

OPEN ACCESS

EDITED BY

Vishwa S. Parekh, University of Texas Health Science Center at Houston, United States

REVIEWED BY

Pranav Kulkarni, University of Maryland, United States Guangyao Zheng, Johns Hopkins University, United States

*CORRESPONDENCE

Qiong Li

☑ liqiong@sysucc.org.cn Xiangmeng Chen ☑ 3897001254@qq.com

RECEIVED 16 July 2025 ACCEPTED 09 September 2025 PUBLISHED 26 September 2025

CITATION

Chen Y, Liu L, Feng B, Chen Y, Xu J, Lin H, Li K, Chen X, Ke Y, Zhou H, Hu Q, Jin Q, Long W, Li Q and Chen X (2025) A meta-learning-based robust federated learning for diagnosing lung adenocarcinoma and tuberculosis granulomas. Front. Oncol. 15:1666937. doi: 10.3389/fonc.2025.1666937

COPYRIGHT

© 2025 Chen, Liu, Feng, Chen, Xu, Lin, Li, Chen, Ke, Zhou, Hu, Jin, Long, Li and Chen. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.

A meta-learning-based robust federated learning for diagnosing lung adenocarcinoma and tuberculosis granulomas

Yuyao Chen¹, Lei Liu², Bao Feng², Yehang Chen², Jun Xu², Huan Lin³, Kunwei Li⁴, Xiaodong Chen⁵, Yuting Ke⁵, Haoyang Zhou², Qinghui Hu², Qinggeng Jin⁶, Wansheng Long⁷, Qiong Li^{8*} and Xiangmeng Chen^{1*}

¹Nanxishan Hospital of Guangxi Zhuang Autonomous Region, Guilin, China, ²Guilin University of Aerospace Technology, Guilin, China, ³Guangdong Provincial People's Hospital, Guangzhou, China, ⁴Fifth Affiliated Hospital of Sun Yat-Sen University Department of Radiology, Zhuhai, China, ⁶Affiliated Hospital of Guangdong Medical University Department of Radiology, Zhanjiang, China, ⁶Guangxi University School of Electrical Engineering, Nanning, China, ⁷Jiangmen Central Hospital, Jiangmen, China, ⁸Sun Yat-Sen University Cancer Center Department of Radiotherapy, Guangzhou, China

Background: Differentiating between lung adenocarcinoma (LAC) and tuberculosis granuloma (TBG) of solitary pulmonary solid nodules (SPSNs) based on CT images alone is a daunting task for clinical diagnosis. Thus, it is crucial to fully utilize CT imaging data to explore effective noninvasive diagnostic methods to improve the identification of TBG and LAC.

Purpose: This study aimed to leverage CT imaging datasets from multiple hospitals for the diagnosis of TBG and LAC in SPSNs. It achieved this by deploying a meta-learning method within a federated learning framework while protecting data privacy.

Methods: A total of 1,026 patients, along with their CT images of solitary pulmonary solid nodules (SPSNs) and corresponding clinical data, were collected from six medical institutions. Subsequently, the data from these six institutions were systematically partitioned into five cohorts. Each cohort was divided into two parts: the training set and the test set. A meta-learning-based robust federated learning model by training set data was proposed to construct personalized federated learning signatures (PFLS) without uploading raw data from each medical institutions. Receiver operating characteristic curve (ROC), area under curve (AUC), decision curve analysis (DCA), net reclassification improvement (NRI) and integrated discrimination improvement (IDI) are used to analyze the performance of the PFLS.

Results: The PFLS trained by the proposed meta-learning-based robust federated learning framework shows superior performance compared to alternative methods. The AUC range on the training sets of the five cohorts is 0.866-0.939, AUC range on the testing sets is 0.808-0.927). The significant difference of AUC between the proposed method and the clinical model was demonstrated by the NRI and IDI. The decision curves indicated a higher net benefit of our proposed method.

Conclusion: The PFLS mitigates overfitting issues arising from limited sample size in local hospitals. It also alleviates the problem that a single global model is not applicable to all hospitals due to the heterogeneity of data distribution among different hospitals.

