathology/breast_histo"
4033
           },
4034
4035
               "Dataset Name": "BreaKHis_400X",
4036
               "Dataset
                   Description": "Breast cancer images on histopathology slides.
4037
                    The BreaKHis database contains microscopic biopsy images
4038
                   benign and malignant breast tumors in separate subfolders.",
4039
               "Dataset Path": "histopathology/BreaKHis_400X"
           },
4040
4041
               "Dataset Name": "lung_and_colon",
4042
               "Dataset Description": "Lung and Colon Cancer Histopathological
4043
                   Images: 25000 images of 5 classes: Lung benign tissue, Lung
                   adenocarcinoma, Lung squamous cell carcinoma, Colon
4044
                   adenocarcinoma, Colon benign tissue in separate subfolders.",
4045
               "Dataset Path": "histopathology/lung_and_colon"
4046
4047
4048
               "Dataset Name": "gastric_cancer",
               "Dataset Description": "Gastric Cancer Histopathology Tissue
4049
                   Image Dataset focuses on the tumor microenvironment (TME) and
```

```
4050
                    includes images categorized into eight distinct tissue types
4051
                   : ADI: Adipose (fat tissue), BACK: Background (non-tissue
4052
                   areas), DEB: Debris (cellular waste), LYM: Lymphocytes (
4053
                   immune cells), MUC: Mucus (protective secretion), MUS: Smooth
                   Muscle (muscle tissue), NORM: Normal Colon Mucosa (healthy
4054
                   tissue for reference), STR: Cancer-associated Stroma (
4055
                   connective tissue around the tumor), TUM: Tumor (cancerous
4056
                   tissue) - all arranged in different subfolders. ",
4057
               "Dataset Path": "histopathology/gastric_cancer"
4058
           },
4059
               "Dataset Name": "gastro_cancer_msi_vs_mss",
4060
               "Dataset Description": "The dataset contains histological images
4061
                   for MSI vs MSS classification in gastrointestinal cancer
4062
                   arranged in different sub-folders.",
4063
               "Dataset Path": "histopathology/gastro_cancer_msi_vs_mss"
           },
4064
4065
               "Dataset Name": "breast_cancer_segmentation",
4066
               "Dataset Description": "Breast Cancer Cell Segmentation with
4067
                   corresponding images and masks in separate subfolders.",
               "Dataset Path": "histopathology/breast_cancer_segmentation"
4068
           },
4069
4070
               "Dataset Name": "ovarian_cancer",
4071
               "Dataset Description": "Ovarian Cancer & Subtypes Dataset
4072
                   Histopathology: This dataset includes histopathology images
4073
                   of 4 subtypes of Ovarian cancer and also non cancerous
                  histopathological images organized in separate subfolders",
4074
               "Dataset Path": "histopathology/ovarian_cancer"
4075
           },
4076
4077
               "Dataset Name": "breast_cancer_histo",
               "Dataset Description": "breast cancer histopathology. JPG images
4078
                   with classifications benign or malignant organized as
4079
                   separate subfolders",
4080
               "Dataset Path": "histopathology/breast_cancer_histo"
4081
4082
               "Dataset Name": "BreCaHAD",
4083
               "Dataset
4084
                   Description": "a dataset for breast cancer histopathological
4085
                   annotation and diagnosis with images belonging to six classes
4086
                   , namely mitosis, apoptosis, tumor nuclei, non-tumor nuclei,
                   tubule, and non-tubule arranged in separate subfolders",
4087
               "Dataset Path": "histopathology/BreCaHAD"
4088
           },
4089
4090
               "Dataset Name": "melanoma",
4091
               "Dataset Description": "This dataset is a melanoma specific
                   dataset with nuclei and tissue annotations along with
4092
                   original images in separate subfolders.",
4093
               "Dataset Path": "histopathology/melanoma"
4094
           },
4095
4096
               "Dataset Name": "choledoch",
               "Dataset Description": "This is a database for both microscopy
4097
                   hyperspectral and color images of cholangiocarcinoma,
4098
                   including 880 scenes among which 689 scenes are samples with
4099
                   part of cancer areas (L), 49 scenes full of cancer areas (N),
4100
                   and 142 scenes without cancer areas (P) organized as
4101
                   separate subfolders",
               "Dataset Path": "histopathology/choledoch"
4102
           },
4103
           {
```

```
4104
               "Dataset Name": "histopath-sn",
4105
               "Dataset Description": "This is a Kaggle dataset, with the task
4106
                   to classify patches: Bronchus and lung samples in image
4107
                   folder along with labels in separate csv file.",
               "Dataset Path": "histopathology/histopath-sn"
4108
           },
4109
4110
               "Dataset Name": "ULMS",
4111
               "Dataset Description": "Uterine leiomyosarcoma (ULMS) dataset
4112
                   comprises mitosis count, necrosis, and nuclear atypia with
                   labels in separate csv file",
4113
               "Dataset Path": "histopathology/ULMS"
4114
4115
           },
4116
               "Dataset Name": "MonuSeg",
4117
               "Dataset Description": "The dataset comprises nuclei from seven
4118
                   organs with associated annotations in csv file.",
4119
               "Dataset Path": "histopathology/MonuSeg"
4120
           },
4121
               "Dataset Name": "kmc_kidney",
4122
               "Dataset Description": "The introduced KMC kidney histopathology
4123
                  dataset includes non-cancerous (Grade-0) and cancerous (Grade
4124
                   -1 to Grade-4) images of the Renal Clear Cell Carcinoma
4125
                  organized as separate subfolders",
4126
               "Dataset Path": "histopathology/kmc_kidney"
4127
           },
4128
               "Dataset Name": "histo-img-text",
4129
               "Dataset Description": "This is a kaggle dataset with
4130
                  histopathology image-text pairs",
4131
               "Dataset Path": "histopathology/histo-img-text"
4132
           },
4133
               "Dataset Name": "cellnet",
4134
               "Dataset Description": "CellNet is a meticulously curated dataset
4135
                    featuring over 120,000 high-quality medical images
4136
                   representing over 20 organ/cancer classes arranged as
                   different subfolders. ",
4137
               "Dataset Path": "histopathology/cellnet"
4138
4139
4140
               "Dataset Name": "PanNuke",
4141
               "Dataset Description": "Nuclei instance segmentation and
                   classification dataset with exhaustive nuclei labels across
4142
                   19 different tissue types. In total the dataset contains
4143
                   205,343 labeled nuclei, each with an instance segmentation
4144
                  mask in separate datasets.",
4145
               "Dataset Path": "histopathology/PanNuke"
4146
           },
4147
               "Dataset Name": "NPC-88k-Public",
4148
               "Dataset Description": "88k histopathology patches of normal,
4149
                   lymphoid hyperplasia (LHP), nasopharyngeal inflammation (NPI)
4150
                   , and nasopharyngeal carcinoma (NPC) organized in separate
4151
                   subfolders.",
               "Dataset Path": "histopathology/NPC-88k-Public"
4152
4153
4154
               "Dataset Name": "EBHI",
4155
               "Dataset
4156
                  Description": "The dataset encompasses various categories,
                   including normal (76 images and 76 ground truth images),
4157
                  polyp (474 images and 474 ground truth images), low-grade
```

```
intraepithelial neoplasia (639 images and 639 ground truth
  images), high-grade intraepithelial neoplasia (186 images and
  186 ground truth images), serrated adenoma (58 images and 58
  ground truth images), and adenocarcinoma (795 images and 795
  ground truth images) arranged in different subfolders",
  "Dataset Path": "histopathology/EBHI"
}
```

# C.3 DETECTING AND ADDRESSING DATA QUALITY ISSUES FOR DATA PRE-PROCESSING AGENT

One of the primary steps in data pre-processing involves identifying data quality issues and removing samples that negatively impact the overall data quality. In this work, we address three key data quality issues *viz.* **off-topic samples**, **near duplicates**, and **label errors** Gröger et al. (2025; 2024; 2023) each of which can significantly compromise the reliability of machine learning models, particularly in sensitive domains like medical imaging.

- Off-topic samples refer to irrelevant inputs mistakenly included in the dataset (e.g., from unrelated modalities or corrupted acquisitions). These introduce noise, distort evaluation metrics, and hinder model convergence.
- **Near duplicates** are different views of the same object, including exact copies. Their presence artificially reduces the diversity of the training set, introduces redundancy, and may lead to data leakage between training and evaluation sets.
- Label errors are incorrectly annotated examples that can misguide both model training and
  evaluation, leading to degraded performance and spurious generalization.

The dataset is formalized as  $\mathcal{X} = \{(x_i, l_i) \mid i \in \mathcal{I}\}$ , where each  $x_i$  is a sample,  $l_i$  is its label among L classes, and  $\mathcal{I} = \{1, \dots, N\}$  the index set. For each issue type, a scoring function  $s(\cdot)$  is defined that maps individual samples or sample pairs to a score in [0, 1], where lower values indicate higher likelihood of an issue. Ranking the samples by these scores yields a prioritized list for inspection or automated filtering based on a pre-defined threshold.

#### REPRESENTATION LEARNING

A deep feature extractor  $f(\cdot;\theta)$  was trained using self-supervised learning (SSL) methods (SimCLR or DINO), both of which were implemented with a Vision Transformer (ViT) backbone. Each sample  $x_i$  was embedded into a latent space as  $e_i = f(x_i;\theta) \in \mathbb{R}^D$ , where D denotes the feature dimension. To ensure consistent geometry across methods,  $\ell_2$ -normalization was applied so that all embeddings lie on a unit hypersphere.

Cosine similarity was adopted to define the distance metric:

$$\mathrm{sim}(e_i,e_j) = \frac{e_i^\top e_j}{\|e_i\|_2 \|e_j\|_2}, \quad \mathrm{dist}(e_i,e_j) = \frac{1-\mathrm{sim}(e_i,e_j)}{2}.$$

## **ISSUE-SPECIFIC DETECTION STRATEGIES**

**Off-topic Detection.** Off-topic samples were identified using agglomerative clustering with single linkage in the representation space. The merging behavior of clusters was analyzed, and samples that were merged at higher distances or at later stages with larger clusters were considered more likely to be anomalous. A scoring function  $s_{\text{OT}}(e_i)$  was constructed based on merge depth and inter-cluster distance dynamics.

**Near Duplicate Detection.** Candidate near-duplicate pairs were detected by evaluating pairwise distances between all sample embeddings. A simple ranking function was applied:

$$s_{ND}(e_i, e_j) = \operatorname{dist}(e_i, e_j),$$

where smaller distances were interpreted as a higher likelihood of duplication.

**Label Error Detection.** Label errors were inferred based on a ratio between intra-class and interclass distances. For each sample  $e_i$ , the following definitions were used:

$$\begin{split} m_{=}(e_i) &= \min_{j \in \mathcal{I}, \ l_j = l_i} \mathrm{dist}(e_i, e_j), \quad m_{\neq}(e_i) = \min_{j \in \mathcal{I}, \ l_j \neq l_i} \mathrm{dist}(e_i, e_j), \\ s_{\mathrm{LE}}(e_i) &= \frac{m_{\neq}^2(e_i)}{m_{=}^2(e_i) + m_{\neq}^2(e_i)}. \end{split}$$

Lower scores were interpreted as indicating a higher likelihood of mislabeling, particularly when the nearest neighbor belonged to a different class.

In all three cases, the local structure of the embedding space was leveraged by the cleaning function used in Tool 9 of the Listing 1. Cluster distances were evaluated using only the nearest neighbors for off-topic detection, proximity among sample pairs was assessed for duplicate identification, and comparative distances to same- and different-class neighbors were exploited to detect label errors.

## C.4 COLLECTION OF FEDERATED LEARNING ALGORITHMS

Federated Learning (FL) has evolved significantly beyond its initial formulation of model averaging, with numerous algorithmic innovations developed to address practical challenges such as data heterogeneity, personalization, privacy preservation, and limited client resources McMahan et al. (2017); Tan & Wang; Tan et al. (2023). In this work, we utilize a set of **40** key federated learning (FL) algorithms, covering core, personalized, generalizable, and adaptive methods, as summarized in Tables 2-4. The algorithm description required by server-based federated training agents for FL algorithm selection is provided in Listing 9.

The selected algorithms reflect the diversity and progression of research in FL across three main axes:

## 1. Foundational and General-Purpose Methods:

We begin with core algorithms such as *FedAvgM*, and *FedProx*, which establish the baseline principles of client-server aggregation and account for statistical and system heterogeneity. These methods are essential for benchmarking and provide the backbone upon which many subsequent algorithms are built.

## 2. Personalization-Oriented Methods:

Recognizing the need to adapt to non-IID data across clients, we include algorithms like *FedRep*, *FedPer*, *Ditto*, *pFedHN*, and *Per-FedAvg*. These approaches personalize part of the model (e.g., classifier heads or entire layers), use meta-learning, or leverage client-specific adaptation strategies. Methods such as *pFedMe* and *FedEM* extend this personalization through bi-level optimization and mixture modeling, respectively.

## 3. Robustness, Adaptivity, and Generalization:

To tackle challenges of out-of-distribution generalization and domain shifts, we incorporate algorithms like *FedIIR*, *FedSR*, and *ADCOL*, which emphasize invariant representation learning and adversarial feature alignment. Techniques such as *FedDyn*, *FedFomo*, and *FedRoD* introduce dynamic regularization and adaptive weighting to stabilize optimization in heterogeneous environments. Moreover, algorithms like *FedBN* and *FedAP* address domain-specific normalization challenges, particularly in healthcare contexts.

## 4. Emerging and Specialized Directions:

The inclusion of recent methods such as *Floco*, *FedAS*, and *PeFLL* highlights advancements in adaptive aggregation, inter-client relationship modeling, and meta-learned personalization. Additionally, *MOON*, *FedGen*, and *CCVR* represent innovative uses of contrastive learning, data-free distillation, and virtual representation calibration.

The rationale for selecting this curated list is threefold:

- **Comprehensiveness:** The algorithms span from classic to state-of-the-art methods, ensuring broad coverage of the field.
- Modular Design Potential: These algorithms are suitable for integration into modular federated learning pipelines, facilitating agent-based automation and tool invocation.

4267

4268

4269 4270

4271

4272 4273 • Relevance to Real-World Scenarios: Many chosen methods address constraints encountered in practical deployments, including label imbalance, resource limitations, domain adaptation, and personalization needs.

This comprehensive collection enables systematic benchmarking, comparative evaluation, and modular composition in our federated learning framework *FedAgentBench*. Each method contributes unique strengths and trade-offs, making them valuable candidates for real-world and research applications.

Listing 9: Federated Learning Algorithm Descriptions for Server-based algorithm selector agents

