

915	Contents	
916	A Broader Impact	24
917	B Compute and Environment Configuration	24
918	C Platform Repository	24
919	D Author Statement	24
920	E Datasheet for Datasets	25
921	E.1 Motivation	25
922	E.2 Composition	25
923	E.3 Collection process	26
924	E.4 Preprocessing/cleaning/labeling	27
925	E.5 Uses	27
926	E.6 Distribution	27
927	E.7 Maintenance	28
928	F Training Task – Lung Opacity Detection	29
929	F.1 Task Description	29
930	F.2 License and Ethics	29
931	F.3 Access and Preprocessing	29
932	F.4 Data Samples	29
933	G Training Task – COVID-19 Detection	30
934	G.1 Task Description	30
935	G.2 License and Ethics	30
936	G.3 Access and Preprocessing	30
937	G.4 Data Samples	30
938	H Training Task – ECG Abnormal Detection	31
939	H.1 Task Description	31
940	H.2 License and Ethics	31
941	H.3 Access and Preprocessing	31
942	H.4 Data Samples	31
943	I Training Task – Mortality Prediction	32
944	I.1 Task Description	32
945	I.2 License and Ethics	32
946	I.3 Access and Preprocessing	32

947	I.4	Data Samples	32
948	J	Validation Task – Enlarged Cardiomendiastinum Detection	34
949	J.1	Task Description	34
950	J.2	Access and Preprocessing	34
951	J.3	Data Samples	34
952	K	Validation Task – Sepsis Prediction	35
953	K.1	Task Description	35
954	K.2	License and Ethics	35
955	K.3	Access and Preprocessing	35
956	K.4	Data Samples	35
957	L	Validation Task – MedVQA	36
958	L.1	Task Description	36
959	L.2	License and Ethics	36
960	L.3	Access and Preprocessing	36
961	L.4	Data Samples	36
962	M	Validation Task – Signal Noise Clarification	37
963	M.1	Task Description	37
964	M.2	License and Ethics	37
965	M.3	Access and Preprocessing	37
966	M.4	Data Samples	37
967	N	Implementation Details	38
968	N.1	Model Details	38
969	N.2	Optimizer Hyperparameters	38
970	N.3	Task Prompts	38
971	N.4	Baselines	38

972 **A Broader Impact**

973 The Federated Medical Knowledge Injection (FMKI) platform introduces a transformative approach
974 in healthcare AI, addressing critical issues of data privacy and accessibility by leveraging federated
975 learning to inject medical knowledge into foundation models. This method not only complies
976 with stringent health regulations, thereby protecting patient confidentiality, but also enhances the
977 scalability and adaptability of medical foundation models. By enabling these models to utilize diverse,
978 multi-modal medical data without direct data sharing, FMKI significantly broadens the potential
979 applications of AI in healthcare, offering improved diagnostic accuracy and personalized treatment
980 options. Furthermore, the platform facilitates equitable technology access, allowing institutions with
981 varying resources to participate in and benefit from cutting-edge medical AI developments. This
982 innovative approach not only promises to improve global healthcare outcomes but also sets new
983 benchmarks in the ethical development and deployment of AI technologies in sensitive sectors.

984 **B Compute and Environment Configuration**

985 All experiments are conducted on an NVIDIA A100 with CUDA version 12.0, running on a Ubuntu
986 20.04.6 LTS server. More details can be found in the GitHub repository.

987 **C Platform Repository**

988 We have established a GitHub repository, available at [https://github.com/psudslab/](https://github.com/psudslab/FEDMEKI)
989 `FEDMEKI`. This repository includes resources for data processing, baselines, environmental setup,
990 our proposed platform, and sample execution scripts. All the details have been documented at the
991 ReadMe file³. We are committed to continuously updating this repository with additional modalities,
992 datasets, and tasks.

993 **D Author Statement**

994 As authors of this repository and article, we bear all responsibility in case of violation of rights and
995 licenses. We have added a disclaimer on the repository to invite original dataset creators to open
996 issues regarding any license-related matters.

³<https://github.com/psudslab/FEDMEKI/blob/main/README.md>

E Datasheet for Datasets

E.1 Motivation

- **For what purpose was the dataset created?**

This work investigates a novel yet practical task – scaling existing medical foundation models by injecting diverse medical knowledge with distributed private medical data. However, no available datasets are suitable for evaluation. Thus, we curated a new multi-site, multi-modal, and multi-task dataset, including five training tasks and three validation tasks and covering six commonly used medical modalities.

- **Who created the dataset (e.g., which team, research group) and on behalf of which entity (e.g., company, institution, organization)?**

The authors of this paper.

- **Who funded the creation of the dataset? If there is an associated grant, please provide the name of the granter and the grant name and number.**

This work is partially supported by the National Science Foundation under Grant No. 2238275, 2333790, 2348541, and the National Institutes of Health under Grant No. R01AG077016.

E.2 Composition

- **What do the instances that comprise the dataset represent (e.g., documents, photos, people, countries)?**

The FEDMEKI data suit contains medical images and the corresponding annotations for the pneumonia detection and COVID-19 detection tasks; ECG signals and the labels for the ECG abnormal detection task; 48 clinical features for the mortality prediction and sepsis prediction tasks; medical text (questions, candidate answers, document collections, ground truths) for the MedQA task; ECG signals, questions and answers for the signal noise clarification task; and medical images, questions, and answers for the MedVQA task.

- **How many instances are there in total (of each type, if appropriate)?**

The Lung Opacity Detection task has 18,406 samples, the ECG Abnormal Detection task has 21,797 samples, and the Mortality Prediction task has 38,129 samples. The COVID-19 Detection task has 13,808 samples. The MedVQA, Signal Noise Clarification and Sepsis Prediction tasks each contain 1,000 samples. Additionally, the Enlarged Cardiomeastinum Detection task has 234 samples. Detailed information about the data can be found in Table 1.

- **Does the dataset contain all possible instances or is it a sample (not necessarily random) of instances from a larger set?**

The ECG Abnormal Detection task includes all available samples from its corresponding database. The Lung Opacity Prediction, COVID-19 Detection, and Mortality Prediction tasks encompass all data samples with available binary labels, making them subsets of the original dataset. For validation tasks without a predefined test set or with an excessively large test set, we randomly selected 1,000 samples for testing. These tasks include MedVQA, Signal Noise Clarification, and Sepsis Prediction. For the Enlarged Cardiomeastinum Detection task, the original database provided a small test set of 234 samples, which we have retained.