KEYWORDS

lung adenocarcinoma, tuberculosis granuloma, solitary pulmonary solid nodules, SPSNs, meta-learning, federated learning, CT images, personalized federated learning signatures

Introduction

The prevalence of CT has been led to a significant upsurge in the detection rate of Solitary Pulmonary Solid Nodules (SPSNs) (1). Clinically, SPSNs can be bifurcated into benign and malignant categories. Lung Adenocarcinoma (LAC) is the most common pathological type of malignant SPSNs, while Tuberculosis Granulomas (TBG) is the common pathological type of benign SPSNs (2, 3). However, the treatment regimens and clinical outcomes for lung adenocarcinoma and tuberculous granulomas are entirely different. Radical surgical resection is the preferred treatment for the former, while the latter is often managed with anti-tuberculosis medications (4). Misdiagnosis can lead to uncontrollable disease progression and a poor prognosis in patients with lung adenocarcinoma. Conversely, it may also result in overtreatment for those with tuberculous granulomas (5).

Although CT scans can identify SPSNs, the differentiation between LAC and TBG based on CT images alone presents a daunting task for clinical diagnosing. This is primarily because LAC and TBG both exhibit similar lobulated and spiculated features, and there is a lack of effective contrast agents to aid in distinguishing TBG from LAC (6, 7). Most patients with SPSNs detected by CT undergo biopsy diagnosis to guide the treatment plan. However, when the lesion is small and difficult to locate, the difficulty and related risks increase significantly (8, 9). Consequently, it is crucial to fully utilize CT imaging data to explore new effective non-invasive diagnostic methods to improve TBG and LAC identification.

Deep learning, as a data-driven technology for model performance, has shown great potential in image classification. Previous studies have demonstrated that deep learning models can extract features from raw medical images at various levels of abstraction (10, 11). Applying deep learning techniques to computer-aided diagnostic systems holds promise for improving the accuracy of TBG and LAC differentiation. However, due to the need for medical data privacy protection, medical centers are generally not allowed to share data, which limits the scale of the data. Unfortunately, robust and accurate deep learning models require a large amount of data for training; otherwise, overfitting is prone to occur, leading to a decline in the generalization ability of deep learning models.

Federated learning facilitates multi-clients collaborative training by aggregating local model parameters of each client into the shared global model, without sharing data from different clients (12). This approach fully utilizes information of each hospital without sharing raw CT image data, thus addressing privacy concerns and limiting overfitting. The federated averaging algorithm of most federated learning methods weights the parameters of each local model according to the sample sizes of different medical institution (13, 14). However, Additionally, data heterogeneity caused by differences in data collection across medical institutions (such as scanning equipment, imaging parameters, population characteristics, etc.) significantly restricts the performance of federated learning models in multi-medical institution medical image analysis (15, 16). Therefore, when there are differences in the data distributions across multiple centers, it is challenging for a single global model obtained merely by aggregating the parameters of each local model to perform consistently well across all centers (17, 18).

In this paper, a meta-learning-based robust federated learning approach is proposed to leverage heterogeneous CT imaging datasets from multiple medical institutions for the diagnosis of TBG and LAC in SPSNs. The reptile algorithm of meta-learning is deployed to aggregate gradients of parameters of each local model. This improves the performance and robustness of the global model on data from each local medical institution. Finally, each center fine-tunes the global model based on local data to complete model personalization.

Materials and methods

Patients

This retrospective study was approved by the Institutional Review Boards of the participating hospitals, with a waiver of informed consent. Detailed inclusion and exclusion criteria are provided in Supplementary S1. Finally, a total of 1,026 samples from six medical institutions. Since one medical institution has only 17 cases, we merged the data of this medical institution into another hospital, so there are a total of 5 cohorts. These five cohorts include:

cohort 1 (2014–2020): 270 patients (training set: 161; test set: 109), cohort 2 (2013–2016): 87 patients (training set: 51; test set: 36), cohort 3 (2014–2019): 119 patients (training set: 70; test set: 49), cohort 4 (2011–2020): 471 patients (training set: 282; test set: 189), cohort 5 (2018–2020): 79 patients (training set: 46; test set: 33).

CT image acquisition and evaluation

Chest CT images were acquired from six different scanners (Siemens, Toshiba, GE, Philips) with patients in the supine position, covering the entire chest from the thoracic inlet to the adrenal glands during a breath-hold. Scans were performed in spiral mode with a tube voltage of 120 kVp and automatic mAs adjustment. Images were reconstructed with standard and high-resolution algorithms at 1.0–3.0 mm slice thickness and 0.8–3.0 mm interslice gap. Two independent chest radiologists, blinded to clinical information, assessed the images using lung and mediastinal window settings, evaluating nodule location, size, margin, lobulation, and spiculation; discrepancies were resolved by consensus. Detailed information is provided in Supplementary S2.