```
4274
4275
4276
               "algorithm": "FedAvg",
4277
               ""description"": "The foundational algorithm in federated
4278
                   learning, where clients perform multiple steps of local
4279
                   stochastic gradient descent (SGD) and periodically average
                   their models on a central server. It is simple and
4280
                   communication-efficient but struggles with non-IID data
4281
                   distributions."
4282
           },
4283
4284
               "algorithm": "FedAvqM",
               ""description"": "An extension of FedAvg that integrates server-
4285
                   side momentum during model aggregation. This is a classical
4286
                   federated learning approach that stabilizes training and
4287
                   improves convergence in the presence of data heterogeneity
4288
                   across clients."
4289
           },
4290
               "algorithm": "FedProx",
               ""description"": "Classical federated learning algorithm that
4292
                  enhances FedAvg by adding a proximal term to the local
4293
                   objective functions, discouraging local updates from drifting
4294
                   too far from the global model. This addresses system and
                   statistical heterogeneity among clients."
4295
           },
4296
4297
               "algorithm": "SCAFFOLD",
4298
               ""description"": "Classical federated learning algorithm that
                   incorporates control variates to correct client-drift caused
4299
                  by non-IID data. Each client maintains local control
4300
                  variables to align updates with the global objective,
4301
                   improving convergence stability."
4302
           },
4303
               "algorithm": "MOON",
4304
               ""description"": "Traditional Federated learning algorithm that
4305
                   implements model-level contrastive learning by aligning
4306
                  current local models with the global model while contrasting
4307
                  them with past local models. This enhances representation
                   learning under non-IID settings."
4309
           },
4310
               "algorithm": "FedDyn",
4311
               ""description"": "Regularization-based federated learning
4312
                   approach that introduces a dynamic regularization term into
4313
                   local objectives that evolves over time to better match the
4314
                  global objective. This mechanism helps mitigate divergence
                   and stabilizes training in heterogeneous environments."
4315
           },
4316
4317
               "algorithm": "FedLC",
4318
               ""description"": "Classical federated learning algorithm that
4319
                   applies logits calibration techniques during local training
                   to address label distribution skew. This helps balance
```

Table 4: Overview of Federated Learning Algorithms (Part 1)

4322		Table 4: Ove	erview of Federated Lear	ning Algorithms (Part 1)	
4323	Method	Source	Key Idea	Strengths	Limitations
4324	FedAvg	McMahan et	Clients perform lo-	Simple and	Degrades with non-
4325	McMa-	al., 2016	cal SGD and periodi-	communication-	IID data due to client
4326	han et al.		cally average with the	efficient.	drift.
4327	(2017)		server.		
4328	FedAvgM	Hsu et al.,	Adds server-side mo-	Improves conver-	Requires careful mo-
4329	Hsu et al.	2019	mentum to FedAvg.	gence on non-IID	mentum tuning.
4330	(2019)	T 1	TT 11' 1	data.	D ' 11'
4331	FedMD Li	Li et al.,	Uses public dataset	Supports diverse ar-	Requires public
4332	& Wang (2019)	NeurIPS 2019	for knowledge distil- lation across hetero-	chitectures.	dataset.
4333	(2019)		geneous models.		
4334	FedPer	Arivazhagan	Uses client-specific	Balances global and	Designing layer split
4335	Arivazha-	et al., arXiv	layers with shared	local learning.	is non-trivial.
4336	gan et al.	2019	global layers.	iour iouring.	15 11011 1111 1111
4337	(2019)		8		
4338	LG-	Liang et al.,	Aggregates global	Preserves local per-	Complex model syn-
4339	FedAvg	NeurIPS 2019	layers, retains local	sonalization.	chronization.
4340	Liang et al.	Workshop	ones.		
4341	(2020)				
4342	CFL Sat-	Sattler et al.,	Clusters clients and	Addresses data het-	Doesn't scale well
4343	tler et al.	arXiv 2019	trains separate mod-	erogeneity.	with many clusters.
4344	(2019)	T: 1 2020	els.	TT 11	N/ 1 1
4345	FedProx	Li et al., 2020	Adds proximal term	Handles statisti-	May slow down con-
4346	Li et al. (2020b)		to local loss.	cal/system heterogeneity.	vergence.
4347	SCAFFOLD	Karimireddy	Uses control variates	Better convergence	Extra storage and
4348	Karim-	et al., 2020	to correct drift.	on non-IID data.	computation.
4349	ireddy	ct al., 2020	to confect difft.	on non no data.	computation.
4350	et al.				
4351	(2020)				
4352	APFL	Deng et al.,	Adaptive mixing of	Combines generaliza-	Requires careful mix-
4353	Deng et al.	arXiv 2020	global and local mod-	tion and personaliza-	ing parameter tuning.
4354	(2020)		els.	tion.	
4355	Per-	Fallah et al.,	Combines FL with	Enables fast personal-	Needs second-order
4356	FedAvg	NeurIPS 2020	MAML.	ization.	gradients.
4357	Fallah				
4358	et al. (2020)				
4359	pFedMe	Dinh et al.,	Uses Moreau en-	Fast convergence and	Requires tuning of
4360	Dinh et al.	NeurIPS 2020	velopes for bi-level	good personalization.	regularization.
4361	(2022)	1400111 5 2020	optimization.	5000 personanzación.	105ului 12ulioii.
4362	MOON	Li et al.,	Aligns local and	Strong representation	Needs previous
4363	Li et al.	CVPR 2021	global models via	learning.	model storage.
4364	(2021a)		contrastive loss.		<i>5</i> ·
4365	FedDyn	Acar et al.,	Dynamic regulariza-	Mitigates client drift.	More complex opti-
4366	Acar et al.	ICLR 2021	tion to align objec-	-	mization.
4367	(2021)		tives.		
4368	FedGen	Zhu et al.,	Uses synthetic data	Enables data-free	Depends on generator
4369	Zhu et al.	ICML 2021	for knowledge distil-	generalization.	quality.
4370	(2021)	D 11' : 1	lation.	E // 11	TT
4371	FedOpt	Reddi et al.,	Uses adaptive opti-	Fast/stable conver-	Hyperparameter tun-
4372	Reddi et al.	ICLR 2021	mizers (Adam/Yogi)	gence.	ing required.
4373	(2021)		in FL.		

Table 5: Overview of Federated Learning Algorithms (Part 2)

Method	Source	Key Idea	Strengths	Limitations
CCVR Luo et al. (2021)	Wang et al., NeurIPS 2021	Virtual representations for calibration.	No real data sharing needed.	Relies on distribution approximations.
FedEM Marfoq et al. (2022)	Marfoq et al., NeurIPS 2021	Mixture model for multi-task personalization.	Captures cross-client distributions.	Assumes shared latent structure.
Ditto Li et al. (2021c)	Li et al., ICML 2021	Maintains global and personalized models.	Robust and fair personalization.	Needs dual model training.
FedRep Collins et al. (2023)	Collins et al., ICML 2021	Shared encoder with local classifiers.	Combines global and local strengths.	Coordination needed for shared layer.
pFedHN Shamsian et al. (2021)	Shamsian et al., ICML 2021	Hypernetworks generate personalized models.	Communication efficient.	Complex hypernet- work training.
FedFomo Zhang et al. (2021)	Zhang et al., ICLR 2021	Aggregates based on client similarity.	Personalization without raw data.	Similarity computation overhead.
FedBN Li et al. (2021d)	Li et al., ICLR 2021	Local BN layers for domain adaptation.	Improves performance on non-IID data.	No global BN normalization.
FedLC Zhang et al. (2022)	Zhang et al., ICML 2022	Logits calibration to handle label skew.	Effective on imbalanced datasets.	Needs label distribution estimation.
MetaFed Chen et al. (2023b)	IJCAI 2022	Cyclic knowledge distillation across federations.	Enhances collaboration.	Federation coordination required.
FedRoD Chen & Chao (2022)	ICLR 2022	Adaptive aggregation for balancing general/personal models.	Personalized and generalizable.	May fail under high heterogeneity.
FedProto Tan et al. (2022)	AAAI 2022	Prototype-based feature alignment.	Preserves global semantics.	Quality depends on prototypes.
pFedLA Ma et al. (2022)	Ma et al., CVPR 2022	Layer-wise model aggregation.	Fine-grained personalization.	Management complexity.
FedBABU Oh et al. (2022)	Oh et al., ICLR 2022	Aggregates body and keeps local heads.	Improves representation learning.	Less consistent pre- dictions.
FedAP Lu et al. (2022)	Chen et al., IEEE 2022	Adaptive BN for healthcare FL.	Handles domain shift.	Sensitive to BN statistics.

Table 6: Overview of Federated Learning Algorithms (Part 3)

Method	Source	Key Idea	Strengths	Limitations
FedSR	NeurIPS 2022	Domain generaliza-	Lightweight and sim-	May fail in extreme
Nguyen		tion via representa-	ple.	domain shift.
et al.		tion regularization.		
(2022a)				
FedALA	AAAI 2023	Adaptive local aggre-	Relevance-aware up-	Unstable weight esti-
Zhang		gation weights.	dates.	mation.
et al.				
(2023)				
FedFed	Yang et al.,	Distills critical fea-	Improves generaliza-	Needs good feature
Yang et al.	NeurIPS 2023	tures.	tion.	selection.
(2023)	Ci.		7.1	
Elastic Ag-	Chen et al.,	Sensitivity-based up-	Balances adapta-	Adds computation.
gregation	CVPR 2023	date weighting.	tion/stability.	
Chen et al.				
(2023a)	ICM 12022		XX 11 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	
ADCOL	ICML 2023	Adversarial align-	Handles domain shift.	Adversarial training
Li et al.		ment of features.		instability.
(2023b)	ICMI 2022	T	0. 1.	NT 1
FedIIR	ICML 2023	Learns invariant re-	Strong generalization.	Needs assumptions on invariance.
Guo et al. (2023)		lationships for OOD generalization.	tion.	on invariance.
pFedSim	Tan et al.,	Similarity-based ag-	Emphlos managenelias	Hard to measure sim-
Tan et al.	arXiv 2023	gregation.	Enables personalization.	ilarity.
(2023)	al XIV 2025	gregation.	tion.	marity.
PeFLL	ICLR 2024	Meta-learns to per-	Fast client adaptation.	High computation
Scott et al.	ICLK 2024	sonalize clients.	rasi cheni adaptation.	cost.
(2025)		sonanze chents.		cost.
FLUTE	ICML 2024	Efficient rep learning	Resource efficient.	May sacrifice expres-
Liu et al.	1CIVIL 2021	under underparame-	Resource emerent.	sivity.
(2024a)		terization.		51,10,1
FedAS	CVPR 2024	Reduces global-local	More consistent up-	More complex train-
Yang et al.		inconsistency.	dates.	ing.
(2024)		,		<b>O</b> .
Floco	NeurIPS 2024	Uses connected	Leverages inter-client	Needs client connec-
Grinwald		modes to model	structure.	tivity info.
et al.		clients.		*
(2025)				

Table 7: Categorization of FL Algorithms

Category	Algorithms						
(i) Classical FL algorithms	FedAvg, FedAvgM, FedProx, SCAFFOLD, MOON,						
	FedLC						
(ii) Personalized FL algo-	Per-FedAvg, pFedMe, FedRep, FedPer, FedBN,						
rithms	pFedLA, pFedHN, FedFomo, LG-FedAvg, APFL,						
	FedEM, pFedSim, FedBABU, CCVR						
(iii) Regularization-based ap-	Ditto, FedDyn, FedRoD, FedAS, SCAFFOLD, pFedMe						
proaches							
(iv) Knowledge Distillation-	FedGen, FedMD, FedFed, MetaFed						
based methods							
(v) Domain generalization	FedSR, FedIIR, ADCOL, FedProto, FedAP						
techniques							
(vi) Optimization and	FedOpt, FedAvgM, FedALA, Elastic						
scheduling variants	Aggregation, FLUTE, PeFLL, CFL						

```
4482
                   prediction confidence and improve accuracy on imbalanced or
4483
                   non-IID datasets."
4484
           },
4485
               "algorithm": "FedGen",
4486
               ""description"": "Personalized Federated Learning leveraging
4487
                   knowledge distillation that uses a server-side generative
4488
                   model to synthesize data representations for knowledge
4489
                   distillation, enabling model personalization without
4490
                   requiring access to client data. This preserves privacy while
                    supporting generalization."
4491
           },
4492
4493
               "algorithm": "CCVR",
4494
               ""description"": "Personalized Federated Learning that uses
4495
                   virtual representations drawn from approximated data
                   distributions to calibrate classifiers. This approach
4496
                   improves generalization in non-IID scenarios without needing
4497
                   to exchange actual data between clients."
4498
           },
4499
               "algorithm": "FedOpt",
4500
               ""description"": "Federated adaptive optimization scheme that
4501
                   extends FedAvg by integrating adaptive gradient methods like
4502
                   FedAdam, FedYogi, and FedAdagrad, which dynamically adjust
4503
                   learning rates and enhance convergence performance in diverse
4504
                    federated settings."
4505
4506
               "algorithm": "Elastic Aggregation",
4507
               ""description"": "Classical federated optimization scheme that
4508
                   introduces elasticity in the aggregation process by assigning
4509
                   dynamic weights to client updates based on the sensitivity
4510
                   of model parameters. This balances stability and adaptability
                   , improving performance on heterogeneous datasets."
4511
           },
4512
4513
               "algorithm": "FedFed",
4514
               ""description"": "Federated learning algorithms that allows
                   partial feature sharing between clients and server and
4515
                   mitigates data heterogeneity by distinguishing between
4516
                   performance-sensitive and performance-robust features and
4517
                   selectively distilling the former. This allows clients to
4518
                   retain useful features while benefiting from cross-client
4519
                   generalization."
4520
           },
4521
               "algorithm": "pFedSim",
4522
               ""description"": "Personalized Federated Learning Algorithm that
4523
                   enhances personalization by aggregating client models based
4524
                   on the similarity of their data distributions. Clients with
                   more similar data contribute more significantly to each other
4525
                   's updates, enabling customized learning without explicit
4526
                   data sharing."
4527
           },
4528
4529
               "algorithm": "FedMD",
               ""description"": "Personalized Federated Learning Algorithm that
4530
                   supports clients with heterogeneous architectures by
4531
                   performing knowledge distillation using a shared public
4532
                   dataset. Clients align on output predictions rather than
4533
                  model parameters, enabling collaborative training without
4534
                   requiring architectural uniformity."
4535
           },
```

```
4536
               "algorithm": "APFL",
4537
               ""description"": "Personalized Federated Learning Algorithm that
4538
                   implements an adaptive mixing strategy where each client
4539
                  maintains both a local and a global model. The final model
                   output is a weighted combination, and the mixing coefficient
4540
                   is learned during training to achieve optimal personalization
4541
4542
           },
4543
4544
               "algorithm": "LG-FedAvg",
               ""description"": "Personalized Federated Learning Algorithm that
4545
                  decomposes models into local and global components, where
                   only the global part is aggregated across clients. This
4547
                  preserves local knowledge while benefiting from global trends
4548
                   , supporting personalized learning in non-IID settings."
4549
           },
4550
               "algorithm": "FedBN",
4551
               ""description"": "Personalized Federated Learning Algorithm that
4552
                   keeps batch normalization layers local to each client while
4553
                   sharing the rest of the model globally. This enables
                   adaptation to client-specific feature distributions and
4554
                   enhances performance under feature heterogeneity."
4555
           },
4556
4557
               "algorithm": "FedPer",
4558
               ""description"": "Personalized Federated Learning Algorithm that
                   introduces personalization by partitioning the model into a
4559
                   globally shared base and a locally updated head. This
4560
                   structure allows clients to fine-tune their models based on
4561
                   local data while retaining shared representations."
4562
           },
4563
               "algorithm": "FedRep",
4564
               ""description"": "Personalized Federated Learning Algorithm that
4565
                   learns a common feature extractor shared across clients and
4566
                   allows each client to train its own classifier head. This
4567
                   separation supports personalization without requiring full
4568
                  model updates across the federation."
4569
           },
4570
               "algorithm": "Per-FedAvg",
4571
               ""description"": "Personalized Federated Learning Algorithm that
4572
                   combines meta-learning (specifically MAML) with federated
4573
                   learning to learn a global initialization that can be rapidly
                   personalized to each clients local data, enabling quick
4574
                   adaptation with limited samples."
4575
           },
4576
4577
               "algorithm": "pFedMe",
4578
               ""description"": "Personalized Federated Learning Algorithm that
                   formulates personalized federated learning as a bi-level
4579
                   optimization problem using Moreau envelopes, which allows
4580
                   decoupling global and local updates. This improves
4581
                   convergence and supports better personalization."
4582
           },
4583
               "algorithm": "FedEM",
4584
               ""description"": "Personalized Federated Learning Algorithm that
4585
                   performs multi-task learning. It treats each clients model
4586
                   as part of a mixture of distributions and trains them via
4587
                   the Expectation-Maximization algorithm. This enables multi-
4588
                   task personalization by modeling shared and unique components
4589
                    across clients."
           },
```

```
4590
4591
               "algorithm": "Ditto",
4592
               ""description"": "Personalized Federated Learning Algorithm that
4593
                   simultaneously trains a global model for generalization and a
                    personalized model for each client, ensuring fairness and
4594
                   robustness through dual-objective optimization."
4595
           },
4596
4597
               "algorithm": "pFedHN",
               ""description"": "Personalized Federated Learning Algorithm that
4598
                  utilizes a central hypernetwork that generates personalized
4599
                  model weights for clients, enabling parameter sharing while
4600
                   allowing client-specific adaptations."
4601
           },
4602
               "algorithm": "pFedLA",
4603
               ""description"": "Personalized Federated Learning Algorithm that
4604
                   performs layer-wise model aggregation, assigning personalized
4605
                    importance to each layer across clients to improve fine-
4606
                   grained adaptation in non-IID environments.'
4607
           },
4608
               "algorithm": "CFL",
4609
               ""description"": "Federated Learning algorithm that clusters
4610
                   clients based on model or data similarity and trains distinct
4611
                   models per cluster to effectively manage heterogeneity
4612
                   across groups."
4613
           },
4614
               "algorithm": "FedFomo",
4615
               ""description"": "Personalized Federated Learning Algorithm that
4616
                  maintains a personalized model by aggregating updates from
4617
                   peer clients weighted by similarity scores, using a first-
4618
                  order gradient approximation to ensure communication
                  efficiency."
4619
           },
4620
4621
               "algorithm": "FedBabu",
4622
               ""description"": "Personalized Federated Learning Algorithm that
                   improves personalized learning by aggregating only the shared
4623
                   body (feature extractor) of the model while keeping client-
4624
                   specific heads independent."
4625
           },
4626
4627
               "algorithm": "FedAP",
               ""description"": "Personalized Federated Learning Algorithm that
4628
                   employs adaptive batch normalization to tailor models to
4629
                  healthcare clients, effectively handling distribution shifts
4630
                   across medical institutions."
4631
           },
4632
               "algorithm": "MetaFed",
4633
               ""description"": "Personalized Federated Learning Algorithm that
4634
                   applies a cyclic knowledge distillation framework across
4635
                   federated groups, improving model generalizability without
4636
                   raw data exchange and without necessity of a server."
4637
           },
4638
               "algorithm": "FedRoD",
4639
               ""description"": "Regularization-based Federated Learning
4640
                   approach that balances the benefits of generalization and
4641
                   personalization by adaptively mixing global and local model
4642
                   components using regularized dual objectives."
4643
           },
```