- **What data does each instance consist of?**

The Lung Opacity Prediction, COVID-19 Detection, and Enlarged Cardiomeastinum Detection tasks involve radiological images, while the ECG Abnormal Detection task involves 12-channel, 10-second ECG signals. The Mortality Prediction and Sepsis Prediction tasks cover temporal features involving vital signs, lab events, and input/output data. The Signal Noise Clarification task includes signal-text pairs, while the MedVQA task comprises image-text pairs.

- **Is there a label or target associated with each instance?**

The answer (label) is provided for each instance.

- 1047 • **Is any information missing from individual instances? If so, please provide a descrip-**
 1048 **tion, explaining why this information is missing (e.g., because it was unavailable). This**
 1049 **does not include intentionally removed information, but might include, e.g., redacted**
 1050 **text.**
 1051 No.
- 1052 • **Are relationships between individual instances made explicit (e.g., users' movie ratings,**
 1053 **social network links)?**
 1054 No.
- 1055 • **Are there any errors, sources of noise, or redundancies in the dataset?**
 1056 Questions are created by filling the slots in the templates with pre-defined values and records
 1057 from the database. Thus, some questions can be grammatically incorrect but not critical
 1058 (e.g., verb tense).
- 1059 • **Is the dataset self-contained, or does it link to or otherwise rely on external resources**
 1060 **(e.g., websites, tweets, other datasets)?**
 1061 The proposed dataset depends on several open-source databases: RSNA [17], COVQU [18],
 1062 PTB-XL [19], MIMIC-III [23], CheXpert [21], VQA-RAD [11], and ECG-QA [22].
- 1063 • **Does the dataset contain data that might be considered confidential (e.g., data that is**
 1064 **protected by legal privilege or by doctor-patient confidentiality, data that includes the**
 1065 **content of individuals' non-public communications)?**
 1066 No.
- 1067 • **Does the dataset contain data that, if viewed directly, might be offensive, insulting,**
 1068 **threatening, or might otherwise cause anxiety?**
 1069 No.
- 1070 • **Does the dataset relate to people?**
 1071 Yes.
- 1072 • **Does the dataset identify any subpopulations (e.g., by age, gender)?**
 1073 No.
- 1074 • **Does the dataset contain data that might be considered sensitive in any way (e.g.,**
 1075 **data that reveals race or ethnic origins, sexual orientations, religious beliefs, political**
 1076 **opinions or union memberships, or locations; financial or health data; biometric or**
 1077 **genetic data; forms of government identification, such as social security numbers;**
 1078 **criminal history)?**
 1079 No. The source datasets are already de-identified.

1080 E.3 Collection process

- 1081 • **How was the data associated with each instance acquired?**
 1082 We directly used the original data instance to curate our own dataset.
- 1083 • **What mechanisms or procedures were used to collect the data (e.g., hardware appara-**
 1084 **tuses or sensors, manual human curation, software programs, software APIs)?**
 1085 We mainly used Python scripts to collect, process and label the data.
- 1086 • **If the dataset is a sample from a larger set, what was the sampling strategy (e.g.,**
 1087 **deterministic, probabilistic with specific sampling probabilities)?**
 1088 The random sampling involved in this study relies on specific seed (42), thus becomes
 1089 deterministic.
- 1090 • **Who was involved in the data collection process (e.g., students, crowd workers, con-**
 1091 **tractors), and how were they compensated (e.g., how much were crowd workers paid)?**
 1092 The data collection process was fully performed by the study's authors.
- 1093 • **Over what timeframe was the data collected?**
 1094 N/A

- 1095 • **Were any ethical review processes conducted (e.g., by an institutional review board)?**
1096 N/A.
- 1097 • **Does the dataset relate to people?**
1098 Yes.
- 1099 • **Did you collect the data from the individuals in question directly, or obtain it via third parties or other sources (e.g., websites)?**
1100 All data are collected through open-source database without interaction with individuals.
1101
- 1102 • **Were the individuals in question notified about the data collection?**
1103 N/A.
- 1104 • **Did the individuals in question consent to the collection and use of their data?**
1105 N/A.
- 1106 • **If consent was obtained, were the consenting individuals provided with a mechanism to revoke their consent in the future or for certain uses?**
1107 N/A.
1108
- 1109 • **Has an analysis of the potential impact of the dataset and its use on data subjects (e.g., a data protection impact analysis) been conducted?**
1110 The dataset does not have individual-specific information.
1111

1112 **E.4 Preprocessing/cleaning/labeling**

- 1113 • **Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucketing, tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing of missing values)?**
1114 Yes. The preprocessing on MIMIC-III data follows existing work [20].
1115
- 1116 • **Was the “raw” data saved in addition to the preprocess/cleaned/labeled data (e.g., to support unanticipated future uses)?**
1117 N/A.
1118
- 1119 • **Is the software that was used to preprocess/clean/label the data available?**
1120 Preprocessing, cleaning, and labeling are done via Python.
1121

1122 **E.5 Uses**

- 1123 • **Has the dataset been used for any tasks already?**
1124 No.
- 1125 • **Is there a repository that links to any or all papers or systems that use the dataset?**
1126 No.
- 1127 • **What (other) tasks could the dataset be used for?**
1128 While the dataset is curated for research on federated medical knowledge injection problem,
1129 other studies concerning developing centralized medical foundation model can also leverage
1130 the dataset.
- 1131 • **Is there anything about the composition of the dataset or the way it was collected and preprocessed/cleaned/labeled that might impact future uses?**
1132 N/A.
1133
- 1134 • **Are there tasks for which the dataset should not be used?**
1135 N/A.

1136 **E.6 Distribution**

- 1137 • **Will the dataset be distributed to third parties outside of the entity (e.g., company, institution, organization) on behalf of which the dataset was created?**
1138 No.
1139

- 1140 • **How will the dataset be distributed?**
1141 The preprocessing code is available at [https://github.com/psudslab/](https://github.com/psudslab/FEDMEKI/tree/main/data_preprocess)
1142 [FEDMEKI/tree/main/data_preprocess](https://github.com/psudslab/FEDMEKI/tree/main/data_preprocess). Users can download corresponding
1143 dataset and utilize the preprocessing scripts for generating the final dataset used in this study.
- 1144 • **Will the dataset be distributed under a copyright or other intellectual property (IP)**
1145 **license, and/or under applicable terms of use (ToU)?**
1146 The dataset is released under MIT License.
- 1147 • **Have any third parties imposed IP-based or other restrictions on the data associated**
1148 **with the instances?**
1149 No.
- 1150 • **Do any export controls or other regulatory restrictions apply to the dataset or to**
1151 **individual instances?**
1152 No.