Pathological diagnosis

All samples were fixed in formalin and subsequently stained with hematoxylin and eosin (HE). The experienced pathologists performed the pathological analysis of the surgical specimens in accordance with the 2011 International Association for the Study of Lung Cancer/American Thoracic Society/European Respiratory Society classification system, and the 2015 World Health Organization (WHO) classification of lung neoplasms (19, 20). These pathologists were blinded to the CT findings.

Image preprocessing

For neural network processing, preprocessing operations are applied to the CT images. An experienced radiologist utilizes a rectangular bounding box to crop the region of interest (ROI) from each CT slice initially. All ROIs are then interpolated and standardized to 224×224 pixels. Next, the ROIs from three sequential single-channel CT slices for the same patient are merged to form a three-channel image with the dimensions 224×224×3. Finally, these three-channel images are used as input data for the neural network. Detailed information is provided in Supplementary Figure S1.

Building the meta-learning-based personalized federated learning signature

In order to adapt to the data situation of each medical institution, we train a personalized federated learning signature

for each medical institution. This usually involves three steps: feature extraction, feature selection, and classifier training.

During the feature extraction process, a federated learning based on model agnostic meta-learning is used to extract the CT features of each hospital. The entire training process of federated learning encompassed three stages.

FedAvg stage: In the initial iteration of the FedAvg stage, both local and global models start with identical parameters pre-trained on ImageNet. Each local client trains its model using its own dataset. After all local clients complete training, they upload their model gradients to the global server. The global server aggregates these gradients by weighting them according to each client's sample size relative to the total samples across all clients. The global model is then updated using these weighted gradients and distributed back to the local clients as the initial parameters for the next iteration. This process repeats for several iterations, and the final global model parameters are passed to the Reptile stage.

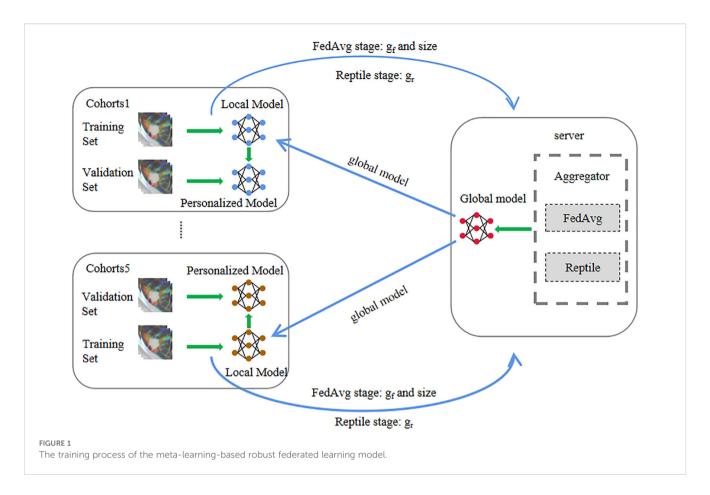
Reptile stage: Unlike the FedAvg stage, the Reptile stage employs the Adam optimizer for local model updates. The aggregation method also differs: instead of sample-size-based weighting, the global server treats each local client as a distinct meta-learning task and applies the Reptile algorithm to compute the combined gradient direction. The global model is then updated with momentum based on this aggregated gradient. After multiple iterations, the final global model parameters are delivered to each local client for the subsequent personalized stage.

Personalized stage: In the Personalize stage, the local clients do not share any data to the global server, and only fine-tune the local models with their own data sets based on the stochastic gradient descent algorithm. And the initial parameters of the local models are the final parameters of the global model of the Reptile stage.

During the whole training process, the raw data of a local client or hospital is never shared with the global server and other local clients, which ensures the security and privacy of the local data. The global server performs aggregation operations on the parameters of the local model so that the local clients can share the training results, effectively avoiding overfitting when the data samples of a single client are too small. The Reptile stage creates ideal conditions for rapid fine-tuning of the local model, and the Personalize stage can effectively solve the problem of data heterogeneity among different hospitals or centers.

More detailed information regarding the model and training details can be found in Supplementary S3 and Figure 1. The local hospitals utilizes the robust personalized local models, trained by the proposed method, to extract 3904 features from the CT images at each layer of the ROI. Subsequently, the features from all layers are fused (refer to Supplementary S4 for detailed information).