```
4644
               "algorithm": "FedProto",
4645
               ""description"": "Personalized and generalizable Federated
4646
                   learning algorithm that aligns client features through the
4647
                   use of global class prototypes, promoting semantic
                   consistency while preserving personalization."
4648
           },
4649
4650
               "algorithm": "FedALA",
4651
               ""description"": "Personalized Federated learning algorithm that
4652
                   aggregates local models adaptively by learning relevance-
                   based weights for each client, enabling better
4653
                   personalization through dynamic influence modeling."
4654
           },
4655
4656
               "algorithm": "PeFLL",
               ""description"": "Personalized Federated learning algorithm that
4657
                   incorporates meta-learning to personalize model updates for
4658
                   each client by learning an optimal initialization that
4659
                   generalizes quickly to local tasks."
4660
           },
4661
               "algorithm": "FLUTE",
4662
               ""description"": "Personalized Federated learning algorithm that
4663
                   addresses model underparameterization in resource-constrained
4664
                    environments by learning efficient global and local decoders
4665
                    for distributed representation learning."
4666
           },
4667
               "algorithm": "FedAS",
4668
               ""description"": "Personalized Federated learning algorithm using
4669
                    regularization-based approach that aligns global and local
4670
                   model updates using adaptive strategies to reduce
4671
                   inconsistency and improve convergence in personalized
                   federated learning."
4672
           },
4673
4674
               "algorithm": "Floco",
4675
               ""description"": "Personalized Federated learning algorithm that
4676
                   models client relationships using a graph of local modes and
                   clusters them for collaborative training, leveraging shared
4677
                   structure without central data."
4678
           },
4679
4680
               "algorithm": "FedSR",
4681
               ""description"": "Federated domain generalization-based technique
                    that applies simple regularization across domain
4682
                   representations to improve out-of-distribution generalization
4683
                    in federated settings."
4684
           },
4685
               "algorithm": "ADCOL",
4686
               ""description"": "Federated domain generalization-based technique
4687
                    that uses adversarial learning to align feature spaces
4688
                   across clients, enabling domain generalization under non-IID
4689
                   conditions."
4690
           },
4691
               "algorithm": "FedIIR",
4692
               ""description"": "Federated domain generalization-based technique
4693
                    that identifies and leverages invariant relationships across
4694
                    domains to enhance generalization to out-of-distribution
4695
                   data in federated settings."
4696
4697
```

### C.5 LLMs as the agent core components

#### MODEL SELECTION JUSTIFICATION

To assess the reasoning, planning, and tool-use capabilities of large language model (LLM) agents in the context of real-world federated learning workflows, we evaluate a set of 24 LLMs on the FedAgentBench suite. The selected models span both proprietary and open-source categories, ensuring broad coverage across scale, training data diversity, and model access paradigms.

We include 10 proprietary LLMs from leading industrial labs such as OpenAI and Anthropic, including multiple variants of GPT-4. These models represent the current frontier of general-purpose foundation models, often topping benchmarks in instruction-following, tool use, and reasoning. Their inclusion allows us to benchmark state-of-the-art commercial performance in the agentic FL setting.

We particularly include a range of GPT-family models developed by OpenAI to cover both ends of the performance-efficiency spectrum in proprietary large language models (LLMs). The rationale is threefold:

## (i) Proven Instruction-Following and Reasoning Abilities:

GPT-4 and its variants have consistently demonstrated state-of-the-art performance across multiple benchmarks involving instruction following, task decomposition, and multi-step reasoning capabilities essential for evaluating LLM agents in complex federated learning pipelines such as FedAgentBench.

#### (ii) Variants across Performance Tiers and Costs:

The selection spans high-end models (e.g., GPT-4.1, GPT-40) and lightweight alternatives (e.g., GPT-4.1-mini, GPT-o3-mini). This allows us to study the trade-offs between agent reasoning quality and computational/resource efficiency, particularly relevant for real-world FL deployment where cost and inference speed matter.

#### (iii) Industry Adoption and API Availability:

These models are widely adopted in both academic and industrial applications and offer stable, reproducible APIs. This ensures consistent evaluation and compatibility with tool-augmented LLM agent frameworks.

Besides, we evaluate 14 open-source LLMs across four major families: LLaMA, DeepSeek, Qwen, and Gemma. These models are chosen for their state-of-the-art performance in open benchmarks, availability in multiple parameter scales (from 9B to 685B), and varying architectural innovations (e.g., distillation in DeepSeek, instruction tuning in Qwen, and scalability in Gemma). This selection ensures a representative spectrum of recent advances in open-source LLM development, and provides insight into how scale, family, and fine-tuning affect FL-agent performance.

By including both proprietary and open models across diverse sizes and pretraining paradigms, our evaluation is designed to offer fair, scalable, and realistic comparisons, while informing the community of strengths and limitations across model categories in complex multi-agent settings like FedAgentBench.

## PROPRIETARY MODELS

**GPT-4.1** OpenAI (2024a) is OpenAI's 2024 flagship model, offering exceptional reasoning, planning, and tool-use capabilities. It supports chain-of-thought reasoning, structured JSON generation, and API-based tool calling. In FedAgentBench, GPT-4.1 serves as a strong proprietary baseline for task decomposition and system orchestration. However, its high inference latency and cost make it better suited for final evaluation rather than real-time agent interactions. It supports only textual inputs (with multimodal functionality present in GPT-4o), offers excellent token consistency, and performs strongly on multilingual tasks. Architecturally, it is assumed to be a dense transformer with proprietary design and trained on a broad, high-quality dataset spanning web, academic, and code sources.

**GPT-40** OpenAI (2024c) is OpenAI's latest multimodal model, offering native support for both vision and text. It is significantly faster and cheaper than GPT-4 while preserving high-quality output. In FedAgentBench, it is leveraged for tasks that require vision-language reasoning or real-time tool use. Its stable token formatting, rapid inference, and strong multilingual capabilities make it highly

Table 8: Descriptions for Proprietary LLMs in FedAgentBench

Model	Description	Capabilities	Use Rationale	Caveats / Notes
GPT-4.1	Latest high-performance	Chain-of-thought reason-	Reference propri-	High cost and latency; not
	model from OpenAI with	ing, tool use, structured	etary agent for	ideal for real-time execu-
	advanced reasoning and	outputs.	end-to-end work-	tion.
	planning.		flows.	
GPT-40	Multimodal flagship	Multilingual, tool calling,	Evaluated for vision	New model; some outputs
	model supporting vision-	multimodal reasoning.	+ tool scenarios.	may vary between calls.
	language tasks.			
GPT-4	Original GPT-4 model	Long-context, reasoning,	Used as baseline for	Slower than turbo and
	with top-tier generaliza-	structured outputs.	reasoning accuracy.	newer variants.
	tion.			
GPT-4-Turbo	Faster and cheaper ver-	Efficient inference, simi-	Preferred when cost	Slightly less coherent out-
	sion of GPT-4 for API	lar capabilities to GPT-4.	is a concern.	puts.
	use.			
GPT-4.1-mini	Distilled variant opti-	Good single-step logic,	Used in real-time as-	Weaker on edge-case and
	mized for fast inference.	mid-range planning.	sistant agents.	ambiguous tasks.
GPT-4o-mini	Smaller variant of GPT-	Vision-language support,	Benchmarked in	Reduced performance in
	40 with multimodal sup-	low-latency.	low-resource multi-	logic-intensive tasks.
	port.		modal agents.	
GPT-o4-mini	Lightweight GPT-4 style	Text generation and sim-	Ablation studies	Unclear origin; may alias
	model.	ple instructions.	for low-cost GPT	other mini variants.
			agents.	
GPT-o3-mini	GPT-3.5-based efficient	Very fast, single-turn	Used for comparison	Weak reasoning; not reli-
	variant.	chat.	with older architec-	able for planning.
			tures.	
GPT-3.5 Turbo	Predecessor to GPT-4,	Fast, capable for basic in-	Low-cost reference	Token alignment issues in
	cheaper and widely used.	struction and QA.	for proprietary	structured tasks.
			agents.	
Claude-3.7 Sonnet	Mid-size model from An-	Safety-aligned genera-	Benchmarked	Slightly lower fluency
	thropic with alignment	tion, multilingual, tool	against non-OpenAI	than top Claude variants.
	tuning.	use.	proprietary model.	_

Table 9: Descriptions for Open-Source LLMs in FedAgentBench

Model	Description	Capabilities	Use Rationale	Caveats / Notes
LLaMA-4 Maverick	Latest LLaMA release (2025) with top-tier accuracy in reasoning and instruction following.	Instruction following, long-context reasoning, coding tasks.	Used for evaluating high- end open-source agents.	Resource heavy; slower than lighter LLaMA variants.
LLaMA-4 Scout	2025 LLaMA-4 variant optimized for cost-efficient inference.	Balanced reasoning and fast response for system agents.	Used as mid-range open- source agent in system and logic tasks.	Less expressive than Maverick.
LLaMA-3 70B	Flagship LLaMA model (2024) with extensive instruction tuning.	Reasoning, multilingual tasks, tool use.	Used for top-tier open- source evaluation.	Less performant than newer LLaMA- 4 variants.
LLaMA-3 8B	Smaller variant of LLaMA-3 for constrained environments.	General understanding, good for fast responses.	Used in real-time bench- marking of lighter agents.	Limited capacity in multi-hop reasoning.
DeepSeek-V3	Latest release from DeepSeek with strong Chinese-English capabil- ity.	Multilingual chat, code, reasoning.	Used to test multilingual and cross-domain agents.	Less stable tool usage.
DeepSeek-R1	General purpose 2024 DeepSeek model.	Basic LLM tasks, reason- ing.	Baseline open-source ref- erence.	Lower precision un- der stress tests.
DeepSeek-R1- Distill	Distilled version of DeepSeek-R1 on LLaMA-70B.	Fast inference, low-resource usage.	Used in lightweight eval- uations.	Lower performance ceiling.
Qwen 3 235B	Massive MoE model by Alibaba; high capacity and strong multilingual.	Multilingual, few-shot generalization, long context.	Benchmarked as high- capacity open-source agent.	Costly to run, sparse documentation.
Qwen QwQ 32B	Intermediate-sized multi- lingual Qwen model.	Instruction following, QA, multilingual chat.	Used as cost-performance mid-range Qwen agent.	Less stable tool usage.
Qwen 3 30B	Well-balanced Qwen variant.	Reliable output, struc- tured reasoning.	Used in systems requiring stable decoding.	Reduced multilingual coverage vs 235B.
Qwen 3 14B	Smaller Qwen for lightweight use.	Quick single-turn tasks.	Used in sub-agents and pre-filtering roles.	Shallow reasoning, poor long-context.
Gemma 3 27B Instruct	Instruction-tuned model by Google.	Tool use, summarization, chat.	Tested for logic tasks.	Less capable in multi-modal do- mains.
Gemma 3 12B Instruct	Smaller Gemma variant.	Common NLP tasks.	System-level fast agent.	May misfire struc- tured outputs.
Gemma 2 9B Instruct	Previous generation Gemma model.	Lightweight inference.	Tested in low-cost agent scenarios.	Lowest instruction accuracy among Gemmas.

versatile. GPT-40 is believed to use a dense architecture, trained on diverse visual and textual corpora. The model is ideal for agent scenarios requiring fast multimodal perception and flexible interaction.

**GPT-4** OpenAI (2023b) is the original release in the GPT-4 series, known for its top-tier reasoning and planning abilities. It supports structured input/output and long-context tasks, making it valuable in FedAgentBench for evaluating logical and multi-step reasoning pipelines. Despite its strengths, it is relatively slower and more costly than newer variants. The model supports only text and maintains high token consistency. Though its architecture remains undisclosed, it is believed to be a dense transformer trained on an extensive web-scale corpus.

**GPT-4-Turbo** OpenAI (2024b) is an optimized variant of GPT-4 that retains strong performance while offering faster responses and lower API cost. It is used in FedAgentBench for latency-sensitive tasks where full GPT-4 performance is not critical. The model performs well in multilingual and structured generation scenarios and maintains good token consistency. Though slightly less coherent than GPT-4 in edge cases, it is well-suited for system-level automation. It is assumed to be a compressed variant of GPT-4 trained on the same high-diversity corpus.

**GPT-4.1-mini** OpenAI (2025a) is a distilled, lighter-weight version of GPT-4.1, optimized for efficiency. It is ideal for deploying system-level agents that require faster responses and lower resource consumption. Used in FedAgentBench for agents like data pre-processors, it supports mid-range reasoning but struggles with highly complex logic. The model is purely textual, offers moderate-to-good token formatting, and has a reduced training footprint drawn from GPT-4.1's corpus.

**GPT-4o-mini** OpenAI (2024d) brings the benefits of multimodal reasoning into a low-cost, fast-inference package. It is used in FedAgentBench to simulate real-world constraints like limited compute. The model supports vision-language inputs and outputs, performs well in multilingual tasks, and maintains good JSON formatting for tool use. Architecturally, it is assumed to be a miniaturized MoE variant with a training corpus inclusive of visual and textual data.

**o4-mini** OpenAI (2025c) is a presumed lighter version of the GPT family with unclear lineage, possibly a prototype or early mini model. It is included in FedAgentBench as a lower-bound proprietary benchmark, offering fast and low-cost inference. It lacks multimodal support, has only fair token consistency, and offers limited reasoning depth. Its training data is assumed to be a subset of GPT-4.1.

**GPT-o3-mini** OpenAI (2025b) is an earlier mini variant from the GPT-3 or GPT-4 lineage, optimized for speed over reasoning complexity. It serves as a baseline ablation reference in FedAgentBench. While offering very fast inference, it struggles with deep reasoning and structured outputs. It is a text-only model with moderate token consistency and basic multilingual support.

**GPT-3.5 Turbo** OpenAI (2023a) is a highly popular model with good general-purpose performance and fast response times. It is used in FedAgentBench as a mid-tier proprietary baseline. While it offers strong text generation, it can be inconsistent with structured formats like JSON. The model supports only text, has limited multilingual capabilities, and is based on OpenAI's GPT-3 architecture with training data from broad internet and code sources.

Claude-3.7 Sonnet Anthropic (2024) is Anthropic's mid-tier model, emphasizing safety, reasoning, and alignment. It is included in FedAgentBench to benchmark performance beyond OpenAI's ecosystem. The model supports multimodal input, generates well-aligned structured responses, and maintains strong multilingual performance. It uses a dense transformer architecture influenced by constitutional AI principles and is trained on a diverse dataset spanning web, academic, and code domains.

#### **OPEN-SOURCE MODELS**

**LLaMA-4 Maverick** Meta-LLaMA (2025) is Meta's most advanced model in the LLaMA family, expected to reach or surpass GPT-4 performance in reasoning and alignment. It is used in FedAgentBench for high-fidelity planning and tool execution in open-source settings. The model is purely text-based but offers strong structured output consistency and robust multilingual capabilities. Architecture-wise, it uses a dense transformer and is trained on a vast mix of web, academic, and code datasets. Due to its size and complexity, inference cost and latency are relatively high.

**LLaMA-4 Scout** Meta-LLaMA (2025) is a lighter-weight sibling of Maverick, designed for faster inference with modest performance trade-offs. It supports most of the same reasoning and generation capabilities but with lower memory and compute demands. In FedAgentBench, Scout is used for scenarios balancing performance with latency. It is purely text-based, has good JSON consistency, and performs well across multiple languages. It is built on a dense architecture and inherits a subset of Maverick's training data.

**LLaMA-3 70B** Meta-LLaMA (2024) is a powerful open-source alternative to proprietary giants, known for high reasoning quality and multilingual fluency. In FedAgentBench, it serves as the default open-source baseline for high-end agent tasks. It supports structured outputs and follows instructions reliably. Though limited to text, it demonstrates high token consistency and robust zero-shot capabilities. Architecturally, it's a dense decoder-only transformer trained on high-quality web and code data.

**LLaMA-3 8B** Meta-LLaMA (2024) is a compact version suited for environments where computational resources are constrained. It is used in FedAgentBench for resource-sensitive deployments. While it has limited reasoning depth compared to its 70B counterpart, it is cost-effective and maintains good multilingual and instruction-following capabilities. It does not support multimodality and is less reliable for long context or multi-step planning.

**DeepSeek-V3** DeepSeek (2024b) is a state-of-the-art open-source LLM with strong multilingual proficiency, especially in Chinese-English bilingual tasks. It is used in FedAgentBench for cross-lingual task evaluation and instruction-following. V3 demonstrates high token consistency, supports function calling, and performs competitively with models twice its size. Though it is a text-only model, its dense architecture enables fast inference. The training set includes code, academic papers, and bilingual web content.

**DeepSeek-R1** DeepSeek (2024a) emphasizes efficiency and alignment, making it suitable for interactive, tool-augmented agents in FedAgentBench. DeepSeek-R1 performs well in text generation, tool use, and multilingual dialogue. It does not support images but excels in structured response formatting. The model balances inference cost with reasonable output quality, making it ideal for mid-range agents in the benchmark.

**DeepSeek-R1-Distill-LLaMA-70B** DeepSeek-AI (2025) is a distilled version of the LLaMA-3 70B model. It maintains much of the performance while significantly reducing latency and memory usage. It is used in FedAgentBench for simulating efficient client-side agents. Although it lacks vision capabilities, it provides consistent JSON formatting and supports Chinese-English reasoning tasks effectively. It uses a compressed architecture optimized for deployment.

**Qwen 3 235B** Alibaba (2025) is Alibaba's largest open-source MoE (mixture-of-experts) model, offering cutting-edge multilingual and reasoning performance. It supports structured tool use and long-context generation. In FedAgentBench, it is used for evaluating large-scale open-source agents that can rival proprietary models. It lacks native vision support but performs exceptionally well in code generation and Chinese-English tasks. Due to its MoE design, only a subset of parameters is active per token, balancing performance and compute.

**Qwen QwQ 32B** Alibaba (2025) is optimized for fast inference while preserving high reasoning capabilities. It supports structured outputs and bilingual fluency and is used in FedAgentBench for tasks requiring instruction-following and tool use. Although not as powerful as the 235B variant, it has a more favorable speed-to-quality ratio. The model is purely text-based and trained on code-rich and multilingual corpora.

**Qwen 3 30B** Alibaba (2025) is a dense decoder-only transformer known for instruction tuning and reasoning efficiency. It performs well on tool-augmented tasks in FedAgentBench and has high token consistency. It does not support vision inputs but is trained on a diverse corpus including code, web, and academic sources. It is used for open-source system-level agents that balance size and capability.

**Qwen 3 14B** Alibaba (2025) is the smallest Qwen model used in FedAgentBench, yet it retains strong instruction-following and token generation properties. It is deployed in tasks that require fast turnaround or operate in limited environments. The model lacks multimodal support but remains reliable for multilingual and structured outputs. It is a dense transformer trained on similar data to its larger siblings.

Table 10: Comparison of open-source and proprietary LLM agents across different stages of federated learning: Client Selection (Client-Sel), Data Pre-processing (Data-Pre), Label Harmonization (Label-Harm), and Federated Training (Fed-Train) in **Dermatology** environment based on skin cancer detection task. a/b refers to the proportion of successful runs 'a' out of the total number of runs 'b'

Model		Fine-g	grained guidance			Goal-oriented guidance					
	Client-Sel	Data-Pre	Label-Harm	Fed-Train	Overall	Client-Sel	Data-Pre	Label-Harm	Fed-Train	Overall	
	$a_1, a_2, a_3$	$a_4$	$a_5$	$a_6, a_7$		$a_1, a_2, a_3$	$a_4$	$a_5$	$a_6, a_7$		
				Proprietar	y Models						
GPT-4.1	5/5, 5/5, 5/5	5/5	3/5	5/5, 5/5	94.29	5/5, 4/5, 5/5	5/5	3/5	4/5, 5/5	88.57	
GPT-4o	5/5, 3/5, 5/5	5/5	1/5	1/5, 5/5	71.43	5/5, 1/5, 5/5	5/5	1/5	1/5, 5/5	65.71	
GPT-4	5/5, 4/5, 5/5	0/5	1/5	3/5, 5/5	65.71	5/5, 1/5, 5/5	0/5	0/5	2/5, 5/5	51.43	
GPT-4-Turbo	5/5, <mark>3/5</mark> , 5/5	2/5	1/5	3/5, 5/5	68.57	5/5, 3/5, 5/5	5/5	1/5	2/5, 5/5	74.29	
GPT-4.1-mini	5/5, 5/5, 5/5	5/5	3/5	3/5, 5/5	88.57	5/5, 5/5, 5/5	3/5	3/5	3/5, 5/5	82.86	
GPT-4o-mini	5/5, 1/5, 3/5	5/5	3/5	3/5, 4/5	68.57	5/5, 0/5, 3/5	5/5	1/5	2/5, 4/5	57.14	
GPT-o4-mini	5/5, 4/5, 5/5	5/5	3/5	3/5, 5/5	85.71	5/5, 3/5, 5/5	4/5	2/5	3/5, 4/5	74.29	
GPT-o3-mini	5/5, 3/5, 5/5	0/5	2/5	3/5, 5/5	65.71	5/5, 1/5, 5/5	0/5	2/5	3/5, 5/5	60.00	
GPT-3.5-Turbo	5/5, 0/5, 0/5	0/5	0/5	1/5, 3/5	25.71	5/5, 0/5, 0/5	2/5	0/5	1/5, 3/5	31.43	
Claude-3-7-Sonnet	5/5, 2/5, 3/5	2/5	1/5	2/5, 3/5	51.42	5/5, 2/5, 3/5	2/5	1/5	2/5, 5/5	57.14	
				Open-sour	ce Models						
				Huge M	Iodels						
DeepSeek-V3	5/5, 1/5, 5/5	5/5	5/5	4/5, 5/5	85.71	5/5, 1/5, 5/5	4/5	4/5	4/5, 5/5	80.00	
DeepSeek-R1	5/5, 0/5, 5/5	0/5	0/5	0/5, 5/5	42.86	5/5, 0/5, 5/5	0/5	0/5	0/5, 5/5	42.85	
Qwen3 235B	5/5, 0/5, 5/5	0/5	0/5	0/5, 5/5	42.86	5/5, 0/5, 5/5	0/5	0/5	0/5, 5/5	42.85	
LLaMA-4 Maverick	5/5, 1/5, 4/5	5/5	3/5	2/5, 5/5	71.43	5/5, 1/5, 4/5	5/5	3/5	3/5, 5/5	74.29	
LLaMA-4 Scout	5/5, 1/5, 5/5	5/5	3/5	2/5, 5/5	74.29	5/5, 2/5, 5/5	5/5	3/5	2/5, 5/5	77.14	
				Large N	Aodels						
DeepSeek-R1-70B	5/5, 0/5, 5/5	0/5	0/5	1/5, 5/5	45.71	5/5, 0/5, 5/5	0/5	0/5	0/5, 5/5	42.86	
LLaMA-3-70B	5/5, 0/5, 5/5	1/5	1/5	2/5, 5/5	54.29	5/5, 0/5, 5/5	2/5	2/5	1/5, 5/5	57.14	
				Medium	Models						
Owen OwO 32B	5/5, 4/5, 5/5	5/5	4/5	4/5, 5/5	91.43	5/5, 4/5, 5/5	5/5	3/5	3/5, 5/5	85.71	
Owen3-30B	5/5, 0/5, 5/5	0/5	0/5	1/5, 5/5	45.71	5/5, 0/5, 5/5	0/5	0/5	1/5, 5/5	45.71	
Gemma3-27B-instruct	5/5, 0/5, 0/5	0/5	0/5	0/5, 0/5	14.29	5/5, 0/5, 0/5	0/5	0/5	0/5, 0/5	14.29	
			,	Small M	Iodels			,			
Gemma-2-9B	5/5, 0/5, 5/5	1/5	1/5	1/5, 5/5	51.43	5/5, 0/5, 5/5	1/5	1/5	1/5, 5/5	51.43	
LLaMA-3-8B	5/5, 0/5, 5/5	5/5	2/5	1/5, 5/5	65.71	5/5, 0/5, 5/5	5/5	2/5	1/5, 5/5	65.71	
Owen-3-14B	5/5, 0/5, 5/5	0/5	0/5	0/5, 5/5	42.86	5/5, 0/5, 5/5	0/5	0/5	0/5, 4/5	40.00	
Gemma3-12B-instruct	5/5, 0/5, 0/5	0/5	0/5	0/5, 0/5	14.29	5/5, 0/5, 0/5	0/5	0/5	0/5, 0/5	14.29	

Gemma 3 27B Instruct DeepMind (2024) is Google's latest instruction-tuned model offering high-quality responses in English and multilingual contexts. It is used in FedAgentBench as an open-source baseline for structured agent reasoning. It provides consistent tool output and supports extended context. Though it is a text-only model, its training set is diverse and includes curated multilingual, code, and web data. It is efficient for its size and performs competitively with larger models.

**Gemma 3 12B Instruct** DeepMind (2024) is a smaller variant of Gemma used for lower-latency deployments in FedAgentBench. It supports basic reasoning, instruction tuning, and structured outputs, though with limitations in depth and contextuality. It is well suited for client-side agents and lower-resource inference. The model maintains high token consistency and multilingual robustness despite being smaller.

**Gemma 2 9B Instruct** DeepMind (2024) is a previous-generation lightweight model that still performs reliably for text classification, generation, and control tasks. It is used in FedAgentBench for fast-response ablation studies and scenarios where minimal hardware is available. The model is trained on public datasets including multilingual and code domains, making it a versatile baseline.

#### D RESULTS AND DISCUSSIONS

We conducted extensive evaluations of both proprietary and open-source LLM agents across 6 environments, out of which the success rates for Histopathology have been mentioned in the main paper. The success rates for the remaining 5 environments *viz.*, Dermatology, Ultrasound, MRI, Fundus and X-Ray environments are reported here. The results of these experiments are presented in Tables 10-15. These tables capture performance under two paradigms: fine-grained multi-step guidance and goal-oriented single-shot instruction, revealing consistent trends across modalities. Notably, the independent script generation setting in Table 12 illustrates a sharp decline in performance for most agents, underscoring the challenges of long-horizon task planning without explicit decomposition. Overall Time-requirement metrics for task resolution are summarized in Table 16, providing a holistic view of capability and practicality across LLM variants. Figs 9-35 show snippets of different phases of the FL workflow with various LLMs and different imaging modalities which help to understand their success and failure modes.

Table 11: Comparison of open-source and Proprietary LLM agents in Ultrasound environment

Model		Fine-g	rained guidance			Goal-oriented guidance					
	Client-Sel	Data-Pre	Label-Harm	Fed-Train	Overall	Client-Sel	Data-Pre	Label-Harm	Fed-Train	Overall	
	$S_1, C_1, S_2$	$C_2$	$C_3$	$S_3, S_4$		$S_1, C_1, S_2$	$C_2$	$C_3$	$S_3, S_4$		
				Proprietar	y Models						
GPT-4.1	5/5, 3/5, 5/5	5/5	5/5	5/5, 5/5	94.29	5/5, 3/5, 5/5	5/5	5/5	5/5, 5/5	94.29	
GPT-40	5/5, 0/5, 5/5	5/5	3/5	1/5, 5/5	68.57	5/5, 0/5, 5/5	5/5	2/5	1/5, 5/5	65.71	
GPT-4	5/5, 3/5, 5/5	1/5	1/5	3/5, 5/5	65.71	5/5, 3/5, 5/5	0/5	1/5	3/5, 5/5	62.86	
GPT-4-Turbo	5/5, <mark>3/5</mark> , 5/5	1/5	1/5	3/5, 5/5	65.71	5/5, 3/5, 5/5	4/5	1/5	3/5, 5/5	74.29	
GPT-4.1-mini	$5/5, \frac{3}{5}, 5/5$	5/5	3/5	4/5, 5/5	85.71	5/5, 2/5, 5/5	3/5	4/5	3/5, 5/5	77.14	
GPT-4o-mini	5/5, 1/5, 3/5	5/5	3/5	3/5, 5/5	71.43	5/5, 1/5, 3/5	5/5	1/5	5/5, 5/5	71.43	
GPT-o4-mini	5/5, 3/5, 5/5	5/5	3/5	4/5, 5/5	85.71	5/5, 3/5, 5/5	4/5	3/5	4/5, 5/5	82.86	
GPT-o3-mini	5/5, 2/5, 5/5	1/5	1/5	3/5, 4/5	60.00	5/5, 1/5, 5/5	1/5	2/5	3/5, 5/5	62.86	
GPT-3.5-Turbo	5/5, 0/5, 0/5	0/5	0/5	1/5, 3/5	25.71	5/5, 0/5, 0/5	2/5	0/5	1/5, 4/5	34.29	
Claude-3-7	5/5, 2/5, 3/5	2/5	1/5	3/5, 3/5	54.29	5/5, 2/5, 3/5	2/5	1/5	3/5, 3/5	54.29	
				Open-sour	ce Models						
				Huge M	Iodels						
DeepSeek-V3	5/5, 3/5, 5/5	5/5	5/5	4/5, 5/5	91.43	5/5, 2/5, 5/5	4/5	5/5	4/5, 5/5	85.71	
DeepSeek-R1	5/5, 0/5, 5/5	0/5	0/5	0/5, 5/5	42.86	5/5, 0/5, 5/5	0/5	0/5	0/5, 5/5	42.86	
Qwen3 235B	$5/5, \frac{0}{5}, \frac{5}{5}$	0/5	0/5	0/5, 5/5	42.86	$5/5, \frac{0}{5}, \frac{5}{5}$	0/5	0/5	0/5, 5/5	42.86	
LLaMA-4 Maverick	$5/5, \frac{3}{5}, \frac{5}{5}$	5/5	4/5	4/5, 5/5	88.57	$5/5, \frac{3}{5}, \frac{5}{5}$	5/5	3/5	3/5, 5/5	82.86	
LLaMA-4 Scout	$5/5, \frac{3}{5}, \frac{5}{5}$	5/5	4/5	3/5, 5/5	85.71	5/5, 1/5, 5/5	5/5	3/5	2/5, 5/5	74.28	
				Large N	Aodels	•					
DeepSeek-R1-70B	$5/5, \frac{3}{5}, 5/5$	3/5	1/5	2/5, 5/5	74.28	5/5, 1/5, 5/5	3/5	0/5	2/5, 5/5	68.57	
LLaMA-3-70B	5/5, 3/5, 5/5	2/5	2/5	2/5, 5/5	68.57	5/5, 3/5, 5/5	2/5	2/5	2/5, 5/5	68.57	
				Medium	Models						
Qwen QwQ 32B	5/5, 3/5, 5/5	4/5	4/5	4/5, 5/5	85.71	5/5, 3/5, 5/5	2/5	4/5	4/5, 5/5	80.00	
Qwen3-30B	5/5, 0/5, 5/5	0/5	0/5	1/5, 5/5	45.71	5/5, 0/5, 5/5	0/5	0/5	1/5, 5/5	45.71	
Gemma3-27B-instruct	5/5, 0/5, 0/5	0/5	0/5	0/5, 0/5	14.29	5/5, 0/5, 0/5	0/5	0/5	0/5, 0/5	14.29	
				Small M	Iodels						
Gemma-2-9B	5/5, 0/5, 4/5	2/5	1/5	1/5, 5/5	51.43	5/5, 0/5, 4/5	1/5	1/5	1/5, 5/5	48.57	
LLaMA-3-8B	5/5, 0/5, 4/5	4/5	2/5	1/5, 5/5	60.00	5/5, 0/5, 4/5	4/5	2/5	1/5, 5/5	60.00	
Qwen-3-14B	5/5, 0/5, 0/5	0/5	0/5	0/5, 5/5	28.57	5/5, 0/5, 0/5	0/5	0/5	0/5, 4/5	25.71	
Gemma3-12B-instruct	5/5, 0/5, 0/5	0/5	0/5	0/5, 0/5	14.29	5/5, 0/5, 0/5	0/5	0/5	0/5, 0/5	14.29	

Table 12: Comparison of open-source and Proprietary LLM agents in **Ultrasound** environment on independent script generation for solving individual task. Striking reduction in performance is noted for all agents

Model	Client-Sel	Data-Pre	Label-Harm	Fed-Train	Overall					
	Pro	prietary Mode	els	•						
GPT-4.1	5/5, 0/5, 5/5	0/5	0/5	0/5, 5/5	42.86					
GPT-40	5/5, 0/5, 5/5	0/5	0/5	0/5, 5/5	42.86					
GPT-4	5/5, 0/5, 5/5	0/5	0/5	0/5, 5/5	42.86					
GPT-4-Turbo	5/5, 0/5, 5/5	0/5	0/5	0/5, 5/5	42.86					
GPT-4.1-mini	5/5, 0/5, 5/5	0/5	0/5	0/5, 5/5	42.86					
GPT-4o-mini	5/5, 0/5, 3/5	0/5	0/5	0/5, 5/5	37.14					
GPT-o4-mini	5/5, 0/5, 5/5	0/5	0/5	0/5, 5/5	42.86					
GPT-o3-mini	5/5, 0/5, 5/5	0/5	0/5	0/5, 5/5	42.86					
GPT-3.5-Turbo	5/5, 0/5, 0/5	0/5	0/5	0/5, 4/5	25.71					
Claude-3-7	5/5, 0/5, 3/5	0/5	0/5	0/5, 3/5	31.43					
	Ope	n-source Mod	els	•	•					
Huge Models										
DeepSeek-V3	5/5, 0/5, 5/5	0/5	0/5	0/5, 5/5	42.86					
DeepSeek-R1	5/5, 0/5, 5/5	0/5	0/5	0/5, 5/5	42.86					
Qwen3 235B	5/5, 0/5, 5/5	0/5	0/5	0/5, 5/5	42.86					
LLaMA-4 Maverick	5/5, 0/5, 5/5	0/5	0/5	0/5, 5/5	42.86					
LLaMA-4 Scout	$5/5, \frac{0}{5}, \frac{5}{5}$	0/5	0/5	0/5, 5/5	42.86					
	I	arge Models								
DeepSeek-R1-70B	5/5, 0/5, 5/5	0/5	0/5	0/5, 5/5	42.86					
LLaMA-3-70B	5/5, 0/5, 5/5	0/5	0/5	0/5, 5/5	42.86					
	M	edium Model	S							
Qwen QwQ 32B	5/5, 0/5, 5/5	0/5	0/5	0/5, 5/5	42.86					
Qwen3-30B	5/5, 0/5, 5/5	0/5	0/5	0/5, 5/5	42.86					
Gemma3-27B-instruct	5/5, 0/5, 0/5	0/5	0/5	0/5, 0/5	14.29					
	S	mall Models								
Gemma-2-9B	5/5, 0/5, 4/5	0/5	0/5	0/5, 5/5	40.0					
LLaMA-3-8B	5/5, 0/5, 4/5	0/5	0/5	0/5, 5/5	40.0					
Qwen-3-14B	5/5, 0/5, 0/5	0/5	0/5	0/5, 4/5	25.71					
Gemma3-12B-instruct	5/5, 0/5, 0/5	0/5	0/5	0/5, 0/5	14.29					

### D.1 DISCUSSION ON AGENTIC PERFORMANCE IN INDIVIDUAL HEALTHCARE ENVIRONMENT

Table 10 shows the performance of open-source and proprietary LLM agents in the **Dermatology environment**. Across both guidance regimes, proprietary models outperform open-source baselines. **GPT-4.1** is the top performer, reaching near-ceiling success with consistent 5/5 across most subtasks and the highest **Overall** scores (**94.29** under fine-grained guidance; **88.57** under goal-oriented guidance). A second tier of proprietary models follows: **GPT-4.1-mini** (88.57 / 82.86), **GPT-04-mini** (85.71 / 74.29), **GPT-40** (71.43 / 65.71), **GPT-4-Turbo** (68.57 / 74.29), **GPT-03-mini** (65.71 / 60.00), and **GPT-4** (65.71 / 51.43). **GPT-3.5-Turbo** lags markedly (25.71 / 31.43), and **Claude-3-7-Sonnet** 

Table 13: Comparison of open-source and Proprietary LLM agents in MRI environment

Model		Fine-g	grained guidance				Goal-o	riented guidance	,	
	Client-Sel	Data-Pre	Label-Harm	Fed-Train	Overall	Client-Sel	Data-Pre	Label-Harm	Fed-Train	Overall
GPT-4.1	5/5, 5/5, 5/5	5/5	5/5	5/5, 5/5	100.00	5/5, 5/5, 5/5	5/5	5/5	5/5, 5/5	100.00
GPT-40	$5/5, \frac{3/5}{5}, 5/5$	5/5	4/5	1/5, 5/5	71.43	$5/5, \frac{3/5}{5}, 5/5$	5/5	3/5	1/5, 5/5	68.57
GPT-4	5/5, 5/5, 5/5	1/5	2/5	3/5, 5/5	71.43	$5/5, \frac{4/5}{5}, 5/5$	0/5	1/5	3/5, 5/5	65.71
GPT-4-Turbo	5/5, 5/5, 5/5	1/5	2/5	3/5, 5/5	71.43	$5/5, \frac{4/5}{5}, 5/5$	4/5	1/5	3/5, 5/5	77.14
GPT-4.1-mini	$5/5, \frac{4/5}{5}, 5/5$	5/5	3/5	4/5, 5/5	88.57	$5/5, \frac{3/5}{5}, 5/5$	3/5	3/5	3/5, 5/5	77.14
GPT-4o-mini	5/5, 3/5, 3/5	5/5	3/5	3/5, 5/5	77.14	5/5, 2/5, 3/5	5/5	2/5	5/5, 5/5	74.29
GPT-o4-mini	5/5, 5/5, 5/5	5/5	3/5	4/5, 5/5	91.43	$5/5, \frac{4/5}{5}, 5/5$	4/5	2/5	4/5, 5/5	85.71
GPT-o3-mini	5/5, 5/5, 5/5	1/5	1/5	4/5, 4/5	71.42	$5/5, \frac{4/5}{5}, 5/5$	1/5	1/5	4/5, 5/5	74.29
GPT-3.5-Turbo	5/5, 0/5, 0/5	0/5	1/5	1/5, 3/5	25.71	5/5, 0/5, 0/5	2/5	0/5	1/5, 4/5	34.29
Claude-3-7	5/5, 4/5, 3/5	2/5	1/5	4/5, 3/5	57.14	5/5, 3/5, 3/5	2/5	1/5	3/5, 3/5	57.14
				Open-sour	ce Models	•				
				Huge M	Iodels					
DeepSeek-V3	$5/5, \frac{4/5}{5}, 5/5$	5/5	5/5	4/5, 5/5	94.29	5/5, <mark>3/5</mark> , 5/5	4/5	5/5	4/5, 5/5	88.57
DeepSeek-R1	$5/5, \frac{2}{5}, \frac{5}{5}$	0/5	0/5	0/5, 5/5	42.86	$5/5, \frac{0}{5}, \frac{5}{5}$	0/5	0/5	0/5, 5/5	42.86
Qwen3 235B	$5/5, \frac{2}{5}, \frac{5}{5}$	0/5	0/5	0/5, 5/5	42.86	5/5, 1/5, 5/5	0/5	0/5	0/5, 5/5	42.86
LLaMA-4 Maverick	5/5, 5/5, 5/5	5/5	4/5	4/5, 5/5	94.29	$5/5, \frac{4/5}{5}, 5/5$	5/5	3/5	3/5, 5/5	85.71
LLaMA-4 Scout	$5/5, \frac{4}{5}, \frac{5}{5}$	5/5	4/5	2/5, 5/5	85.71	$5/5, \frac{3}{5}, 5/5$	5/5	3/5	2/5, 5/5	74.29
				Large N	Aodels					
DeepSeek-R1-70B	5/5, 5/5, 5/5	3/5	1/5	2/5, 5/5	74.29	$5/5, \frac{4/5}{5}, \frac{5}{5}$	3/5	0/5	2/5, 5/5	68.57
LLaMA-3-70B	$5/5, \frac{4/5}{5}, 5/5$	2/5	2/5	2/5, 5/5	71.43	$5/5, \frac{4/5}{5}, 5/5$	2/5	2/5	2/5, 5/5	71.43
				Medium	Models					
Qwen QwQ 32B	$5/5, \frac{4}{5}, \frac{5}{5}$	4/5	4/5	4/5, 5/5	88.57	$5/5, \frac{4/5}{5}, \frac{5}{5}$	2/5	4/5	4/5, 5/5	82.86
Qwen3-30B	$5/5, \frac{2}{5}, \frac{5}{5}$	0/5	0/5	1/5, 5/5	48.57	$5/5, \frac{0}{5}, \frac{5}{5}$	0/5	0/5	1/5, 5/5	45.71
Gemma3-27B-instruct	5/5, 1/5, 2/5	0/5	0/5	0/5, 0/5	14.29	5/5, 0/5, 2/5	0/5	0/5	0/5, 0/5	14.29
				Small N	Iodels					
Gemma-2-9B	5/5, 1/5, 4/5	2/5	1/5	1/5, 5/5	51.43	5/5, 1/5, 4/5	1/5	1/5	1/5, 5/5	48.57
LLaMA-3-8B	5/5, 3/5, 4/5	4/5	2/5	1/5, 5/5	62.86	5/5, 2/5, 4/5	4/5	2/5	1/5, 5/5	60.00
Qwen-3-14B	5/5, 1/5, 2/5	0/5	0/5	0/5, 5/5	28.57	5/5, 0/5, 2/5	0/5	0/5	0/5, 4/5	25.71
Gemma3-12B-instruct	5/5, 1/5, 2/5	0/5	0/5	0/5, 0/5	14.29	5/5, 0/5, 2/5	0/5	0/5	0/5, 0/5	14.29

Table 14: Comparison of open-source and Proprietary LLM agents in Fundus environment

Model		Fine-g	rained guidance				Goal-o	riented guidance		
	Client-Sel	Data-Pre	Label-Harm	Fed-Train	Overall	Client-Sel	Data-Pre	Label-Harm	Fed-Train	Overall
	$a_1, a_2, a_3$	$a_4$	$a_5$	$a_6, a_7$		$a_1, a_2, a_3$	$a_4$	$a_5$	$a_6, a_7$	
				Proprietar	y Models					
GPT-4.1	5/5, 5/5, 5/5	5/5	5/5	4/5, 5/5	97.14	5/5, 4/5, 5/5	5/5	5/5	4/5, 5/5	94.29
GPT-40	$5/5, \frac{2}{5}, \frac{5}{5}$	5/5	3/5	1/5, 5/5	74.29	$5/5, \frac{2}{5}, \frac{5}{5}$	5/5	3/5	1/5, 5/5	74.29
GPT-4	$5/5, \frac{4/5}{5}, \frac{5}{5}$	1/5	1/5	3/5, 5/5	68.57	$5/5, \frac{4/5}{5}, 5/5$	0/5	1/5	3/5, 5/5	65.71
GPT-4-Turbo	$5/5, \frac{4/5}{5}, \frac{5}{5}$	1/5	1/5	3/5, 5/5	68.57	$5/5, \frac{4/5}{5}, 5/5$	4/5	1/5	3/5, 5/5	77.14
GPT-4.1-mini	$5/5, \frac{4/5}{5}, \frac{5}{5}$	5/5	3/5	4/5, 5/5	88.57	$5/5, \frac{2}{5}, \frac{5}{5}$	3/5	4/5	3/5, 5/5	77.14
GPT-4o-mini	$5/5, \frac{3}{5}, \frac{3}{5}$	5/5	3/5	3/5, 5/5	77.14	5/5, 2/5, 3/5	5/5	1/5	4/5, 5/5	71.43
GPT-o4-mini	$5/5, \frac{4/5}{5}, \frac{5}{5}$	5/5	3/5	4/5, 5/5	88.57	$5/5, \frac{4/5}{5}, 5/5$	4/5	3/5	4/5, 5/5	85.71
GPT-o3-mini	$5/5, \frac{4}{5}, \frac{5}{5}$	1/5	1/5	4/5, 4/5	68.57	$5/5, \frac{4}{5}, \frac{5}{5}$	1/5	2/5	4/5, 5/5	74.29
GPT-3.5-Turbo	5/5, 0/5, 0/5	0/5	0/5	1/5, 3/5	25.71	5/5, 0/5, 0/5	2/5	0/5	1/5, 4/5	34.29
Claude-3-7	5/5, 3/5, 3/5	2/5	1/5	3/5, 3/5	57.14	5/5, 3/5, 3/5	2/5	1/5	3/5, 3/5	57.14
				Open-sour	ce Models	•				
				Huge M	Iodels					
DeepSeek-V3	5/5, 4/5, 5/5	5/5	5/5	4/5, 5/5	94.29	5/5, 3/5, 5/5	4/5	5/5	4/5, 5/5	88.57
DeepSeek-R1	5/5, 0/5, 5/5	0/5	0/5	0/5, 5/5	42.86	5/5, 0/5, 5/5	0/5	0/5	0/5, 5/5	42.86
Qwen3 235B	5/5, 0/5, 5/5	0/5	0/5	0/5, 5/5	42.86	5/5, 0/5, 5/5	0/5	0/5	0/5, 5/5	42.86
LLaMA-4 Maverick	5/5, 4/5, 5/5	5/5	4/5	4/5, 5/5	91.43	5/5, 4/5, 5/5	5/5	3/5	3/5, 5/5	85.71
LLaMA-4 Scout	$5/5, \frac{4/5}{5}, \frac{5}{5}$	5/5	4/5	2/5, 5/5	85.71	5/5, 1/5, 5/5	5/5	3/5	2/5, 5/5	74.28
				Large N	Aodels	•				
DeepSeek-R1-70B	5/5, 4/5, 5/5	3/5	1/5	2/5, 5/5	71.43	5/5, 4/5, 5/5	3/5	0/5	2/5, 5/5	68.57
LLaMA-3-70B	5/5, 4/5, 5/5	2/5	2/5	2/5, 5/5	71.43	5/5, 4/5, 5/5	2/5	2/5	2/5, 5/5	71.43
				Medium	Models					
Qwen QwQ 32B	5/5, 4/5, 5/5	4/5	4/5	4/5, 5/5	88.57	5/5, 4/5, 5/5	2/5	4/5	4/5, 5/5	82.86
Qwen3-30B	5/5, 1/5, 5/5	0/5	0/5	1/5, 5/5	48.57	5/5, 0/5, 5/5	0/5	0/5	1/5, 5/5	45.71
Gemma3-27B-instruct	5/5, 0/5, 0/5	0/5	0/5	0/5, 0/5	14.29	5/5, 0/5, 0/5	0/5	0/5	0/5, 0/5	14.29
				Small N	Iodels					
Gemma-2-9B	5/5, 0/5, 4/5	2/5	1/5	1/5, 5/5	51.43	5/5, 0/5, 4/5	1/5	1/5	1/5, 5/5	48.57
LLaMA-3-8B	5/5, 1/5, 4/5	4/5	2/5	1/5, 5/5	62.86	5/5, 0/5, 4/5	4/5	2/5	1/5, 5/5	60.00
Qwen-3-14B	5/5, 0/5, 0/5	0/5	0/5	0/5, 5/5	28.57	5/5, 0/5, 0/5	0/5	0/5	0/5, 4/5	25.71
Gemma3-12B-instruct	5/5, 0/5, 0/5	0/5	0/5	0/5, 0/5	14.29	5/5, 0/5, 0/5	0/5	0/5	0/5, 0/5	14.29

sits mid-pack (51.42 / 57.14). Stage-wise, stronger models are most reliable in **Client Selection** and **Federated Training**, frequently solving both steps (5/5, 5/5). In contrast, **Data Pre-processing** and especially **Label Harmonization** reveal brittleness for mid-tier and smaller models—many cells show 0/5 or 1/5, which depresses the Overall even when other steps succeed. Among open-source systems, **DeepSeek-V3** leads (85.71 / 80.00) with comparatively balanced performance, while **LLaMA-4 Maverick/Scout** form the next cluster (71.43–85.71 depending on guidance). Less-aligned or smaller open-source models (e.g., **DeepSeek-R1**, **Qwen3 235B**, **Gemma3-12B-instruct**) often fail early pipeline steps and remain near the floor (Overall  $\approx 14.29$ –42.86). Finally, **fine-grained guidance** consistently yields higher Overall scores than **goal-oriented guidance** for the same model, indicating

Model

GPT-4.1

507650775078

Table 15: Comparison of open-source and Proprietary LLM agents in X-Ray environment

Model		Fine-g	rained guidance				Goal-o	riented guidance		
	Client-Sel	Data-Pre	Label-Harm	Fed-Train	Overall	Client-Sel	Data-Pre	Label-Harm	Fed-Train	Overal
	$a_1, a_2, a_3$	$a_4$	$a_5$	$a_6, a_7$		$a_1, a_2, a_3$	$a_4$	$a_5$	$a_6, a_7$	
				Proprietar	y Models	•				
GPT-4.1	5/5, 5/5, 5/5	5/5	5/5	5/5, 5/5	100.00	5/5, 5/5, 5/5	5/5	5/5	5/5, 5/5	100.00
GPT-40	5/5, 1/5, 5/5	5/5	3/5	1/5, 5/5	71.43	5/5, 1/5, 5/5	5/5	2/5	1/5, 5/5	68.57
GPT-4	5/5, 5/5, 5/5	1/5	1/5	3/5, 5/5	71.43	5/5, <mark>4/5</mark> , 5/5	0/5	1/5	3/5, 5/5	65.71
GPT-4-Turbo	5/5, 5/5, 5/5	1/5	1/5	3/5, 5/5	71.43	$5/5, \frac{4/5}{5}, 5/5$	4/5	1/5	3/5, 5/5	77.14
GPT-4.1-mini	$5/5, \frac{4/5}{5}, \frac{5}{5}$	5/5	3/5	4/5, 5/5	88.57	$5/5, \frac{2}{5}, \frac{5}{5}$	3/5	4/5	3/5, 5/5	77.14
GPT-4o-mini	$5/5, \frac{3}{5}, \frac{3}{5}$	5/5	3/5	3/5, 5/5	77.14	5/5, 2/5, 3/5	5/5	1/5	5/5, 5/5	74.29
GPT-o4-mini	5/5, 5/5, 5/5	5/5	3/5	4/5, 5/5	91.43	$5/5, \frac{4/5}{5}, 5/5$	4/5	3/5	4/5, 5/5	85.71
GPT-o3-mini	5/5, 5/5, 5/5	1/5	1/5	4/5, 4/5	71.42	$5/5, \frac{4/5}{5}, 5/5$	1/5	2/5	4/5, 5/5	74.29
GPT-3.5-Turbo	5/5, 0/5, 0/5	0/5	0/5	1/5, 3/5	25.71	5/5, 0/5, 0/5	2/5	0/5	1/5, 4/5	34.29
Claude-3-7	$5/5, \frac{3}{5}, \frac{3}{5}$	2/5	1/5	3/5, 3/5	57.14	5/5, 3/5, 3/5	2/5	1/5	3/5, 3/5	57.14
				Open-sour	ce Models	•				
				Huge M	Iodels					
DeepSeek-V3	5/5, 4/5, 5/5	5/5	5/5	4/5, 5/5	94.29	5/5, 3/5, 5/5	4/5	5/5	4/5, 5/5	88.57
DeepSeek-R1	5/5, 0/5, 5/5	0/5	0/5	0/5, 5/5	42.86	5/5, 0/5, 5/5	0/5	0/5	0/5, 5/5	42.86
Qwen3 235B	5/5, 0/5, 5/5	0/5	0/5	0/5, 5/5	42.86	5/5, 0/5, 5/5	0/5	0/5	0/5, 5/5	42.86
LLaMA-4 Maverick	5/5, 5/5, 5/5	5/5	4/5	4/5, 5/5	94.29	5/5, 4/5, 5/5	5/5	3/5	3/5, 5/5	85.71
LLaMA-4 Scout	5/5, 4/5, 5/5	5/5	4/5	2/5, 5/5	85.71	5/5, 1/5, 5/5	5/5	3/5	2/5, 5/5	74.28
				Large N	Aodels					
DeepSeek-R1-70B	5/5, 5/5, 5/5	3/5	1/5	2/5, 5/5	74.28	5/5, 4/5, 5/5	3/5	0/5	2/5, 5/5	68.57
LLaMA-3-70B	5/5, 4/5, 5/5	2/5	2/5	2/5, 5/5	71.43	5/5, 4/5, 5/5	2/5	2/5	2/5, 5/5	71.43
				Medium	Models					
Qwen QwQ 32B	5/5, 4/5, 5/5	4/5	4/5	4/5, 5/5	88.57	5/5, 4/5, 5/5	2/5	4/5	4/5, 5/5	82.86
Qwen3-30B	5/5, 1/5, 5/5	0/5	0/5	1/5, 5/5	48.57	5/5, 0/5, 5/5	0/5	0/5	1/5, 5/5	45.71
Gemma3-27B-instruct	5/5, 0/5, 0/5	0/5	0/5	0/5, 0/5	14.29	5/5, 0/5, 0/5	0/5	0/5	0/5, 0/5	14.29
			-	Small N				-		
Gemma-2-9B	5/5, 0/5, 4/5	2/5	1/5	1/5, 5/5	51.43	5/5, 0/5, 4/5	1/5	1/5	1/5, 5/5	48.57
LLaMA-3-8B	5/5, 1/5, 4/5	4/5	2/5	1/5, 5/5	62.86	5/5, 0/5, 4/5	4/5	2/5	1/5, 5/5	60.00
Owen-3-14B	5/5, 0/5, 0/5	0/5	0/5	0/5, 5/5	28.57	5/5, 0/5, 0/5	0/5	0/5	0/5, 4/5	25.71
Gemma3-12B-instruct	5/5, 0/5, 0/5	0/5	0/5	0/5, 0/5	14.29	5/5, 0/5, 0/5	0/5	0/5	0/5, 0/5	14.29

5099 5100 5101

5102

5103

5104

Table 16: Comparison of average time taken by each agent to solve respective tasks (in seconds) using different LLMs.

Proprietary Models

 $S_1$ 

1.8

 $C_1$ 

64.8

 $S_2$ 

55.5

 $C_2$ 

302.4

 $C_3$ 

130.7

 $S_3$ 

54.1

 $S_4$ 

18.8

5105	
5106	
5107	
5108	
5109	
5110	
5111	
5112	

5117

5118

5126

5127

5128

5129

GPT-40 201.0 1.0 311.3 53.5 58.7 30.9 9.6 GPT-4 2.9 235.2 87.4 172.3 615.5 243.7 31.5 GPT-4-Turbo 1.8 81.2 54.8 259.9 266.7 76.7 16.6 29.9 GPT-4.1-mini 1.0 78.1 183.6 161.5 69.8 9.7 GPT-4o-mini 1.0 73.4 29.3 370.7 292.1 77.0 10.7 127.4 404.6 503.9 168.2 GPT-o4-mini 4.2 164.8 42.5 GPT-o3-mini 4.9 156.2 145.9 177.7 412.3 172.1 44.5 GPT-3.5-Turbo 1.1 51.1 32.8 163.9 199.9 52.7 9.9 Claude-3-7 3.9 231.6 115.5 414.0 457.7 203.0 37.2 Open-source Huge Models 197.2 DeepSeek-V3 4.4 169.3 131.2 554.1 461.5 44.1 162.9 328.0 DeepSeek-R1 8.1 242.1 567.1 134.2 77.4 Qwen3 235B 11.0 180.3 328.8 642.8 440.7 168.9 108.3 Large Models LLaMA-4 Maverick 282.7 118.3 1.2 98.9 124.2 13.6 37.2 LLaMA-4 Scout 2.3 105.3 69.1 172.0 300.4 103.6 24.6 DeepSeek-R1-70B 1.5 96.0 44.4 168.0 312.5 99.0 15.2 LLaMA-3-70B 1.5 93.2 45.4 193.7 257.4 76.3 15.0 Medium Models Qwen QwQ 32B 0.8 77.2 24.0 186.0 253.1 74.2 8.4 Qwen3-30B 2.3 73.9 68.2 164.4 297.4 83.7 24.6 Gemma3-27B-instruct 2.8 140.9 82.8 297.4 535.0 133.9 26.3 Small Models Gemma-2-9B 0.5 116.9 15.4 105.3 283.1 111.2 5.1 LLaMA-3-8B 1.4 155.3 573.4 144.9 13.5 42.6 212.1 Qwen-3-14B 165.1 123.6 520.0 357.3 176.5 45.0 Gemma3-12B-instruct 184.7 94.4 400.2 487.9 195.1 33.5

that explicit, stepwise supervision helps agents navigate the multi-stage federated learning workflow in Dermatology.

Table 11 compares open-source and proprietary LLM agents in the **Ultrasound environment** under two guidance paradigms: fine-grained guidance, where each subtask is explicitly defined (e.g., client selection, data preprocessing), and goal-oriented guidance, where the model is only given the overall objective. Each model's performance is evaluated on four core subtasks, and the final column reports the average normalized score. GPT-4.1 achieves the highest performance (94.29) under both guidance types, demonstrating strong generalization across all subtasks. Smaller models like Gemma3-12B-instruct significantly underperform (14.29), especially when tasks require coherent execution across multiple stages. Open-source models such as DeepSeek-V3 and LLaMA-4 Maverick exhibit competitive performance with proprietary models under fine-grained prompts but see mild drops in goal-oriented execution.

Table 12 evaluates LLM agents' ability to generate independent scripts for each subtask in the Ultrasound environment without any explicit or implicit guidance on the workflow or availability of tools. It is expected to write scripts for completing the tasks. This setup is more challenging than the previous table. A uniform drop in performance is observed across all models, regardless of type or size. Most top proprietary models, such as the GPT-4 series, drop to a common score of 42.86, indicating reliance on guided execution for complex task planning. Mid and small-scale models like Claude-3-7, Gemma3-12B, and Qwen-3-14B perform poorly, with scores as low as 14.29 to 31.43, demonstrating that few-shot context is crucial for robust task understanding.

Table 13 presents evaluation in the **MRI environment**, following the same structure. GPT-4.1 again leads with a perfect score (100.00) under both guidance types. A general trend of better performance under fine-grained guidance than goal-oriented guidance is maintained across most models. Open-source large-scale models such as DeepSeek-V3 and LLaMA-4 Maverick narrow the performance gap significantly, achieving scores above 85 under fine-grained guidance. Models with weaker subtask handling like Gemma3-12B-instruct remain consistently poor performers, struggling to follow multi-step instructions even in highly structured MRI tasks.

Table 14 presents the evaluation of proprietary and open-source LLM agents deployed in the **Fundus environment**. The table demonstrates that proprietary models such as GPT-4.1 and GPT-4.1-mini achieve near-perfect scores across both guidance styles, indicating robust task execution capabilities. GPT-4.1 achieves the highest overall fine-grained score (97.14) and maintains a strong goal-oriented score (94.29), suggesting high generalization capacity even with minimal instruction. In contrast, smaller models like GPT-3.5-Turbo and Gemma3-12B-instruct exhibit major limitations, particularly under goal-oriented prompting, often failing multiple subtasks and scoring below 35.

Among open-source models, DeepSeek-V3 and LLaMA-4 Maverick lead performance under both guidance types, with fine-grained scores above 90 and goal-oriented scores above 85. These models close the gap with top proprietary agents, showcasing the progress of the open-source ecosystem. However, performance drops significantly in lightweight models such as Qwen-3-14B and Gemma3-12B-instruct, which perform well only on the most basic subtasks and fail to coordinate complex operations under goal-driven conditions.

Table 15 compares open-source and proprietary LLM agents in the **XRay environment** Proprietary models continue to dominate the list. **GPT-4.1** reaches ceiling performance with consistent 5/5 across all sub-tasks and the best **Overall** scores (**100.00** under both fine-grained and goal-oriented guidance). A strong second tier includes **GPT-04-mini** (91.43 / 85.71) and **GPT-4.1-mini** (88.57 / 77.14), followed by **GPT-40-mini** (77.14 / 74.29), **GPT-4-Turbo** (71.43 / 77.14), **GPT-4** (71.43 / 65.71), **GPT-40** (71.43 / 68.57), and **GPT-03-mini** (71.42 / 74.29). **Claude-3-7** is mid-pack (57.14 / 57.14), while **GPT-3.5-Turbo** trails (25.71 / 34.29). Open-source leaders close the gap in X-Ray: **DeepSeek-V3** (94.29 / 88.57) and **LLaMA-4 Maverick** (94.29 / 85.71) approach the top proprietary tier, with **LLaMA-4 Scout** (85.71 / 74.28) and **Qwen QwQ 32B** (88.57 / 82.86) forming a competitive middle. In contrast, **DeepSeek-R1** and **Qwen3 235B** remain low (42.86 / 42.86), and **Gemma3-27B-instruct** bottoms out (14.29 / 14.29). Among smaller models, **LLaMA-3-8B** (62.86 / 60.00) outperforms **Gemma-2-9B** (51.43 / 48.57) and **Qwen-3-14B** (28.57 / 25.71). Stage-wise patterns mirror other modalities: top systems reliably solve **Client Selection** and **Federated Training** (5/5, 5/5), whereas weaker models falter on **Data-Pre** and especially **Label-Harm** (frequent 0/5 or 1/5), depressing Overall scores. Across models, **fine-grained guidance** usually yields slightly higher

Overall than **goal-oriented** (with a few reversals, e.g., GPT-4-Turbo and GPT-o3-mini), indicating the benefit of explicit stepwise supervision for end-to-end X-Ray workflows.

Overall, all the tables reveal two key insights: (1) proprietary models consistently outperform opensource ones across both settings, (2) fine-grained prompting benefits all models but especially weaker ones.

518951905191

5187

5188

# D.2 DISCUSSION ON TIME-EFFICIENCY

519251935194

5195

5196

5197

5198

5199

5200

5201

5202

5203

5204

Table 16 compares the average time taken (in seconds) by each agent across the seven subtasks (S1–S4, C1–C3) in the pipeline. GPT-4.1 is among the fastest overall, particularly in inference-heavy subtasks like S1 and S4. Open-source models such as Qwen3-235B and DeepSeek-R1 exhibit significantly higher latency, especially in complex subtasks like C2, where times range from approximately 550 to 640 seconds. Lightweight models such as Qwen QwQ 32B and Gemma-2-9B complete tasks much faster but at the cost of performance, as seen in the other tables. This table complements the prior performance evaluations by highlighting the efficiency–performance tradeoff, which is critical for real-world federated deployments.

We have conducted a comparison of **time–efficiency vs. performance** for each agent role (S1, C1, S2, C2, C3, S3, S4) across model families. Overall, we observe the following:

- C2 (data prep) and C3 (label harmonization) dominate wall-clock time for almost every model. S1/S2/S3/S4 are comparatively light; differences here are smaller and rarely drive total runtime. The best choices balance high stage success and short C2/C3 times. Agent-wise takeaways (cross-model):
- (i) S1 (server task extraction/broadcast). Times are uniformly small. Fastest include Gemma-2-9B (0.5s), QwQ-32B (0.8s), GPT-4o/4o-mini/4.1-mini (1.0s). This stage won't bottleneck overall runtime, so one should prefer models with higher downstream success rather than saving fractions of a second here.
- 5212 (ii) C1 (client selection). A moderate cost stage. GPT-3.5-Turbo (51.1s) and QwQ-32B (77.2s)
  5213 are among the fastest; GPT-4.1 (64.8s) and GPT-4.1-mini (78.1s) are also efficient. Very large
  5214 open-source models (e.g., Qwen3-235B 180s) are slower without clear gains.
- 5215 5216 5217 (iii) S2 (approval/coordination). Also light in terms of time complexity. Gemma-2-9B (15.4s), QwQ-32B (24.0s), GPT-4o/4o-mini/4.1-mini (29–30s) are quickest.
- (iv) C2 (data prep / cleaning). One of the two big time sinks. Fastest include Gemma-2-9B (105s) and LLaMA-4 Maverick (124s); GPT-3.5 (164s), Qwen3-30B (164s), LLaMA-3-70B (194s), QwQ-32B (186s) are solid. GPT-4.1 (302s) and huge open-source (DeepSeek-V3 554s; Qwen3-235B 643s) are slower. LLaMA-4 Maverick and QwQ-32B are strong Pareto options (good success, reasonable C2 time).
- (v) C3 (label harmonization). The other major time sink and the hardest stage. Standout: GPT-4.1 (131s)—both fast and high success. Next tier includes QwQ-32B (253s) and LLaMA-3-70B (257s), which are respectable; GPT-4o (201s) is faster than many but weaker on Label Harmonization accuracy. GPT-4 (616s) and huge open-source (e.g., DeepSeek-V3 462s) are slow here.
- 5227 (vi) S3 (algorithm selection). Lightweight. GPT-3.5 (52.7s), GPT-40 (53.5s), GPT-4.1 (54.1s) are quickest; QwQ-32B (74s) is not far behind. This stage rarely determines end-to-end time.
- (vii) S4 (training trigger/monitor). Very small across models. Gemma-2-9B (5.1s) is fastest; QwQ-32B (8.4s), GPT-4o/4.1-mini (9–10s) are close. Not a driver of total latency.
- We summarize the overall recommendations based on our experiments below:
- Best overall (reliability time): GPT-4.1 exceptional C3 time (130.7s) and top success. Best open-source Pareto: Qwen QwQ 32B C2 186s / C3 253s with strong success; or LLaMA-4 Maverick if faster C2 is needed (124s). Budget/latency-focused orchestration: GPT-4.1-mini or GPT-4o-mini (But need to keep in mind the success drop on C3). Avoid very large open-source for time-critical runs unless one specifically needs open-source + the higher success of DeepSeek-V3

(and can pay the time cost).

 D.3 DISCUSSION ON CLIENT SELECTION, REASONING VS NON-REASONING MODELS AND FAILURE MODES:

**Qualitative analysis of client selection across modalities.** Figures 9–33 present the qualitative agentic performance in the *Client Selection* stage under three clinical modalities—**skin cancer** (dermatology), **histopathology** (breast cancer detection), and **X-Ray** (pneumonia detection)—and contrast *non-thinking/reasoning* and *thinking/reasoning* LLM agents. Across all settings, the figures illustrate *when/how* the server approves or declines prospective clients for federated training. For non-thinking agents (e.g., Figs. 9-11; 20-24; 27-29), the selection is typically concise: the model applies eligibility checks and emits a binary decision (approve/decline) with minimal justification. This often highlights crisp gating on dataset relevance to the target task, basic quality constraints, and coarse client readiness.

**Impact of using thinking/reasoning agents** For **thinking/reasoning** agents (e.g., Figs.12-16, 25-26, 30-33), the server-facing rationale becomes more elaborate. These figures show richer criteria—such as finer judgements about class balance, labeling consistency, or potential contribution to global convergence—before issuing approve/decline decisions. While this often results in clearer, auditable justifications, it can also introduce overhead: Fig. 15 exemplifies *overthinking*, where extended deliberation adds verbosity without changing the final decision. Taken together, the sequences suggest a trade-off: explicit reasoning improves transparency and sometimes catches subtle issues, but may reduce efficiency and occasionally distract from the primary selection objective.

**Failure modes: hallucination and task drift.** Figures 17-18 document characteristic **hallucinations** during client selection with skin cancer datasets. In one case, the model drifts to an *irrelevant task*, attempting to solve something other than client eligibility; in another, it answers in *Russian*, a response channel misaligned with the specified instruction and downstream system expectations. Such behaviors indicate vulnerability to prompt misinterpretation and context leakage even at the pre-training data curation stage. The remaining thinking-model traces (e.g., Fig 19) demonstrate successful recoveries where the agent returns to the approval/decline protocol after structured reasoning.

Consistency across datasets and tasks. Across histopathology (breast cancer) and X-Ray (pneumonia) examples, we observe the same qualitative patterns: non-thinking models provide fast, rule-like triage; thinking models surface nuanced justifications but are susceptible to verbosity and occasional digressions. The figures collectively map the decision boundary between acceptance and rejection—anchored in dataset/task alignment and basic quality signals—while exposing two practical risks for agentic selection: (i) *over-elaboration*, which inflates latency without added value, and (ii) *hallucination/task drift*, which can misroute the pipeline if not caught by server-side validation. These qualitative insights complement the quantitative tables, clarifying *how* different prompting regimes lead to the observed approval/decline outcomes in federated client onboarding.

## E PRIVACY ANALYSIS OF HARMONIZED LABELS AND METADATA

Our benchmark's contribution lies in system-level automation and task performance evaluation, not in proving privacy guarantees. However, since FedAgentBench utilizes harmonized labels and some form of metadata exchange across clients, below, we rigorously analyze the privacy implications of these harmonized labels and transmitted metadata.

### E.1 MUTUAL INFORMATION ANALYSIS

Let X be the original dataset at a client, and M=f(X) represent the harmonized labels and metadata extracted from the local dataset X, where f includes only non-identifying structural information and label taxonomies. In practice, f is a projection or generalization map (e.g., mapping "melanoma" and "BCC" both to "malignant"). To quantify potential data tracing risk, we use Mutual Information (MI):

$$MI(X; M) = H(X) - H(X|M)$$

where H is the Shannon entropy.

To guarantee minimal traceability:

5292

5293 5294

5343

5344 5345

 $MI(X; M) < \delta, \quad \delta \to 0$ 5295 5296 **Proof:** 5298 • By designing the function f (harmonization process), we ensure maximal entropy in 5299 H(X|M). 5300 • Assume f maps multiple distinct datasets  $X_i \in \mathcal{X}$  to a similar M. Let  $|\mathcal{X}| \gg |\mathcal{M}|$ . This 5301 introduces significant ambiguity, thus: 5302 5303  $H(X|M) \approx H(X)$ 5304 5305 which implies: 5306  $MI(X; M) \approx 0$ 5307 5308 Hence, tracing original data through metadata is theoretically negligible. 5309 5310 E.2 DIFFERENTIAL PRIVACY (DP) PROOF 5311 5312 We formalize DP guarantees. 5313 Let  $\mathcal{A}$  be a randomized mechanism (e.g., gradient updates with Gaussian noise), and D, D' two 5314 neighboring datasets differing by one record. A satisfies  $(\epsilon, \delta)$ -DP if: 5315 5316  $\Pr(\mathcal{A}(D) \in S) \le e^{\epsilon} \Pr(\mathcal{A}(D') \in S) + \delta, \quad \forall S \subseteq \operatorname{Range}(\mathcal{A})$ 5317 5318 **Proof Outline:** 5319 5320 • If Gaussian noise  $\mathcal{N}(0, \sigma^2)$  is added to updates during training: 5321 5322  $\mathcal{A}(D) = \nabla f(D) + \mathcal{N}(0, \sigma^2)$ 5323 • For mechanism sensitivity  $\Delta$ , noise variance  $\sigma^2$  satisfies: 5325 5326  $\sigma \ge \frac{\Delta\sqrt{2\ln(1.25/\delta)}}{\epsilon}$ 5327 5328 thus rigorously satisfying DP conditions. 5329 5330 E.3 K-ANONYMITY ANALYSIS 5332 Let  $\mathcal{C}$  be the set of clients. Metadata M ensures k-anonymity if each metadata description transmitted from a client  $m \in M$  is generalized such that it matches at least k indistinguishable clients: 5334 5335 5336  $\forall m \in M, \quad |\{c \in \mathcal{C} : f(X_c) = m\}| \ge k$ 5337 5338 **Proof:** 5339 5340 • By metadata generalization, f is designed such that distinct datasets yield identical or highly similar metadata. 5342

 $|\{c \in \mathcal{C} : f(X_c) = m\}| \ge k$ 

• Given  $|\mathcal{C}| \gg k$ , the number of clients per metadata class is enforced:

thus rigorously satisfying k-anonymity.

#### E.4 PRIVACY-UTILITY TRADE-OFF

Define utility U as the expected accuracy of the trained model, and privacy loss  $\epsilon$  as above. We have:

$$U(\epsilon) = \mathbb{E}[\operatorname{Acc}(M_{\epsilon})] \quad \text{with} \quad \frac{dU}{d\epsilon} > 0$$

implying greater privacy (lower  $\epsilon$ ) results in lower accuracy.

#### **Theoretical Bound:**

 Utility degradation due to noise addition (DP) or generalization (k-anonymity) is bounded by:

$$|U(\epsilon) - U(0)| \le O\left(\frac{1}{\epsilon}\right)$$

This rigorous mathematical analysis demonstrates that harmonized labels and metadata transmission in **FedAgentBench** can achieve stringent privacy guarantees with negligible traceability risks, aligning with formal **differential privacy** and **k-anonymity** standards.

## F BROADER SOCIAL IMPACT

Positive Societal Impacts: FedAgentBench offers a significant advancement toward democratizing access to federated learning (FL) in healthcare by enabling agentic automation of otherwise labor-intensive tasks such as client selection, data cleaning, label harmonization, and FL algorithm configuration. This shift can particularly benefit healthcare institutions in low- and middle-income countries (LMICs) and rural areas, where hiring skilled data engineers is not feasible. By reducing the human expertise required to deploy and manage FL workflows, our benchmark promotes more equitable participation in collaborative medical AI initiatives. Furthermore, the open-source nature of our framework ensures broad accessibility, fostering reproducibility, community contributions, and educational use.

**Negative Societal Impacts:** The automation of critical decision-making steps in FL workflows, especially in high-stakes healthcare applications, raises concerns about over-reliance on large language models (LLMs) that may hallucinate or misinterpret complex clinical instructions. Incorrect data preprocessing or label harmonization could propagate errors into downstream training, leading to biased or unreliable models. Additionally, while the framework supports privacy-preserving learning, if misused without proper oversight, it may inadvertently reinforce data quality issues or amplify existing healthcare disparities. Careful human oversight and robust safety mechanisms must accompany such autonomous systems to ensure responsible deployment in real-world medical settings.

#### G LLM USAGE:

LLMs (GPT-40 and GPT-5) were sometimes used for improving grammar and wording.