1153 E.7 Maintenance

- 1154 • **Who will be supporting/hosting/maintaining the dataset?**
1155 The authors of this paper.
- 1156 • **How can the owner/curator/manager of the dataset be contacted(e.g., email address)?**
1157 Contact the first authors (jqwang@psu.edu and xcwang@psu.edu).
- 1158 • **Is there an erratum?**
1159 No.
- 1160 • **Will the dataset be updated (e.g., to correct labeling erros, add new instances, delete**
1161 **instances)?**
1162 If any corrections are required, our plan is to upload an updated version of the dataset with
1163 comprehensive explanations for the changes. Furthermore, as we broaden our QA scope, we
1164 will consistently update the dataset with new QA templates/instances.
- 1165 • **If the dataset relates to people, are there applicable limits on the retention of the data**
1166 **associated with the instances (e.g., were the individuals in question told that their data**
1167 **would be retained for a fixed period of time and then deleted)?**
1168 N/A
- 1169 • **Will older versions of the dataset continue to be supported/hosted/maintained?**
1170 Primarily, we plan to maintain only the most recent version of the dataset. However, under
1171 certain circumstances, such as significant updates to our dataset or the need for validation
1172 of previous research work using older versions, we will exceptionally preserve previous
1173 versions of the dataset for up to one year.
- 1174 • **If others want to extend/augment/build on/contribute to the dataset, is there a mecha-**
1175 **nism for them to do so?**
1176 Contact the authors of this paper.

1177 **F Training Task – Lung Opacity Detection**

1178 **F.1 Task Description**

1179 In the United States, pneumonia keeps the ailment on the list of top 10 causes of death in the country.
1180 The task is to locate lung opacities on chest radiographs. In this challenge [17], 18,406 images are
1181 annotated as either Lung Opacity or Normal, providing a basis for extracting the binary classification
1182 task. The task is to develop an algorithm to detect visual indicators of pneumonia in medical images.
1183 Specifically, the algorithm needs to identify and localize lung opacities in chest radiographs.

1184 **F.2 License and Ethics**

1185 This dataset is permitted to access and utilize these de-identified imaging datasets and annotations for
1186 academic research, educational purposes, or other commercial or non-commercial uses, provided you
1187 adhere to the appropriate citations.

1188 **F.3 Access and Preprocessing**

1189 The resource is available to access via the official website at [https://www.rsna.org/rsnai/](https://www.rsna.org/rsnai/ai-image-challenge/rsna-pneumonia-detection-challenge-2018)
1190 [ai-image-challenge/rsna-pneumonia-detection-challenge-2018](https://www.rsna.org/rsnai/ai-image-challenge/rsna-pneumonia-detection-challenge-2018) and Kag-
1191 [gle](https://www.kaggle.com/c/rsna-pneumonia-detection-challenge/overview) at [https://www.kaggle.com/c/rsna-pneumonia-detection-challenge/](https://www.kaggle.com/c/rsna-pneumonia-detection-challenge/overview)
1192 [overview](https://www.kaggle.com/c/rsna-pneumonia-detection-challenge/overview). It includes dataset description, annotations, and mapping from RSNA image dataset
1193 to original NIH dataset. The data is organized as a set of patient IDs with corresponding im-
1194 age class annotations, including "No Lung Opacity/Not Normal," "Normal," and "Lung Opac-
1195 ity." We collected images labeled as either "Normal" or "Lung Opacity" and formulated the
1196 problem as a binary classification task. The code for preprocessing is available at https://github.com/psudslab/FEDMEKI/tree/main/data_preprocess.
1197

1198 **F.4 Data Samples**

1199 We provide a random data sample from the dataset and visualize it in Figure 2.

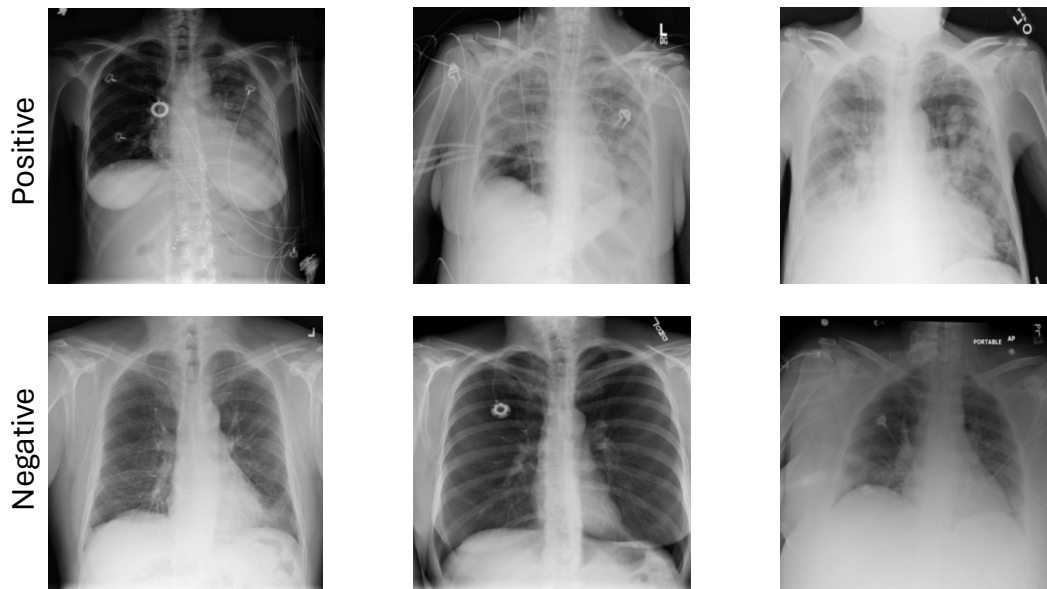


Figure 2: Data sample of lung opacity detection.

1200 **G Training Task – COVID-19 Detection**

1201 **G.1 Task Description**

1202 This task challenges the model to assess whether X-ray images display symptoms of Covid-19,
1203 thereby evaluating the model’s proficiency in interpreting medical imagery. For this purpose, we
1204 employ the COVQU dataset [18].