The classifier can utilize numerous features obtained from the above operations to diagnose TBG or LAC. However, most of these features are not conducive for diagnosing pulmonary nodules and may introduce noise, negatively affecting diagnostic accuracy. Therefore, in the process of feature selection and classifier training, the Mann-Whitney U test is employed to evaluate the diagnostic significance of features, retaining only those with a p-



value < 0.05. Finally, a Bayesian extreme learning machine is employed to building the personalized federated learning signature (PFLS) using the selected features (21). To validate its effectiveness, we further performed systematic comparisons with several widely used classifiers, including logistic regression (LR), support vector machine (SVM), and random forest (RF). These models are representative in medical image analysis and AI classification tasks, covering linear, kernel-based, and ensemble learning approaches, respectively. All models were trained and evaluated under identical data splits and preprocessing settings to ensure fair comparison. The pseudocode of the algorithm is provided in Supplementary S5.

Personalized federated learning signature comparison FedAvg model

The FedAvg (22) model is the federated learning model and the global model trained through the federated averaging algorithm. In the training iteration of FedAvg, each client accepts the global model parameters, initiating local training based on this global model. After training by the local clients on local data, the parameters of the local models are uploaded, and the parameters of each local model are averagely weighted by the global server to achieve collaborative training of the ResNet18 global model. More model parameter Settings can be found in Supplementary S6. Subsequently, each local hospital uses the same global model to

extract 3904 features from the CT images, followed by utilization of Mann-Whitney U test to select features with significant difference from the extracted features. Finally, a Bayesian Extreme Learning Machine is employed for classification.

Personalized federated learning signature comparison independent local models

Independent local models(ILM) are the ResNet18 and trained exclusively with local data, with no data interaction occurring among the local models from other hospitals. The training process of this model first involves pre-training using ImageNet data, and then training respective models with each local dataset. More model parameter Settings can be found in Supplementary S5. Through the ResNet18 model, 3904 deep learning features are extracted from the CT image of each case. Features with significant differences are identified using the Mann-Whitney U test. Eventually, a Bayesian Extreme Learning Machine is applied to perform classification using these selected features.

Personalized federated learning signature comparison building the clinical model

CT image data are collected from a total of five cohorts, with data from each hospital divided into training and test sets. The

patient distribution and clinical features of the CT images are outlined in Table 1. Thus, this study selects clinical features (gender, age, nodule size, shape of lesion margin, lobulated shape, and spiculated sign to build the clinical model(CM) based on Bayesian Extreme Learning Machine.

Personalized federated learning signature comparison merged data centralized model

To validate the necessity and advantages of the proposed Personalized Federated Learning Signature (PFLS) framework, we established a Merged Data Centralized Model (MDCM) as a comparative benchmark. This model integrates training data from all participating centers to train a single deep learning model without any privacy constraints—simulating an ideal scenario where data sharing faces no regulatory or ethical barriers. After training, the centralized MDCM was independently evaluated on the local test sets of each hospital to assess its generalization performance across heterogeneous data distributions. This approach enables a quantitative comparison between the centrally trained model and the personalized federated models, highlighting the impact of data heterogeneity and demonstrating the effectiveness of federated learning in maintaining model performance while preserving data privacy.

Personalized federated learning signature comparison with personalized federated model

To further evaluate the effectiveness of the proposed PFLS framework, we selected several representative personalized federated learning methods for comparison, including FedProx (23), FedBN (24), and Moon (25). FedProx introduces a proximal term into the local objective function to constrain local updates from deviating excessively from the global model, thereby stabilizing the optimization process under non-IID data distributions. FedBN retains the Batch Normalization (BN) parameters locally while aggregating the remaining parameters globally, which alleviates performance degradation caused by feature distribution shifts across centers. Moon incorporates a contrastive learning objective during local training to encourage consistency between local and global representations, thus improving robustness in heterogeneous data scenarios. After federated training, each method employed its respective personalized model to extract features from the local data. Finally, a Bayesian Extreme Learning Machine is employed for classification.

Ablation experiments on PFLS

To quantitatively verify the effectiveness of the Reptile step (26) and the personalization step in PFLS, we designed ablation

experiments. Specifically, we constructed different algorithm variants by selectively removing these two steps: (1) removing the Reptile step while retaining the personalization step, where each site was validated using its own personalized model; (2) removing the personalization step while retaining the Reptile update, where all centers were validated using the global model after the Reptile update. All variants were trained and evaluated under the same experimental settings. By comparing their performance with the complete PFLS, we were able to assess the contribution of each component.