```
5402
5403
5404
5405
5406
5407
      Running AgentPipelines/run.py...
5408
      Human Requirements: I want to train skin cancer classification model. Only use datase
5409
      ts that have skin cancer classes
5410
5411
      communication_1_content Task: Skin cancer classification
5412
      Modality: Image
5413
5414
      Clients, please respond with the name of the selected dataset that contains skin cance
      r classes suitable for this classification task.
5415
      Selector Content: **Dataset Name** : DDI_skin_dataset
5416
5417
      5418
      communication_2_content Approved. Prepare for training
5419
      Selector Content: **Dataset Name** : skinL2 dataset
5420
5421
      5422
      communication_2_content Approved. Prepare for training
5423
      Selector Content: **Dataset Name** : ISIC 2020
5424
      5425
5426
      communication_2_content Approved. Prepare for training
5427
      Selector Content: **Dataset Name** : PH2Dataset
5428
      5429
5430
      communication_2_content Approved. Prepare for training
      Selector Content: **Dataset Name** : Dermnet
5431
5432
5433
      5434
      communication_2_content Approved. Prepare for training
5435
      Selector Content: **Dataset Name** : ISIC2018_HAM10000
5436
      5437
5438
      communication_2_content Approved. Prepare for training
5439
      Selector Content: no dataset: The provided dataset "skin_disease_classification_kaggle
      " contains classes for skin diseases such as acne, eye bags, and redness, but it does
not include any skin cancer classes. Therefore, it is not suitable for the skin cancer
5440
5441
      classification task.
5442
      5443
      communication_2_content Client not needed for the task
5444
```

Figure 9: Client Selection with **skin cancer datasets** for **non-thinking/reasoning** models. It shows when/how the server approves clients for training whereas when/how it declines.