1205 **G.2 License and Ethics**

1206 The licensing and ethical compliance adhere to the regulations established by the original datasets.

1207 **G.3 Access and Preprocessing**

1208 This dataset can be accessed via the link at <https://www.kaggle.com/datasets/tawsifurrahman/covid19-radiography-database>. This dataset, featuring COVID-
1209 19, normal, and other lung infection categories, is being released incrementally. The initial release
1210 comprised 219 COVID-19, 1,341 normal, and 1,345 viral pneumonia chest X-ray (CXR) images.
1211 The first update expanded the COVID-19 category to include 1,200 CXR images. In the second
1212 update, the collection was further enlarged to include 3,616 COVID-19 positive cases, along with
1213 10,192 normal, 6,012 lung opacity (non-COVID lung infection), and 1,345 viral pneumonia images,
1214 complete with corresponding lung masks. We selected normal and COVID-19 positive images to
1215 formulate this task as a binary classification problem.
1216

1217 **G.4 Data Samples**

1218 We provide a random data sample from the dataset and visualize it in Figure 3.

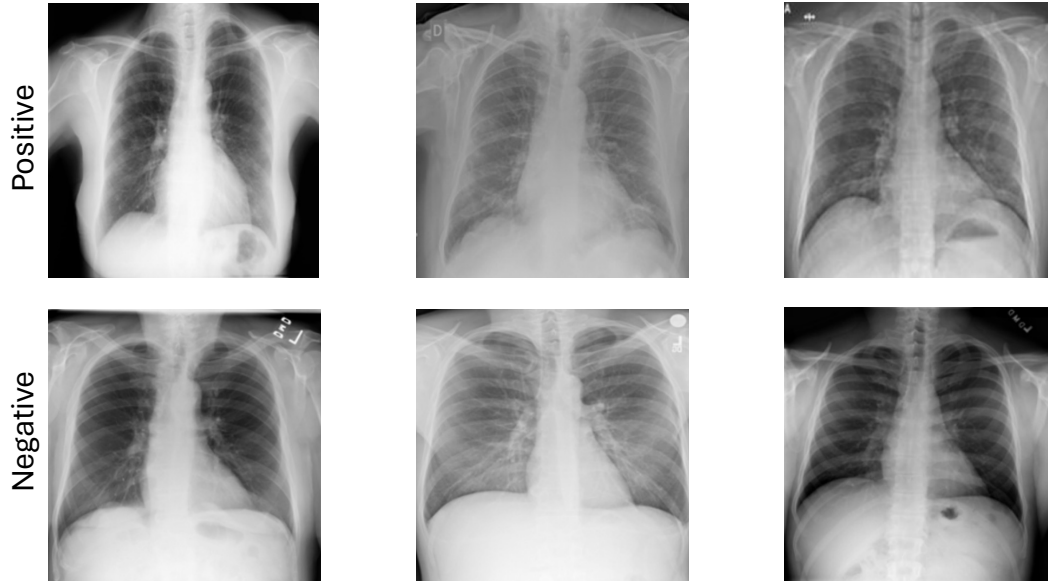


Figure 3: Data sample of Covid-19 detection.

1219 H Training Task – ECG Abnormal Detection

1220 H.1 Task Description

1221 Electrocardiography (ECG) is a crucial diagnostic tool for assessing a patient’s cardiac condition,
1222 and automatic ECG interpretation algorithms offer significant support to medical personnel given
1223 the volume of ECGs routinely performed. The task involves analyzing the PTB-XL ECG dataset to
1224 develop and evaluate automatic ECG interpretation algorithms. We offer training and test set splits
1225 to facilitate algorithm comparability and include extensive metadata on demographics, infarction
1226 characteristics, diagnostic likelihoods, and signal properties, making it a comprehensive resource for
1227 training and evaluating automatic ECG interpretation algorithms.

1228 H.2 License and Ethics

1229 The Institutional Ethics Committee approved the publication of the anonymous data in an open-access
1230 database (PTB-2020-1).

1231 H.3 Access and Preprocessing

1232 The dataset can be directly downloaded with granted permission at <https://physionet.org/content/ptb-xl/1.0.3/> or via the terminal by `wget -r -N -c -np https://physionet.org/files/ptb-xl/1.0.3/`. Raw signal data was recorded
1234 in a proprietary compressed format, encompassing the standard set of 12 leads (I, II, III, AVL, AVR,
1235 AVF, V1, ..., V6) with reference electrodes on the right arm. Corresponding metadata, including
1236 age, sex, weight, and height, was systematically gathered in a database. Each ECG record includes
1237 a report, either generated by a cardiologist or automatically by the ECG device, which was then
1238 translated into a standardized set of SCP-ECG statements (`scp_codes`). For the relevant metadata, it
1239 is saved as one row per record identified by `ecg_id`. Totally, there are 28 columns categorized into
1240 identifiers, general metadata, ECG statements, signal metadata, and cross-validation folds. Additional
1241 details such as the heart’s axis and stages of infarction (if applicable) were also documented. To
1242 ensure privacy and compliance with HIPAA standards, all personal information, including names of
1243 cardiologists and nurses, recording locations, and patient ages (with ages over 89 years reported
1244 within a 300-year range), was pseudonymized.

1246 H.4 Data Samples

1247 We provide a random data sample from the dataset and visualize it in Figure 4. Here, a positive
1248 result indicates the presence of an abnormality in the ECG signal, while a negative result represents a
1249 normal signal. All signals are 12-channel, derived from the standard set of 12 leads (I, II, III, aVL,
1250 aVR, aVF, V1, ..., V6) with reference electrodes on the right arm.

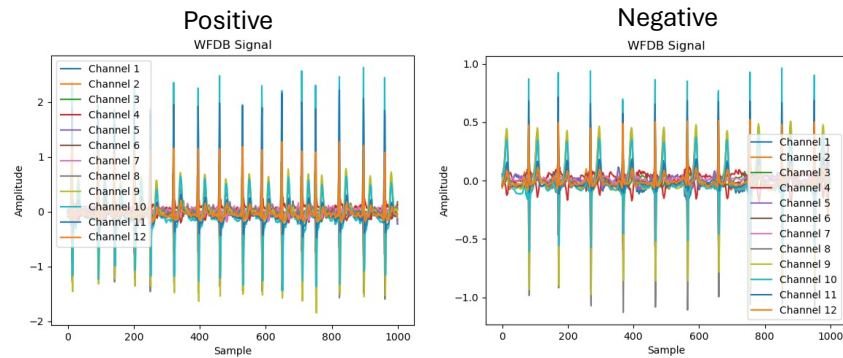


Figure 4: Data sample of ECG abnormal detection.