Statistical analysis

The performance evaluation of the models involved calculating various metrics, including the receiver operating characteristic curve (ROC), area under the curve (AUC), sensitivity, specificity, accuracy, positive probability value (PPV), and negative probability value (NPV). The net reclassification improvement (NRI) and integrated discrimination improvement (IDI) were used to measure the degree of improvement of PFLS in overall discriminative ability compared with FedAvg, ILM, and CM. P-values less than 0.05 were considered a significant difference.

Results

Clinical factors and subjective CT findings analysis

The patient distribution and clinical features of the CT images are outlined in Table 1. The table details various clinical parameters such as gender, age, lesion size, location, margin, lobulated shape, and spiculated sign, with a clear distinction between training and testing sets within each cohort. A notable observation is the inconsistent distribution of these clinical features across different cohorts. For instance, the proportion of males and females varies significantly, with some cohorts having a higher male prevalence (e.g., Cohort 1 and Cohort 4) while others show a more balanced or female-dominant distribution (e.g., Cohort 3). Similarly, the distribution of lesion location, margin, lobulated shape, and spiculated sign further underscores the heterogeneity among cohorts. For example, the presence of lobulated and spiculated lesions varies widely, suggesting differences in disease characteristics or diagnostic practices across cohorts.

The performance of PFLS identifies LAC and TBG

As shown in Table 2, all models were trained and validated using the same feature set to ensure fairness in comparison. The results revealed that SVM and RF experienced severe overfitting during training, as indicated by the large performance gap between the training and test sets. In contrast, logistic regression (LR) did

frontiersin.org

TABLE 1 The patient distribution and clinical features.

			Coh	ort 1					Coh	ort 2					Coh	ort 3		
Clinical information	Trai	ning set (n=	161)	Test	ting set (n=1	.09)	Tra	ining set (n=		Tes	sting set (n=	36)	Trai	ining set (n	=70)	Tes	ting set (n=	49)
Cunical information	TBG	LAC		TBG	LAC		TBG	LAC		TBG	LAC		TBG	LAC		TBG	LAC	
	(n=34)	(n=127)	P value	(n=23)	(n=86)	P value	(n=14)	(n=37)	P valule	(n=10)	(n=26)	P valule	(n=21)	(n=49)	P valule	(n=15)	(n=34)	P valule
Gender																		
Male	21	57	0.08	12	39	0.56	9	20	0.51	5	16	0.529	14	25	0.227	6	24	0.043
Female	13	70	0.08	11	47	0.56	5	17	0.51	5	10	0.529	7	24	0.227	9	10	0.043
Age(mean ± SD, years)																		
	59.2 ±13.8	57.8 ±10.1	0.533	53.0 ±11.5	59.8 ±10.9	0.011	53.8 ±12.1	59.4±9.3	0.087	61.2±9.2	57.8 ±12.2	0.433	47.9 ±15.6	58.7 ±9.2	0.001	50.7 ±11.5	60.7 ±9.7	0.003
Lesion size(mm)																		
	14.5±8.0	20.7±9.7	0.001	12.7±7.4	20.3±9.1	<0.001	18.7 ±6.00	20.9 ±13.3	0.555	20.5±7.9	25.8 ±22.1	0.467	12.1±8.0	14.1 ±5.1	0.206	12.3±5.6	15.5 ±4.8	0.046
Location																		
Upper and middle	17	89	0.028	13	56	0.447	10	27	0.912	6	16	0.932	12	26	0.753	9	25	0.344
Lower	17	38	0.020	10	30	0.117	4	10	0.512	4	10	0.552	9	23	0.755	6	9	0.511
Lesion Margin																		
Irregular	20	117	<0.001	16	71	0.168	12	33	0.731	8	24	0.293	11	42	0.003	9	33	0.001
Regular	14	10		7	15		2	4		2	2	1.07	10	7		6	1	
Lobulated shape																		
Absence	22	18	<0.001	12	24	0.028	5	4	0.037	3	0	0.004	11	8	0.002	8	1	<0.001
Presence	12	109		11	62		9	33		7	26		10	41		7	33	
Spiculated sign																		
Absence	28	70	0.004	20	51	0.013	10	13	0.02	9	7	0.001	18	27	0.014	14	13	<0.001
Presence	6	57	0.001	3	35	0.015	4	24	0.02	1	19	0.001	3	22	0.011	1	21	10.001
			Coh	ort 4					Coh	ort 5								
Clinical information	Trair	ning Set (n=	282)	Test	ting Set (n=1	.89)	Tra	ining Set (n=	:46)	Tes	sting Set (n=	33)						
Curiicat ii ii Oi ii iatioii	TBG	LAC	p valule	TBG	LAC		TBG	LAC	p valule	TBG	LAC	p valule						
	(n=96)	(n=186)	p valule	(n=64)	(n=125)	p valule	(n=22)	(n=24)	p valule	(n=16)	(n=17)	p valule						
Gender																		
Male	61	93	0.03	42	53	0.003	13	11	0.369	10	7	0.221						

Chen et al.