```
5455
5456
5457
5458
      Selector Content: no dataset. The provided dataset "Monkeypox_Skin_Image_Dataset"
     ains skin disease classes such as Chickenpox, Measles, Monkeypox, and Normal, but it does not have skin cancer classes suitable for the skin cancer classification task.
5459
5460
      5461
5462
     communication_2_content Client not needed for the task
      Selector Content: **Dataset Name** : Dermquest
5463
5464
      5465
5466
      communication_2_content Approved. Prepare for training
      Selector Content: **Dataset Name** : augmented_skin_condition_dataset_kaggle
5467
      This dataset contains a skin cancer class named "Carcinoma" suitable for the skin canc
5468
      er classification task.
5469
      5470
5471
      communication_2_content Approved. Prepare for training
5472
     Selector Content: **Dataset Name** : fitzpatrick17k
5473
      5474
5475
      communication_2_content Approved. Prepare for training
5476
      Selector Content: **Dataset Name** : Dermis
5477
      5478
5479
      communication_2_content Approved. Prepare for training
     Selector Content: **Dataset Name** : ISIC_2017
5480
5481
      5482
5483
     communication_2_content Approved. Prepare for training
     Selector Content: **Dataset Name** : Mednode
5484
5485
      5486
      communication_2_content Approved. Prepare for training
5487
      Selector Content: **Dataset Name** : ISIC 2024
5488
5489
      5490
      communication_2_content Approved. Prepare for training
5491
      Selector Content: **Dataset Name** : Derma7P
5492
5493
      5494
      communication_2_content Approved. Prepare for training
5495
      Selector Content: **Dataset Name** : skin_disease_kaggle_dataset
5496
      This dataset contains skin cancer classes such as Basal Cell Carcinoma (BCC) and Melan
      oma, suitable for skin cancer classification.
5497
5498
5499
      5500
      communication_2_content Approved. Prepare for training
5501
```

Figure 10: Client Selection with **skin cancer datasets** for **non-thinking/reasoning** models. It shows when/how the server approves clients for training whereas when/how it declines.