1251 I Training Task – Mortality Prediction

1252 I.1 Task Description

1253 The data source is MIMIC-III(Medical Information Mart for Intensive Care III), which is a substantial,
1254 anonymous, and publicly accessible repository of medical records. Each entry in the dataset contains
1255 ICD-9 codes that categorize the diagnoses and procedures conducted. In our work, we use the
1256 processed dataset to conduct the mortality prediction task.

1257 I.2 License and Ethics

1258 The dataset is available for non-profit use in accordance with the license at <https://www.physionet.org/content/mimiciii/view-license/1.4/>.

1260 I.3 Access and Preprocessing

1261 MIMIC-III can be accessed as a credentialed user on PhysioNet with an approved application at
1262 <https://mimic.mit.edu/>. In our experiment, we follow the ICU-oriented preprocessing
1263 pipeline [72] to process the data and follow the feature extraction pipeline [20] to extract dynamic
1264 features. Features extracted from this MIMIC-III database are listed in Table 5.

1265 I.4 Data Samples

1266 We provide a random data sample from the dataset and visualize it in Figure 5.

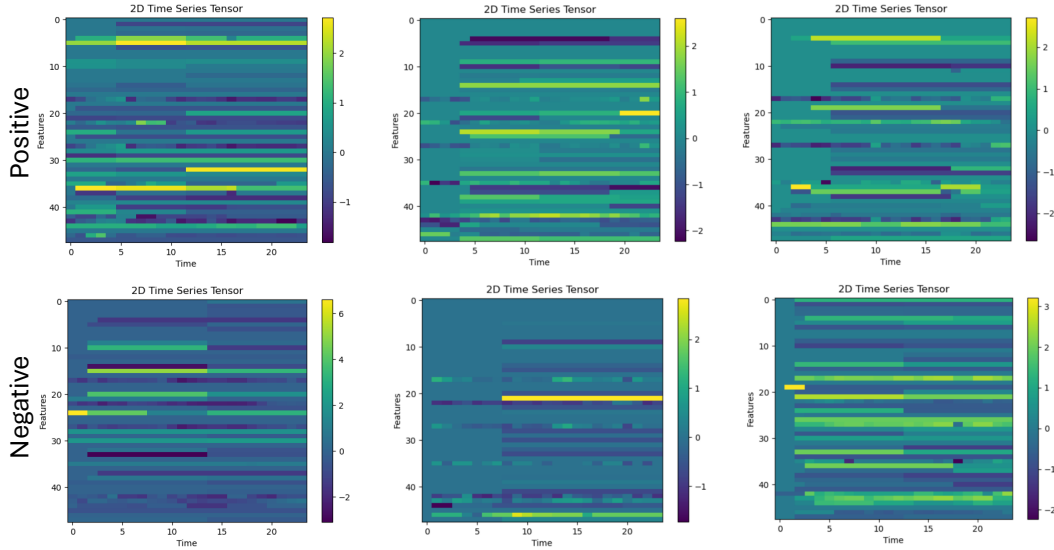


Figure 5: Data sample of mortality prediction.

Table 5: Clinical concepts extracted from MIMIC-III database [23]. The information is based on Table 15 provided in [20].

Feature	RICU	Unit
Blood pressure (systolic)	sbp	mmHg
Blood pressure (diastolic)	dbp	mmHg
Heart rate	hr	beats/minute
Mean arterial pressure	map	mmHg
Oxygen saturation	o2sat	%
Respiratory rate	resp	breaths/minute
Temperature	temp	°C
Albumin	alb	g/dL
Alkaline phosphatase	alp	IU/L
Alanine aminotransferase	alt	IU/L
Aspartate aminotransferase	ast	IU/L
Base excess	be	mmol/L
Bicarbonate	bicar	mmol/L
Bilirubin (total)	bili	mg/dL
Bilirubin (direct)	bili_dir	mg/dL
Band form neutrophils	bnd	%
Blood urea nitrogen	bun	mg/dL
Calcium	ca	mg/dL
Calcium ionized	cai	mmol/L
Creatinine	crea	mg/dL
Creatinine kinase	ck	IU/L
Creatinine kinase MB	ckmb	ng/mL
Chloride	cl	mmol/L
CO2 partial pressure	pco2	mmHg
C-reactive protein	crp	mg/L
Fibrinogen	fgn	mg/dL
Glucose	glu	mg/dL
Haemoglobin	hgb	g/dL
International normalised ratio (INR)	inr_pt	-
Lactate	lact	mmol/L
Lymphocytes	lymph	%
Mean cell haemoglobin	mch	pg
Mean corpuscular haemoglobin concentration	mchc	%
Mean corpuscular volume	mcv	fL
Methaemoglobin	methb	%
Magnesium	mg	mg/dL
Neutrophils	neut	%
O2 partial pressure	po2	mmHg
Partial thromboplastin time	ptt	sec
pH of blood	ph	-
Phosphate	phos	mg/dL
Platelets	plt	1,000 / μ L
Potassium	k	mmol/L
Sodium	na	mmol/L
Troponin T	tnt	ng/mL
White blood cells	wbc	1,000 / μ L
Fraction of inspired oxygen	fio2	%
Urine output	urine	mL

1267 J Validation Task – Enlarged Cardiomediastinum Detection

1268 J.1 Task Description

1269 This task is designed to evaluate the probability of an enlarged cardiomediastinum by using medical
1270 images from clinical assessments. It serves to gauge the model’s ability to interpret radiographs
1271 effectively. The data for this task are sourced from the CheXpert Dataset [21]. CheXpert is a collection
1272 of 224,316 chest radiographs from 65,240 patients who underwent radiographic examinations at
1273 Stanford Health Care from October 2002 to July 2017. These images were gathered from both
1274 inpatient and outpatient centers and include the associated radiology reports.

1275 J.2 Access and Preprocessing

1276 This dataset can be accessed via the link at [https://aimi.stanford.edu/](https://aimi.stanford.edu/chexpert-chest-x-rays)
1277 [chexpert-chest-x-rays](https://aimi.stanford.edu/chexpert-chest-x-rays) and downloaded via the link at [https://stanfordaimi.](https://stanfordaimi.azurewebsites.net/datasets/8cbd9ed4-2eb9-4565-affc-111cf4f7ebe2)
1278 [azurewebsites.net/datasets/8cbd9ed4-2eb9-4565-affc-111cf4f7ebe2](https://stanfordaimi.azurewebsites.net/datasets/8cbd9ed4-2eb9-4565-affc-111cf4f7ebe2).
1279 The training set comprises 224,316 high-quality X-ray images from 65,240 patients, annotated
1280 to automatically identify 14 different observations from radiology reports, reflecting the inherent
1281 uncertainties of radiographic interpretation. The validation set includes 234 images from 200 patients,
1282 each manually annotated by three board-certified radiologists. The test set, which remains unreleased
1283 to the public and is held by the organizers for final assessment, contains images from 500 patients
1284 annotated through the consensus of five board-certified radiologists. CheXpert images have an
1285 average resolution of 2828x2320 pixels.