			Coh	ort 4					Coh	ort 5						
	Trair	ning Set (n=	282)	Test	ing Set (n=1	.89)	Trai	ning Set (n=	:46)	Tes	sting Set (n=	33)				
Clinical information	TBG	LAC		TBG	LAC		TBG	LAC		TBG	LAC					
	(n=96)	(n=186)	p valule	(n=64)	(n=125)	p valule	(n=22)	(n=24)	p valule	(n=16)	(n=17)	p valule				
Gender																
Female	35	93		22	72		9	13		6	10					
Age(mean ± SD, years)																
	51.4 ±12.9	60.8±9.8	<0.001	52.3 ±11.3	60.5 ±10.7	<0.001	53.3±9.2	63.1±8.6	0.001	47.3 ±11.9	59.8±9.9	0.002				
Lesion size(mm)																
	12.3±6.6	17.6±8.2	<0.001	13.1±7.5	17.4±8.5	0.001	10.0±6.4	17.3±9.8	0.005	8.6±3.8	25.1 ±25.4	0.015				
Location		_		J							_	<u>'</u>		'		
Upper and middle	65	125	0.932	44	80	0.515	11	16	0.251	8	14	0.049				
Lower	31	61	0.932	20	45	0.313	11	8	0.231	8	3	0.049				
Lesion Margin																
Irregular	53	160	< 0.001	35	109	<0.001	14	22	0.021	6	17	<0.001				
Regular	43	26	10.001	29	16	10.001	8	2	0.021	10	0	10.001				
Lobulated shape																
Absence	60	32	< 0.001	37	21	<0.001	11	12	1	12	5	0.009				
Presence	36	154		27	104		11	12		4	12					
Spiculated sign																
Absence	80	98	< 0.001	55	65	< 0.001	12	14	0.796	13	11	0.286				
Presence	16	88		9	60		10	10		3	6					

not show obvious overfitting; however, its classification performance was still inferior to that of the Bayesian extreme learning machine (Bayesian ELM). By comparison, Bayesian ELM not only maintained high training efficiency but also demonstrated stronger generalization ability on the test set. The performance of the personalized federated learning signature (PFLS) across the five cohorts is presented in Figures 2, 3 and Table 3. On the test sets, PFLS achieved AUCs of 0.846 (95% CI, 0.748-0.944), 0.889 (95% CI, 0.771-1.000), 0.922 (95% CI, 0.845-0.999), 0.876 (95% CI, 0.825-0.927), and 0.893 (95% CI, 0.770-1.000), with corresponding accuracies of 0.807, 0.861, 0.878, 0.788, and 0.818. These results indicate that PFLS consistently demonstrated excellent predictive performance across different cohorts, effectively distinguishing lung adenocarcinoma from tuberculosis granulomas. Furthermore, decision curve analysis (Figure 4) showed that PFLS provided a higher net benefit than other models in cohorts 2, 3, and 5. Supplementary Table S1 in the supplementary materials for model details.

Comparison of the PFLS with the FedAvg model

As shown in Table 3, the AUC of the FedAvg model on the test sets from the five cohorts are 0.740 (95% CI, 0.637-0.843), 0.823 (95% CI, 0.672-0.974), 0.839 (95% CI, 0.724-0.955), 0.737 (95% CI, 0.664-0.810), and 0.728 (95% CI, 0.537-0.919). The FedAvg model achieved accuracy of 0.560, 0.694, 0.837, 0.698, and 0.667 on the test sets of each cohort. Compared with FedAvg, the AUC of PFLS at each cohort increased by 6.6%-16.5%, and the accuracy rate increased by 4.1%-25.1%. Supplementary Table S2 and Supplementary Table S3 respectively present the results of NRI and IDI. The results of NRI and IDI show that the performance of PFLS on each central test set is improved compared with the FedAvg model. The range values of NRI at each cohort were 0.287 to 1.263. The range values of IDI at each cohort were 0.016 to 0.069. Supplementary Table S4 in the supplementary materials for FedAvg model details.