```
5510
5511
5512
5513
      Selector Content: **Dataset Name** : ISIC 2019
5514
      5515
      communication_2_content Approved. Prepare for training
5516
      Selector Content: no dataset
5517
5518
      None of the available datasets include skin cancer classes (e.g., melanoma, basal cel
      l carcinoma, squamous cell carcinoma). Therefore, there is no suitable dataset for th
5519
      e given task.
5520
      5521
5522
      communication_2_content Client not needed for the task
      Selector Content: **Dataset Name** : PAD UFES 20
5523
5524
      5525
      communication_2_content Approved. Prepare for training
5526
      Selector Content: **Dataset Name** : ISIC_2016
5528
      5529
      communication_2_content Approved. Prepare for training
5530
      Selector Content: **Dataset Name** : skin_disease_3_class
5531
      This dataset includes the class "basal cell carcinoma," which is a type of skin cance
5532
5533
      5534
5535
      communication_2_content Approved. Prepare for training
5536
      Selector Content: no dataset
5537
      The only available dataset ("Dermatology_tabular dataset") is non-image (tabular) and
5538
      does not include dermatology images or specific skin cancer classes, so it is not su
5539
      itable for a skin cancer classification model.
5540
      5541
5542
      communication 2 content Client not needed for the task
5543
      Selector Content: no dataset
5544
      The provided dataset ("skin-infection-disease-datasaet") covers only infectious skin
5545
      conditions (cellulitis, impetigo, athlete's foot, nail fungus, ringworm, cutaneous la
5546
      rva migrans, chickenpox, shingles) and does not include any skin cancer classes.
      5547
5548
      communication_2_content Client not needed for the task
5549
      Selector Content: no dataset
5550
      None of the available datasets include skin cancer classes; the only dataset provided
5551
      (Monkeypox_Skin_Image_Dataset) contains Chickenpox, Measles, Monkeypox, and Normal c
5552
      ategories, which do not cover any form of skin cancer.
5553
      5554
```

Figure 11: Client Selection with **skin cancer datasets** for **non-thinking/reasoning** models. It shows when/how the server approves clients for training whereas when/how it declines.

5611