1286 J.3 Data Samples

1287 We provide a random data sample from the dataset and visualize it in Figure 6.



Figure 6: Data sample of enlarged cardiomediastinum detection.

1288 K Validation Task – Sepsis Prediction

1289 K.1 Task Description

1290 This task focuses on predicting the likelihood of sepsis during ICU stays, assessing the model’s
1291 ability to analyze various clinical data, including lab events, diagnoses, and prescriptions. For
1292 this research, we utilize the MIMIC-III database, extracting features and cohorts through a well-
1293 established preprocessing pipeline [20].

1294 K.2 License and Ethics

1295 This dataset is governed by the license available at the following URL: <https://www.physionet.org/content/mimiciii/view-license/1.4/>.

1297 K.3 Access and Preprocessing

1298 MIMIC-III dataset can be accessed with the approved permission via <https://mimic.mit.edu/>. A random sampling strategy is applied to select a subset with 1,000 samples for testing.

1300 K.4 Data Samples

1301 We provide a random data sample from the dataset and visualize it in Figure 7. It has clinical features
1302 only, including lab events and vital signs. Although the data feature space of both mortality prediction
1303 and sepsis prediction is the same, the feature distributions are significantly different. That is why the
1304 zero-shot inference on this task performs worse than the mortality prediction, as shown in Table 4.

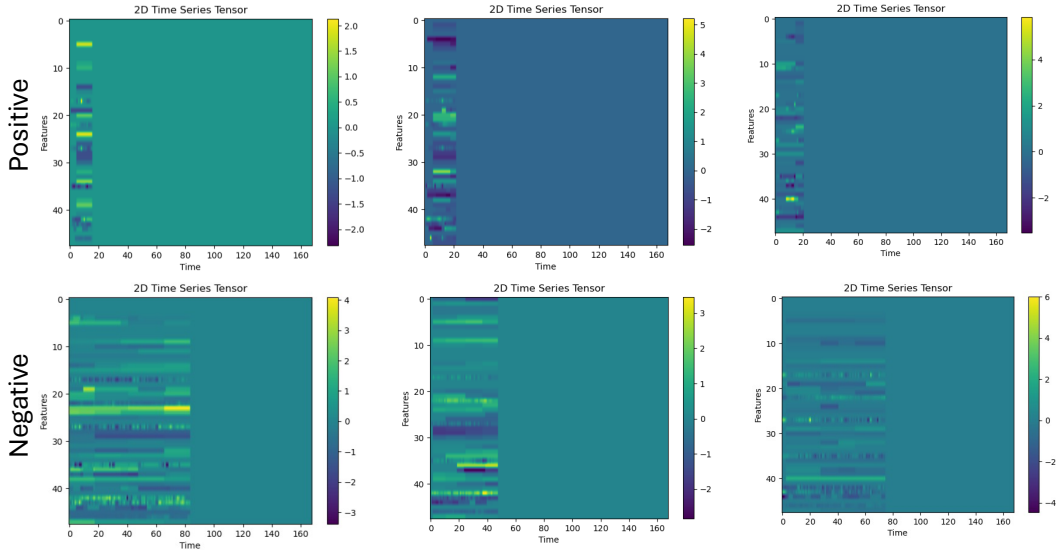


Figure 7: Data sample of sepsis prediction.

1305 L Validation Task – MedVQA

1306 L.1 Task Description

1307 The SLAKE dataset is designed for Medical Visual Question Answering (Med-VQA), integrating
1308 detailed visual and textual annotations with a medical knowledge base. It features semantic segmen-
1309 tation masks and object detection bounding boxes for each radiology image. SLAKE includes both
1310 basic clinical and complex compositional questions, and is uniquely bilingual in English and Chinese.
1311 It expands coverage to more body parts and introduces new question types related to shape and
1312 knowledge graphs, with comparative data provided against the VQA-RAD dataset. In this task, the
1313 model uses both visual context and verbal questions as inputs, requiring answers that integrate textual
1314 questions and visual context. This task tests the model’s ability to align text and image modalities in
1315 the medical domain.

1316 L.2 License and Ethics

1317 Ethical approval was not required as confirmed by the license attached with the open access data in
1318 [12].

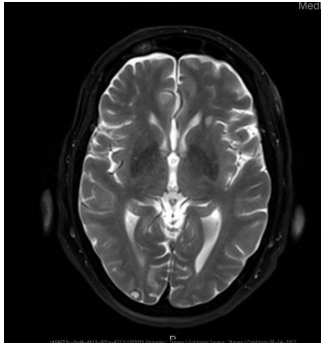
1319 L.3 Access and Preprocessing

1320 This dataset can be accessed at <https://www.med-vqa.com/slake/>. For the image part,
1321 642 images, including 12 diseases and 39 organs, were in the format of CTs and MRIs. With the help
1322 of a constructed knowledge graph, it covers questions with ten different content types and semantic
1323 labels proposed by doctors. We randomly select 1000 samples from the test dataset.

1324 L.4 Data Samples

1325 We provide a random data sample from the dataset. Question-answering pair and corresponding
1326 image (Figure 9) are listed below. This MedVQA dataset contains different types of images except
1327 for chest X-ray images, which are different from the ones we used in the model training. In addition,
1328 the trained FEDMEKI does not use any medical question-answering training tasks. Therefore, its
1329 performance of this “new” task is limited, as shown in Table 4.

Q: How can you tell this is a
T2 weighted image?
A: CSF is white.



Q: Is the heart enlarged?
A: No.



Q: What is causing the widening?
A: Mass

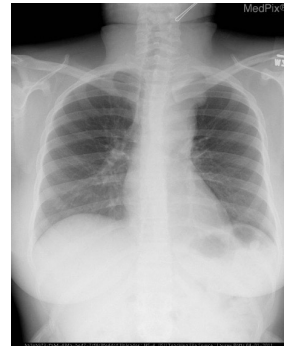


Figure 8: Data sample of MedVQA task.

1330 M Validation Task – Signal Noise Clarification

1331 M.1 Task Description

1332 This task is dedicated to precisely characterizing noise in ECG signals through a question-and-answer
1333 format.

1334 M.2 License and Ethics

1335 The Institutional Ethics Committee approved the publication of the anonymous data in an open-access
1336 database (PTB-2020-1).