Comparison of the PFLS with ILM

The AUC of ILM on the test sets from each cohort are 0.773 (95% CI, 0.672-0.873), 0.808 (95% CI, 0.640-0.975), 0.857 (95% CI, 0.745-0.969), 0.778 (95% CI, 0.707-0.848), and 0.824 (95% CI, 0.676-0.972). The ILM achieved accuracy of 0.716, 0.833, 0.735, 0.698, and 0.636 on the test sets of each cohort. Compared to ILM, the PFLS achieved an AUC improvement of 6.5%- 9.4% across the cohorts, and the PFLS showed an improvement of 2.8%-19.1% in the prediction accuracy for LAC and TBG at each cohort. More details are provided in Table 3 and Figure 2, 3. As shown in Supplementary Table S2 and Supplementary Table S3, the performance of PFLS on each central test set is improved compared with the ILM. The range values of NRI at each cohort were 0.324 to 969. The range values of IDI at each cohort were 0.022

to 0.061. Supplementary Table S5 in the supplementary materials for ILM details.

Comparison of the PFLS with the CM

As shown in Table 3, the AUC of the CMmodel on the test sets from the five cohorts are 0.692(95% CI, 0.569-0.815), 0.719(95% CI, 0.523-0.915), 0.759(95% CI, 0.618-0.900), 0.679(95% CI, 0.587-0.770), and <math>0.665(95% CI, 0.471-0.860). In comparison, the experimental results indicate that the PFLS outperforms the CM in all metrics. As shown in Supplementary Table S2, Supplementary Table S3, in the test sets, PFLS exhibits higher NRI and IDI indices with significant p-values (p < 0.05) compared to CM across cohorts, indicating its superiority.

Comparison of the PFLS with MDCM

As shown in Table 3, the AUC of the MDCM on the test sets from the five cohorts are 0.752, 0.842, 0.825, 0.746, and 0.820. The experimental results indicate that the MDCM exhibits relatively poor consistency in performance across different centers. This performance variability may be attributed to the inherent heterogeneity in the data from various centers. When training on merged data, such heterogeneity can lead to suboptimal model generalization, as the model may be biased toward certain centerspecific characteristics rather than capturing universally representative features. Consequently, the model's ability to perform consistently well across diverse and unseen datasets is compromised.

Comparison of the PFLS with personalized federated model

As shown in Table 4, the three personalized federated models, FedProx, FedBN, and Moon, achieved AUC values ranging from 0.746 to 0.834 and accuracies between 0.619 and 0.749 across the five cohorts. However, all three methods exhibited unstable performance when confronted with heterogeneous multi-center data. In contrast, our proposed PFLS consistently outperformed them, with AUC improvements of approximately 4%-15% and accuracy gains of 6%-22% across different cohorts. The performance of PFLS and other personalized federated learning algorithms on the training and test sets across the five cohorts is shown in Supplementary Figure S2, Supplementary Figure S3. Additionally, the Decision Curve Analysis are shown in Supplementary Figure S4.

Ablation experiments

As shown in Table 5, we present the AUC values on the training and test sets for different ablation variants. When the Reptile step was

TABLE 2 Performance of logistic regression, SVM, random forest and ELM.