```
5564
5565
      Selector Content: **Dataset Name** : skin_disease_3_class
5566
5567
       5568
      communication_2_content Since the client has returned a dataset named "skin_disease
5569
      _3_class", it indicates that they have provided a dataset.
5570
5571
      Therefore, the response should be: "Approved. Prepare for training"
      Selector Content: The dataset description indicates that the dataset is a tabular d
5572
      ataset containing clinical features for diagnosing skin diseases, which doesn't mat
5573
      ch the modality of Dermatology Images. Therefore, this dataset is not suitable for
5574
      the task of Skin Cancer Detection using Dermatology Images.
5575
      Since I have to respond with the name of the selected dataset, I will check if ther
5576
      e's any other dataset that matches the task and modality.
5577
      **No suitable dataset found** : The dataset description doesn't match the modality
5579
      of Dermatology Images.
5580
5581
      5583
      communication_2_content Since the client has returned a dataset, although it's not
      suitable for the task, the client has still responded with a dataset. However, the
5584
      response indicates that the dataset is not suitable.
5585
5586
      The correct interpretation of the task is to respond based on whether the client ha
5587
       s returned any dataset. Since the client has returned a dataset (even though it's n
      ot suitable), the response should be based on the presence of a dataset.
5588
5589
      Therefore, the response should be: "Approved. Prepare for training".
5590
      Selector Content: **Dataset Name** : skin-infection-disease-datasaet
5591
      5592
5593
      communication_2_content Since the client has returned a dataset named "skin-infecti
5594
      on-disease-datasaet", I will return the approval message.
5595
       "Approved. Prepare for training"
      Selector Content: **Dataset Name** : Monkeypox_Skin_Image_Dataset
5597
5598
       5599
      communication 2 content Since the client has returned a dataset, **Dataset Name**
5600
       Monkeypox_Skin_Image_Dataset, we can infer that the client has indeed provided a d
       ataset for the task.
5602
      Therefore, the response should be: "Approved. Prepare for training".
      Selector Content: **Dataset Name** : Dermquest
5604
5605
       5606
5607
      communication_2_content Since the client has returned a dataset named "Dermquest",
      the response should be: "Approved. Prepare for training".
5608
      Therefore, the output is: Approved. Prepare for training
5610
```

Figure 12: Client Selection with **skin cancer datasets** for **thinking/reasoning** models. It shows when/how the server approves clients for training whereas when/how it declines.

```
5619
 5620
5621
                                                                   tet me check the intervention of the control of the
5622
                                                     o, the user mentioned "Derma?" which is a dataset name, since they specified one dataset here, that mest be condition of having at least one dataset. Therefore, the correct response should be "Appro dataset and the condition of having at least one dataset, make represent, which it is. Airlight, that".

A Prepare for training: I need to make sure not to add anything alse. Just the exact nessage as per the instructions. Let me double—check that the dataset name is present, which it is. Airlight, that".
 5623
5624
5625
5626
                                                          y, let me see. The user provided a response where the client mentioned a dataset called "skin_disease_kaggle_dataset". The instructions say that if the client returns one or we and prepare for training. If not, send the other message.
5627
                                                        first, check if there's a dataset name listed. Here, the dataset name is clearly given. The format specifies to look for **Oataset Name**. The user input has exactly that. So the condition is met. Theoretic response is "Approved. Prepare for training".
5628
                                                        need to make sure there's no trick here. Maybe check if the dataset is valid? But the problem doesn't mention validating the da
is present, the answer is straightforward. No need to consider anything else. Just follow the first instruction.
 5629
5630
                                                        proved. Prepare for training
deter Content: "Postages Hamme": no dataset
eaverlable "scin_dataset" includes skin condition categories like 'GROWTH_OM_MOLE' but does not explicitly include **skin cancer subtypes** (e.g., melanona, basal cell carcino
etition of skin cancer. Its focus on broader skin issues (e.g., acme, pigmentation problems, hair loss) makes it unsuitable for the specific task of cancer classification.
 5631
5632
5633
                                                                       t me parse this step by step. The user provided a response from a client agent in a federated learning setup. The client mentioned that their dataset, "scin_dataset," includes categories like 'NOLE' but doesn't explicitly cover specific skin cancer subtypes such as melanoma or basal cell carcinoma. They also noted the dataset focuses more on broader skin issues like acne or hair loss tungstable for cancer classification.
                                                      ne user's instructions are clear: if the client returns one or more applicable datasets, the server should respond with "Approved. Prepare for training." If no suitable dataset is prahould be "Client not needed for the task."
5634
 5635
                                                      merefore, since the provided dataset doesn't fulfill the necessary criteria, the correct response per instructions is to inform the client they aren't meeded. The user's guidelines specify that if no ilcable dataset is returned, the nessage should be 'client not needed for the task. Even though the client did mention a dataset, it's not suitable, so it counts as not providing a needed dataset.
5637
 5639
```

Figure 13: Client Selection with **skin cancer datasets** for **thinking/reasoning** models. It shows when/how the server approves clients for training whereas when/how it declines.

```
5649
5650
5651
                         re, if "no dataset" is the dataset name they are returning, then they have "no dataset" as their dataset name, which would mean they have no suitable dataset. Thus "no datase dataset, so server's response is "client not needed".
5652
                         tively, their having a dataset: "The available dataset (Dernatology...) exists, but it's not suitable. But in their response, they also have "no dataset" as the dataset name, which
5653
                     en that the first line says: Dataset Name: no dataset. The rest explains why
5654
                      s, in this case, if the client's returned dataset name is "no dataset", then the count is zero, so return "Client not needed"
5655
                         stively, the client actually has a dataset (Dermatology...) and the "no dataset" may be an erro
                     nce the user's problem says "available dataset... contains only tabular...", which is a dataset. So the client does have a dataset, but it's unsuitable.
5656
5657
5658
                         stively, perhaps "returned one or more datasets" refers to datasets that comply, and if they don't comply, they aren't counted
5659
5660
                     n. I think there's ambiguity, but based on the user's initial statement, if the client provides any dataset, even unsuitable, the server must resping".
5661
5662
5663
```

Figure 14: Client Selection with **skin cancer datasets** for **thinking/reasoning models**. It shows when/how the server approves clients for training whereas when/how it declines.

```
communication 2_content

Communication 2_conte
```

Figure 15: Client Selection with **skin cancer datasets** for **thinking/reasoning models**. It shows when/how the server approves clients for training whereas when/how it declines. **This is an instance of typical overthinking** 

```
Relector Content: "Potaseet Name": In detaset is, disease, classification kapple" contains classes like 'scne', 'eye bags', and 'redness', which are skin conditions rather than skin cancer subtypes (e.g., nelanoma, bas al cell carcinoma). Lack of natignancy-related classes nakes it mismatched to skin cancer detection.

Communication_Content

Content

Conten
```

Figure 16: Client Selection with **skin cancer datasets** for **thinking/reasoning** models. It shows when/how the server approves clients for training whereas when/how it declines.

Figure 17: Typical example of **hallucination** in Client Selection with **skin cancer datasets** for **thinking/reasoning** models. It shows the model is trying to solve irrelevant task.

```
Rencrymu nawawa moryr сильно влиять на работу, провяваесь различными способами. Вот некоторые из них:

**1. Снижение производительности и концентрация:**

**1. Снижение производительности и концентрация:**

***Трудности с концентрация:** бо время приступа паники сложо сосредоточиться на задачах, что приводит к снижения производительности. Мысли могут быть хаотичными, а внимание рассениями.

***Трудности с концентрация:** бо время приступа паники сложо сосредоточиться на задачах, что приводит к снижения производительности. Мысли могут быть хаотичными, а внимание рассениями.

***Трудности с концентрация:**

***Tpygnoctus:**

***Tpygno
```

Figure 18: Typical example of **hallucination** in Client Selection with **skin cancer datasets** for **thinking/reasoning** models. It shows the model is answering in Russian language which is totally unrelated to the given task.

```
5779
5780
5781
                                    r mentioned a specific file path: /home/chrid-19/Desktop/DPhil work/AI_agents/H3Builder/Referencefiles/DataCard/individual_datasets/dataset_23.json. I should start by reading that file to get the datasets. The function read_files is available for this purpose.
5782
5783
                                               sure I use the correct syntax for the tool call. The file path is provided, so I'll pass that to the read_files function. Then, based on the re
5784
                                          ead_files", "arguments": {"file_paths": ["/home/chri6419/Desktop/DPhil work/A
itent: <think>
5785
                                       Content: «think»
user wants to train a skin cancer detection model using dermatology images. Let me start by understanding the requirements. The task is skin cancer de
ch likely refers to dermoscopic images or similar medical images of the skin.
5786
                             rst, I need to check the available datasets. The user monitoned a specific file path: |home/chric4si9/mestago/pibhl unvix/al_apents/PiBuilder/meferencefiles/platicard/individual_datasets/dataset_17.jos should use the read_files function to read this file and examine its contents. The dataset descriptions will help determine if any of the are suitable for the task.
5787
                             ll call the read files function with the provided file path. Once I have the dataset information, I'll look for entries where the description mentions skin cancer detection or the modality is dermatology images. If there are matching datasets, I'll list their names. If not, I'll inform the user that no suitable dataset was found.
5788
5789
                              name": "read_files", "arguments": ("file_paths": ["/home/chri6419/Desktop/DPhil work/AI_agents/M3Builder/ReferenceFiles
5790
                                   ication_2_content <think>
et me process this step by step. The user is working on a skin cancer detection project using dermatology images. They provided a file path to a dataset description file
et is suitable for their task.
5791
                             cxt, I need to read the contents of the aperticed 2500 file. The file path is /home/chrid=19/Desktop/DPhil work/al_agents/H380itIder/#eferenceFiles/DataCard/individual_datasets/dataset_17.json. Using the read-files function. I'll access the dataset descriptions stored therein
5792
                             ce I have the data, I'll look for any entries that mention "skin cancer detection" in their description. The modality should be dermatology images, which might
kin lestons," or similar. If there's a match, I'll confirm the dataset's suitability. If mot, I'll inform the user that no relevant dataset was found.
5793
                                     also check if the dataset's structure and content align with typical requirements for training a detection model, such as having labeled images of skin linext step is to proceed with training. Otherwise, the user might need to look for alternative datasets or
5794
                                   the user wants to train a skin cancer detection model using dermatology inages. Let me start by understanding the requirements. The task is skin cancer
which likely refers to dermoscopic images or similar.
5795
5796
                             reading the 350M file, I'll look at each dataset's description to see if it matches the task and modality. The goal is to find datasets that contain dermatology images used for skin cancer detection the file isn't found or doesn't contain relevant data, I'll need to inform the user.
5797
5798
5799
```

Figure 19: Client Selection with **skin cancer datasets** for **thinking/reasoning** models. It shows when/how the server approves clients for training whereas when/how it declines.