1337 M.3 Access and Preprocessing

1338 It utilizes data from an established ECG question answering dataset [22] and a related ECG
1339 database [19], which can be accessed via <https://physionet.org/content/ptb-xl/1.0.3/> and <https://github.com/Jwoo5/ecg-qa/tree/master/ecgqa/ptbx1>. The
1340 ECG signals used in this task consist of 12 channels and have a duration of 10 seconds, mirroring the
1341 parameters used in the ECG Abnormal Detection task. We randomly sample 1,000 ECG-question
1342 pairs as the validation data.
1343

1344 M.4 Data Samples

1345 We provide a random data sample from the dataset. Question-answering pair and corresponding
1346 signal (Figure 9) are listed below. Although we have a training task on ECG, the ECG abnormal
1347 detection task is different from this one. This task aims to answer the noise types of ECG signals
1348 according to the input ECG and the question. We can see that the ECG signals in Figure 9 are quite
1349 different from the ones in Figure 4, which increases the difficulty of this task significantly.

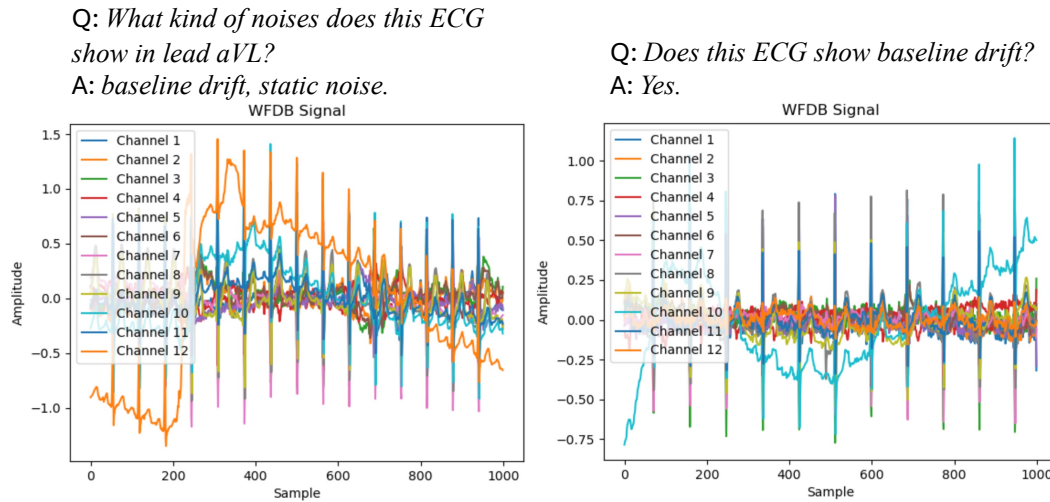


Figure 9: Data sample of signal noise clarification.

N Implementation Details

N.1 Model Details

For each local model \mathbf{W}_n deployed in client C_n , we implement their modality-specific encoders and task-specific decoders. Details about encoders for different modalities can be found in Table N. As all training tasks can be categorized as binary classification, we use MLPs as their task-specific decoders, where decoders for different tasks do not share any parameter.

Table 6: Details of modality-specific encoders.

Modality	Encoder	# of Parameters
Image	Deit-tiny [73]	7.8M
Signal	CNN [74]	4.1M
Diagnosis Code	Transformer [75]	3.7M
Procedure Code	Transformer [75]	3.7M
Drug	Transformer [75]	3.7M
Lab Results	Transformer [75]	3.7M

N.2 Optimizer Hyperparameters

We leverage Adam optimizer [76] for optimizing both local model \mathbf{W}_n and foundation model \mathcal{F} . The number of communication rounds is set to 10. For local model \mathbf{W}_n , we find the learning rate of $1e-4$ for local models achieves a decent convergence, while the learning rate for the foundation model is configured to $5e-4$. The batch size of both the foundation model and local models is set to 64.

N.3 Task Prompts

Task prompts for classification tasks are listed in Table 7. For MedVQA and signal noise clarification, prompts are questions themselves.

Table 7: Task Prompts.

Task Name	Prompt
Lung Opacity Detection	Assess this CT image: should it be classified as lung opacity?
Covid-19 Detection	Based on this image, is the patient COVID-19 positive?
ECG Abnormal Detection	Is the given ECG abnormal?
Mortality Prediction	Based on these clinical features, will mortality occur in this patient?
Enlarged Cardiomeastinum Detection	Does this image show evidence of enlarged cardiomeastinum?
Sepsis Prediction	Based on these clinical features, will sepsis occur in this patient?

N.4 Baselines

To better understand the benchmarks used in the experiments, we use visualizations to demonstrate each approach clearly.

For single-task evaluation, we use FedAvg_s/FedProx_s (Figure 10), FedAvg_s⁺/FedProx_s⁺ (Figure 11), FedAvg_s^{*}/FedProx_s^{*}, and FedAvg_s^F/FedProx_s^F (Figure 12). When training single tasks, we only use each task data as the model input. For multi-task training, we also have eight baselines that are shown from Figure 13 to Figure 15. These models will train all the task data together.