```
5807
 5808
                                                                                             1_content Task: Binary breast cancer detection (benign vs. malignant)
                                                              ents, please respond with the name of the selected dataset that matches this req
ector Content: *Dataset Name** : BreaKHis 400X
 5809
 5810
                                                            mmunication_2_content Approved. Prepare for training.
lector Content: **Dataset Name** : HEROHE
5811
 5812
                                                             nmunication_2_content Approved. Prepare for training.
lector Content: **Dataset Name** : breast_cancer_segmentation
5813
                                                              munication_2_content Approved. Prepare for training ector Content: no dataset
5814
5815
                                                            planation: The "BreCaHAO" dataset is a breast cancer histopathology dataset, but it is focused on multi-class segmentation of histological structures (nitosis, a
ry classification of breast cancer as benign vs. malignant. Therefore, it does not match the requirement for binary breast cancer detection.
 5816
                                                           ommunication_2_content Client not needed for the task.
elector Content: **Dataset Name** : SLN Breast
5818
                                                             mmunication_2_content Approved. Prepare for training.
Lector Content: no dataset
5819
 5820
5821
5822
                                                             elamation: The watilable dataset, "gastric cancer," contains histopathology images of gastric (stomach) cancer, not breast cancer. Therefore, it does not match the requirement for binary breast can
be cection (bening vs. malignant) vs. malignative vs. malignative cancer.
 5823
 5824
 5825
                                                           planston. The available dataset, "ovarian_cancer", contains histopathology inages for ovarian cancer, not breast cancer. Therefore, it does not match the requirement for binary breast cancer detectibeneds us. miletance and the contract of the contract of
5826
```

Figure 20: Client Selection with **histopathology datasets** for **non-thinking/reasoning** models on **breast cancer detection task**. It shows when/how the server approves clients for training whereas when/how it declines.

```
5834
                       Dataset Name** : No sultable dataset for binary breast cancer detec
5835
                           ication_2_content Client not needed for the task
or Content: Based on the file content, the dataset suitable for your task is:
5836
5837
                                et is a breast cancer histopathology dataset for binary classification. It contains two folders: 'IDC' and 'non-IDC', re
ppe of breast cancer. This makes it suitable for your task of binary breast cancer detection from histopathology images.
5838
5839
5840
5841
                              ition_2_content Approved. Prepare for training
Content: Based on the content of the json file, the dataset available is "gastro_cancer_nsi_vs_mss". This dataset is a gastrointestinal histopathology dataset with subfolders representing MSI answers for classification.
5842
                        ever, this dataset does not suit your requirements as it is not a binary breast cancer detection dataset. It is a gastrointestinal cancer dataset and the classes are MSI and MSS, not benign and malig which are required for binary breast cancer detection.
5843
                        efore, there is no suitable dataset available for the given task of bi
5844
                       Dataset Name** : No sultable dataset
5845
                        5846
5847
                                 on_2_content Approved. Prepare for training.
ntent: The dataset in the provided JSDN file is named "tung_and_colon". However, this dataset is not suitable for the task of binary breast cancer detection as it consists of lung histopathol not breast. Therefore, there is no suitable dataset for the clown task in the provided JSDN file.
5848
5849
                       Dataset Name** : No suitable dataset
5850
5851
```

Figure 21: Client Selection with **histopathology datasets** for **non-thinking/reasoning** models on **breast cancer detection task**. It shows when/how the server approves clients for training whereas when/how it declines.

```
5860
5861
5862
                                   r clients, please respond with the name of the selected dataset that contains histopathology images for binary breast cancer detection
ector Content: The dataset that matches the task of binary breast cancer detection using histopathology images is:
5863
5864
5865
5866
                                         cation_2_content Approved. Prepare for training.
or Content: Based on the content of the dataset description file, the suitable dataset for your task is
5867
5868
5869
                                               on 2 content Approved. Prepare for training moment from the dataset **Broca+MDT** contains breast cancer histopathology images. Nowever, it supports multi-class segmentation into six contains breast cancer to the information from the dataset description, the dataset **Broca+MDT** contains proportion, tumor muclei, non-tumor muclei, tubule, and non-tubule. It does not explicitly mention whether it can be used for binary breast cancer detection (i.e., benign and malignant) are nay need further information or preprocessing to use this dataset for the specific task of binary breast cancer detection.
5870
5871
5872
5873
5874
                                           cation_2_content Approved. Prepare for training.

Content: Based on the dataset description, the dataset named "melanoma" contains histopathology images of melanoma, not breast cancer. Therefore, this dataset is not suitable for binary breast
5875
5876
5877
                                              tion_z_content Client not needed for the task.
onleni: From the dataset description, the dataset "gastric_cancer" contains HBE-stained histopathological images from human gastric cancer. However, this dataset is not suitable for your task
were not contain histopathology images for binary breast cancer detection.
5878
5879
5880
```

Figure 22: Client Selection with **histopathology datasets** for **non-thinking/reasoning** models **on breast cancer detection task**. It shows when/how the server approves clients for training whereas when/how it declines.

```
5900
5901
5902
5903
                            ignant classification.
5904
                          .cation_2_content Client not needed for the task
or Content: no dataset: The provided dataset "BreCaHAD" is inten
5905
                              5906
                          .cation_2_content Approved. Prepare for training.
or Content: no dataset: The available dataset "melanoma" is for seç
5907
5908
                            tion 2 content Client not needed for the task.
Content: no dataset - The dataset provided ("gastric_cancer") is focused on gastric cancer histopathological images and does not match the task of binary classification for
he yearloans!
5909
                             5910
5911
                    ommunication_2_content Client not needed for the task
:lector Content: **Dataset Name**: breast histo
5912
5913
                     mmunication_2_content Approved. Prepare for training
lector Content: **Dataset Name** : breast_cancer_histo
5914
                            tion_content Approved. Prepare for training Content: no dataset: The available dataset, "gastro_cancer_msi_vs_mss", is a gastrointestinal histopathology dataset classified into MSI and MSS clarancer detection task (benign vs. malignam).
5915
5916
5917
                         icetion 2 content Client not needed for the eask
or content: no dataset: The provided dataset 'lung and colon' is designed for lung and colon tissue classification (including lung b
the requirement for binary breast cancer detection using histopathology images.
5918
5919
5920
```

Figure 23: Client Selection with **histopathology datasets** for **non-thinking/reasoning models** on **breast cancer detection task**. It shows when/how the server approves clients for training whereas when/how it declines.

```
5943
5944
5945
5946
5947
5948
         **Task:** Train a binary classification model for breast cancer detection.
5949
          *Modality:** Histopathology images
5950
5951
         **Client Agents:** Please respond with the name of the dataset you have that matches this requirement.
5952
         Selector Content: **Dataset Name** : BreaKHis_400X
5953
5954
5955
         communication 2 content Approved. Prepare for training.
5956
         Selector Content: **Dataset Name:** HEROHE
5957
5958
5959
         communication_2_content Approved. Prepare for training.
5960
5961
         Selector Content: **Dataset Name** : breast_cancer_segmentation
5962
         5963
         communication_2_content This response indicates the client has returned a dataset (breast_cancer_segmentation).
5964
5965
         Therefore, the appropriate server response is:
5966
         **Approved. Prepare for training**
5967
         Selector Content: **Dataset Name**: BreCaHAD
5968
5969
5970
         communication_2_content Approved. Prepare for training.
5971
         Selector Content: **Dataset Name**: SLN_Breast
5972
5973
5974
          communication_2_content Approved. Prepare for training.
5975
         Selector Content: no dataset : The provided dataset is for melanoma segmentation, not breast cancer detection.
5976
5977
5978
5979
          communication_2_content Client not needed for the task.
5980
5981
         Selector Content: no dataset
         This dataset contains gastric cancer images and not breast cancer
5982
```

Figure 24: Client Selection with **histopathology datasets** for **non-thinking/reasoning** models on **breast cancer detection task**. It shows when/how the server approves clients for training whereas when/how it declines.

```
communication__content ethink-
day, so I need to figure out how to respond to the client's nessage. The client sent a dataset named AGROBAT_Breast. By task is to determine the appropriate response based on whether they provided a dataset, I should check if the client returned one or nore datasets. In this case, they did provide a dataset, I should check if the client returned one or nore datasets. In this case, they did provide a dataset, I should check if the client returned one or nore datasets. Is not to respond to the client's presented for the task." But since they did provide one, the first response is the correct one.

If they hadn't provided any dataset, I would respond with "client not needed for the task." But since they did provide one, the first response is the correct one.

I should nake sure I'n following the instructions accurately. The client's message clearly states the dataset name, so there's no ambiguity here. Therefore, the response should be "Approved. Prepare for training."

"Approved. Prepare
```

Figure 25: Client Selection with **histopathology datasets** for **thinking/reasoning models** on **breast cancer detection task**. It shows when/how the server approves clients for training whereas when/how it declines.

```
6022
6023
                                                                                                         ink>
train a binary breast cancer detection model using histopathology (mages. First, I need to figure out what exactly they're asking for. They retween benign and malignant tumors. The data modality is histopathology images, which are microscope slides of tissue samples.
6024
                                                                               tract the task and modality from their requirement. The task is binary classification, specifically for breast cancer detection. The modality is histo
6025
6026
6027
6028
6029
6030
6031
                                                           I should understand the instructions clearly. If the client returns one or more datasets, the server should respo
be "Client not needed for the task."
                                                 king at the client's response, they provided the dataset name. That means they have returned at least one dataset. So, according to the instructions, the appropriate resort retaining.
6032
6033
6034
                                              oproved. Prepare for training.
Plector Content: The selected dataset for your task is:
6035
6036
                                                is dataset is suitable because it contains histopathology images specific to breast cancer, with cases labeled as positive (malic
6037
                                                            on task.
                                                           cation__content «think»

cation__content «think»

it received a message from the cilent about a dataset they'we selected. Let me try to understand what they're saying. The dataset name is HEROME, and they mentioned it's suitable because it in this process the selection of the s
6038
6039
                                               m, the client has clearly provided a dataset, so according to the instructions, I should respond with "agroved, Prepare for training." That makes sense because they've net the criteria by selecting learnt dataset. I don't think I need to do anything else here except acknowledge their selection and give the go-ahead for the next step.
6040
                                                should make sure my response is clear and concise. No need for extra information since they've already given the necessary details. Just the approval and the instruction to prepare for training. Airl I think I've got it.
6041
6042
```

Figure 26: Client Selection with **histopathology datasets** for **thinking/reasoning models** on **breast cancer detection task**. It shows when/how the server approves clients for training whereas when/how it declines.

```
6050
                                                                     taset available, named "bone-fracture-dataset," contains X-ray images of fractured and non-fractured bones, which does not match the requirement for pnit is not suitable for the specified task.
6051
6052
6053
                                                           mmunication_2_content Client not needed for the task.
lector Content: **Dataset Name** : anatomy-in-chest-x-rays-pax
6054
 6055
                                                      The dataset "unifesp-xray-bodypart-classification" is designed for identifying 20 body parts using whole body X-ray inages and does not specifically target pneumonia detection from chest X-ray it is not suitable for the task of pneumonia detection from chest X-ray it is not suitable for the task of pneumonia detection from chest X-ray inages.
6056
6057
                                                                     ication_2_content Client not needed for the task.
or Content: **Dataset Name** : nih-chest-x-rays-bbox-version
6058
6059
                                                          mmunication_2_content Approved. Prepare for training.
lector Content: **Dataset Name** : covid19-xray-dataset-train-test-set:
6060
                                                         ommunication_2_content Approved. Prepare for training
elector Content: **Dataset Name** : no dataset
 6061
                                                         se provided dataset "dental-radiography" contains dental X-ray images, which are not suitable for training a pneumonia detection model that requires chest X-ray images. Therefore, there is no app dataset available for this specific task.
6062
 6063
                                                          mmunication_2_content Client not needed for the task.
lector Content: **Dataset Name** : chest-xray-pneumor
6064
6065
                                                           mmunication_2_content Approved. Prepare for training
lector Content: **Dataset Name** : no dataset
6066
                                                         te dataset provided, "Bone break-classification-inage-dataset," focuses on different types of bone fractures in X-ray Images, such as elbow, wrist, and shoulder fractures. It does not contain tata specifically related to pneumonia detection inage-dataset," socuses on one execusive for the task at hand. Therefore, it is not suitable for training a purpose dataset commonly and the provided in the 
 6067
 6068
 6069
```

Figure 27: Client Selection with **X-Ray datasets** for **non-thinking/reasoning** models **on pneumonia detection task**. It shows when/how the server approves clients for training whereas when/how it declines.

```
6078
6079
6080
6081
                                                             ontent Client not needed for the task.
The dataset variable is named *cookid9-and-normal-x-ray-dataset-balanced**. It contains chest X-ray images balanced across COVID-19 positive and normal cases. He data for pneumonia detection. Therefore, it may not be suitable for your task of pneumonia detection from chest XBay images.
6082
6083
                                                          content Approved. Prepaire for training
t: Based on the dataset description from the json file, the dataset "human-bone-fractures-image-dataset-hbfmid" is organized into folders based on fracture type and anatomical location
not seen to be suitable for pneumonia detection as it is focused on bone fractures. Therefore, there is no suitable dataset for the given task.
6084
6085
                                                                tent Client not needed for the task
ased on the dataset description, the dataset available is "osteoporosis-database" which includes knee X-ray images labeled for osteoporosis stages. Unfortunately, this dataset is not
k of premunal achection from chest Xaay images as it contains knee X-ray images and is labeled for osteoporosis stages, not pneumonia.
6086
                                     available dataset does not match the required task and modality. The task requires a dataset of chest XRay images labeled for pneumonia, but the available dataset contains knee XRay images labeled fi
steoporosis. Therefore, we need to look for a different dataset that is more sultable for the task.
6087
6088
                                    reunitation_2 content client on needed for the task, needed for the task, needed for the content client content in detacts available is "Guastes Name": chest-x-ray-lungs-segmentation, However, this dataset is specifically for tuberculosis detection and lung segmentation, not for pn client interfere, it may not be suitable for your task of pneumonia detection from chest XBay images.
6089
6090
                                     munication_2_content Client not needed for the task.
.ector Content: The suitable dataset for your requirement is:
6091
6092
6093
6094
                                              ation_2_content Approved. Prepare for training.
Content: The dataset available is **pataset Name**: lung-segment.
table for your task of pneumonia detection from chest XRay images
6095
6096
```

Figure 28: Client Selection with **X-Ray datasets** for **non-thinking/reasoning models** on **pneumonia detection task**. It shows when/how the server approves clients for training whereas when/how it declines.

```
6104
6105
6106
6107
6108
6109
6110
6111
6112
6113
6114
6115
6116
6117
                         ication 2_content Approved. Prepare for training
or Content: Based on the information from the dataset description, the dataset named "fracatlas" is an annotated fracture image dataset. It does not seem to be suitable for p
**Nay images since it is focused on fractures, not pneumonia.
6118
6120
6121
                     Dataset Name** : No suitable dataset
6122
6123
```

Figure 29: Client Selection with **X-Ray datasets** for **non-thinking/reasoning** models on **pneumonia detection task**. It shows when/how the server approves clients for training whereas when/how it declines.

```
6132
6133
6134
6135
6136
6137
                                      nd with the name of the selected dataset that matches this task and modality.

tents "States haves" in odd states of the contain chest X-ray (mages labeled for pineumonia vs. normal or other lung conditions. The available "fracture-multi-region-x-ray-data" focuses on boggions like links and knees, which is unrelated to pineumonia detection in chest X-rays. No suitable dataset exists in the provided options for this task.
6138
6139
6140
6141
                                         ns say: if the client returns datasets, say "Approved. Prepare for training". If they don't have any, respond "Client not needed for the task
6142
                                  e, the correct response should be "Client not needed for the task"
6143
6144
6145
6146
                          wy, let me process this. The user provided a response indicating that the dataset named "bone-fracture-dataset" isn't suitable for pneumonia detection in chest X-rays. They mentioned no relevant data
6147
                          e instructions say if the Client returns no dataset, I should respond with 'Client not needed for the task'. Since the user states there's no relevant dataset here, that means the client doesn't have necessary data for the task. Herefore, the correct response is to inform them they're not needed.
6148
6149
6150
```

Figure 30: Client Selection with **X-Ray datasets** for **thinking/reasoning models** on **pneumonia detection task**. It shows when/how the server approves clients for training whereas when/how it declines.

Figure 31: Client Selection with **X-Ray datasets** for **thinking/reasoning** models on **pneumonia detection task**. It shows when/how the server approves clients for training whereas when/how it declines.

Figure 32: Client Selection with **X-Ray datasets** for **thinking/reasoning models** on **pneumonia detection task**. It shows when/how the server approves clients for training whereas when/how it declines.

```
**Postaset Name**: NID: Chest X-ray Dataset
aggrowed. Propage for training.

Selector content: "Obseased Name**: condition radiography-database
This dataset is suitable for training a pneumonia detection model as it contains images specifically labeled for pneumonia.

***Communication_2_content ethics**

**Communication_2_content ethics**

**Communicat
```

Figure 33: Client Selection with **X-Ray datasets** for **thinking/reasoning** models on **pneumonia detection task**. It shows when/how the server approves clients for training whereas when/how it declines.

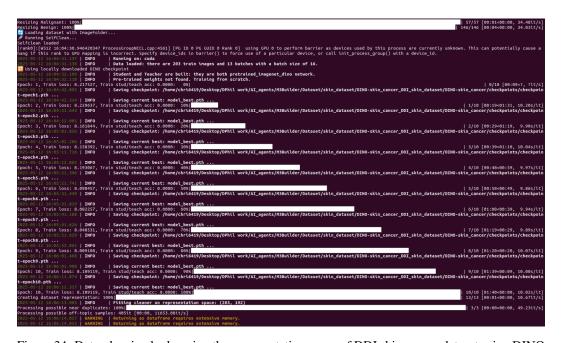


Figure 34: Data-cleaning by learning the representation space of DDI skin cancer dataset using DINO

```
Resiting Keratosis: 100N|
Resiting Carcinoms: 100N|
Resiting Carcinoms: 100N|
Resiting Carcinoms: 100N|
Resiting Carcinoms: 100N|
Resiting Actions: 100N|
Resiting Actions: 100N|
Resiting Actions: 100N|
Resiting Carcinoms: 100N
```

Figure 35: Data-cleaning by learning the representation space of augmented-skin-condition-dataset-kaggle using DINO