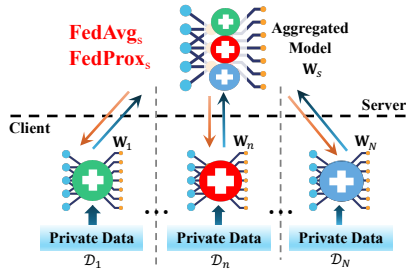


Figure 10: FedAvg_s or FedProx_s

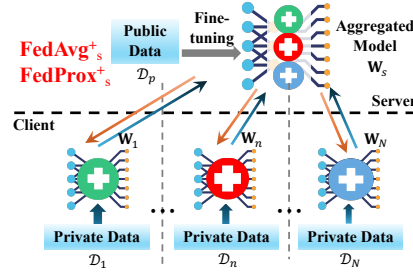


Figure 11: FedAvg_s^+ and FedProx_s^+

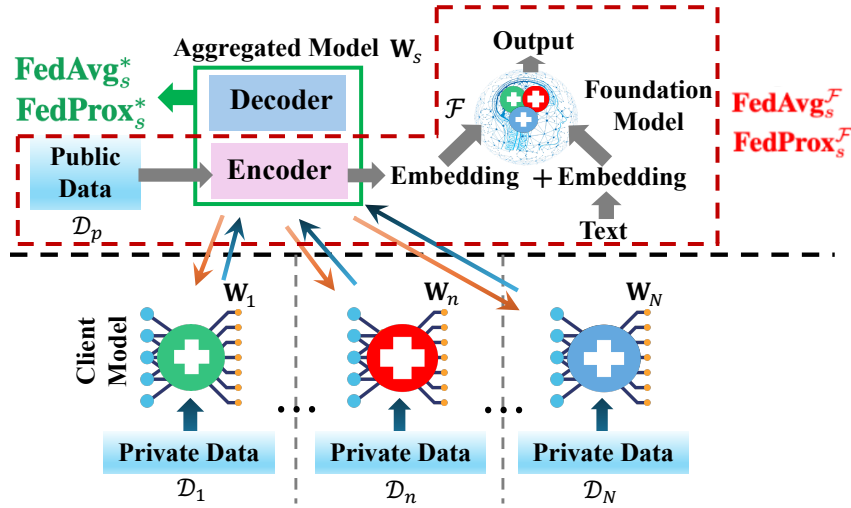


Figure 12: $\text{FedAvg}_s^F / \text{FedProx}_s^F$ (red dot line) and $\text{FedAvg}_s^* / \text{FedProx}_s^*$ (green line).

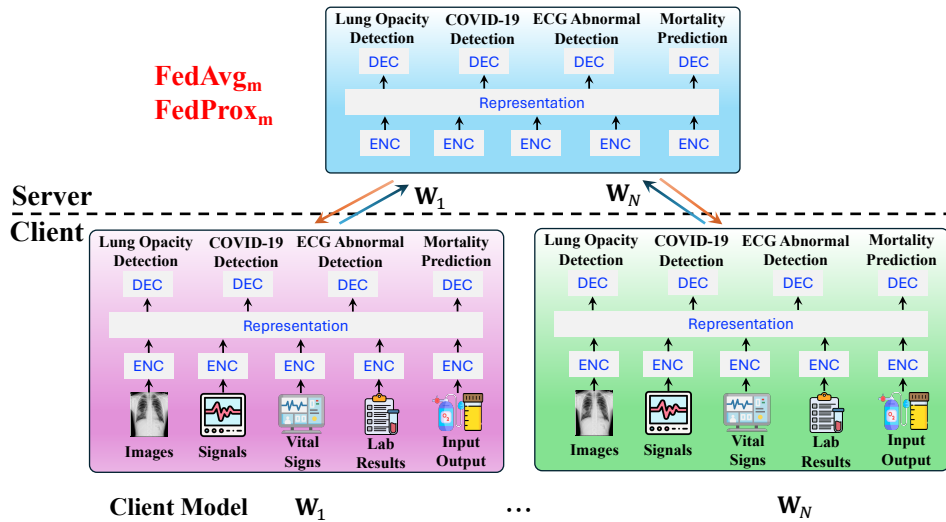


Figure 13: FedAvg_m or FedProx_m

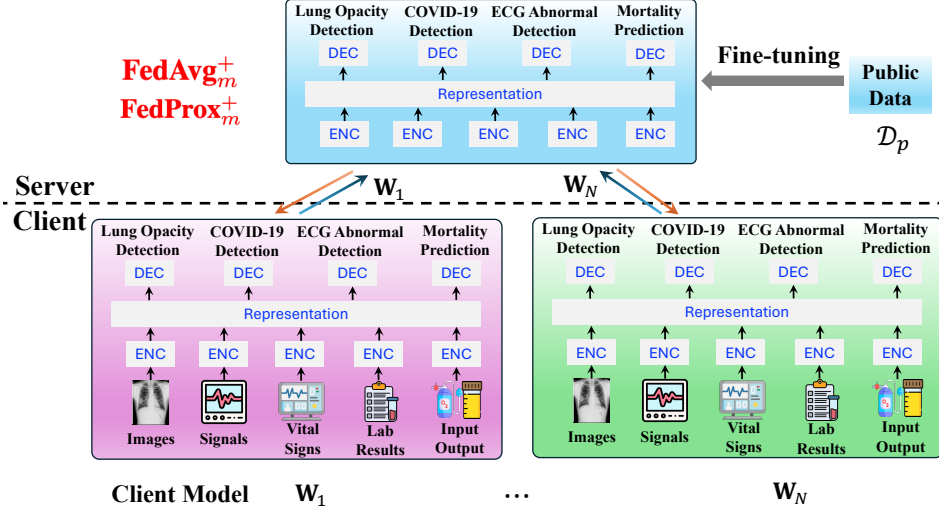


Figure 14: **FedAvg_m⁺** and **FedProx_m⁺**

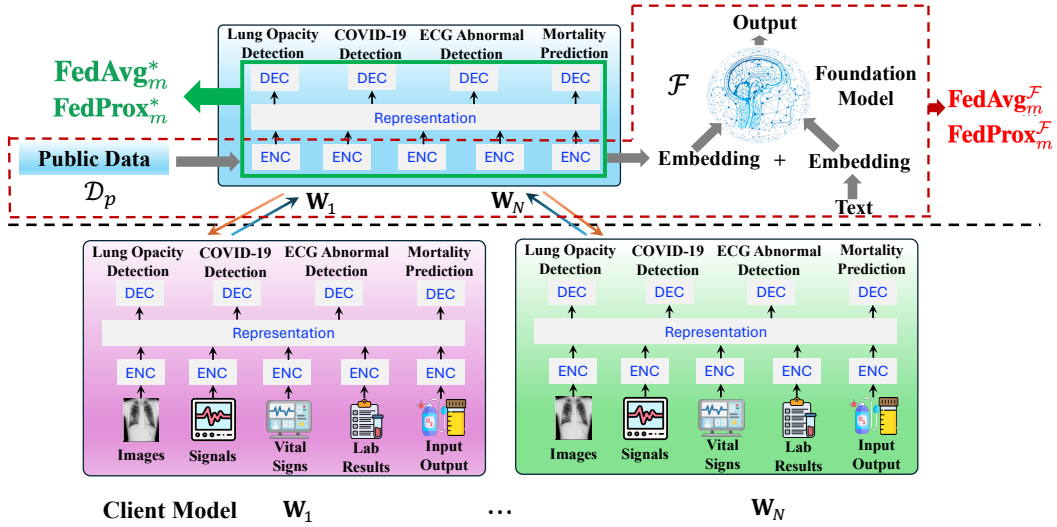


Figure 15: **FedAvg_m^F/FedProx_m^F** (red dot line) and **FedAvg_m^{*}/FedProx_m^{*}** (green line).