Calaaut 1		Trainin	ıg set 1		Testing set 1							
Cohort 1	SVM	RF	LR	ELM	SVM	RF	LR	ELM				
AUC	0.923	0.952	0.822	0.895	0.812	0.852	0.777	0.846				
(95%CI)	(0.867-0.980)	(0.920-0.984)	(0.741-0.903)	(0.833-0.956)	(0.700-0.923)	(0.756-0.949)	(0.662-0.893)	(0.748-0.944)				
Sensitivity	0.819	0.882	0.724	0.906	0.694	0.800	0.682	0.826				
	(104/127)	(112/127)	(92/127)	(115/127)	(59/85)	(68/85)	(58/85)	(71/85)				
Specificity	0.941	0.941	0.765	0.765	0.792	0.708	0.750	0.739				
	(32/34)	(32/34)	(26/34)	(26/34)	(19/24)	(17/24)	(18/24)	(17/24)				
Accuracy	0.845	0.894	0.733	0.876	0.716	0.780	0.697	0.807				
	(136/161)	(144/161)	(118/161)	(141/161)	(78/109)	(85/109)	(76/109)	(88/109)				
PPV	0.981	0.982	0.920	0.935	0.922	0.907	0.906	0.922				
	(104/106)	(112/114)	(92/100)	(115/123)	(59/64)	(68/75)	(58/64)	(71/77)				
NPV	0.582	0.681	0.426	0.684	0.422	0.500	0.400	0.531				
	(32/55)	(32/47)	(26/61)	(26/38)	(19/45)	(17/34)	(18/45)	(17/32)				
Cobout 2		Trainin	g set 2		Testing set 2							
Cohort 2	SVM	RF	LR	ELM	SVM	RF	LR	ELM				
AUC	0.842	0.892	0.838	0.925	0.746	0.900	0.831	0.889				
(95%CI)	(0.702-0.981)	(0.803-0.981)	(0.729-0.947)	(0.853-0.997)	(0.570-0.922)	(0.796-1.000)	(0.689-0.973)	(0.771-1.000)				
Sensitivity	0.757	0.757	0.676	0.892	0.808	0.923	0.808	0.923				
	(28/37)	(28/37)	(25/37)	(33/37)	(21/26)	(24/26)	(21/26)	(24/26)				
Specificity	0.929	1.000	0.929	0.786	0.400	0.600	0.600	0.700				
	(13/14)	(14/14)	(13/14)	(11/14)	(4/10)	(6/10)	(6/10)	(7/10)				
Accuracy	0.804	0.824	0.745	0.863	0.694	0.833	0.750	0.861				
	(41/51)	(42/51)	(38/51)	(44/51)	(25/36)	(30/36)	(27/36)	(31/36)				
PPV	0.966	1.000	0.962	0.917	0.778	0.857	0.840	0.889				
	(28/29)	(28/28)	(25/26)	(33/36)	(21/27)	(24/28)	(21/25)	(24/27)				
NPV	0.591	0.609	0.520	0.733	0.444	0.750	0.545	0.778				
	(13/22)	(14/23)	(13/25)	(11/15)	(4/9)	(6/8)	(6/11)	(7/9)				
Cohort 3		Trainin	g set 3		Testing set 3							
Conort 3	SVM	RF	LR	ELM	SVM	RF	LR	ELM				
AUC	0.964	0.984	0.832	0.939	0.888	0.953	0.790	0.922				
(95%CI)	(0.924-1.000)	(0.948-1.000)	(0.707-0.957)	(0.884-0.994)	(0.784-0.992)	(0.897-1.000)	(0.620-0.960)	(0.845-0.999)				
Sensitivity	0.939	0.959	0.857	0.918	0.794	0.971	0.824	0.882				
	(46/49)	(47/49)	(42/49)	(45/49)	(27/34)	(33/34)	(28/34)	(30/34)				
Specificity	0.952	1.000	0.762	0.857	0.800	0.733	0.800	0.867				
	(20/21)	(21/21)	(16/21)	(18/21)	(12/15)	(11/15)	(12/15)	(13/15)				
Accuracy	0.943	0.971	0.829	0.900	0.796	0.898	0.816	0.878				
	(66/70)	(68/70)	(58/70)	(63/70)	(39/49)	(44/49)	(40/49)	(43/49)				
PPV	0.979	1.000	0.894	0.938	0.900	0.892	0.903	0.938				
	(46/47)	(47/47)	(42/47)	(45/48)	(27/30)	(33/37)	(28/31)	(30/32)				
NPV	0.870	0.913	0.696	0.818	0.632	0.917	0.667	0.765				
	(20/23)	(21/23)	(16/23)	(18/22)	(12/19)	(11/12)	(12/18)	(13/17)				
Cohort 4		Trainin	g set 4		Testing set 4							
Conort 4	SVM RF LR ELM			ELM	SVM	RF	LR	ELM				
AUC	0.931	0.972	0.837	0.886	0.869	0.875	0.834	0.876				
(95%CI)	(0.903-0.960)	(0.956-0.987)	(0.786-0.888)	(0.845-0.926)	(0.814-0.923)	(0.819-0.931)	(0.769-0.9000)	(0.825-0.927)				

TABLE 2 Continued

		Traini	ng set 4			Testin	g set 4	
Cohort 4	SVM	RF	LR	ELM	SVM	RF	LR	ELM
Sensitivity	0.812	0.930	0.726	0.807	0.736	0.824	0.736	0.768
	(151/186)	(173/186)	(135/186)	(150/186)	(92/125)	(103/125)	(92/125)	(96/125)
Specificity	0.906	0.896	0.833	0.833	0.828	0.766	0.797	0.828
	(87/96)	(86/96)	(80/96)	(80/96)	(53/64)	(49/64)	(51/64)	(53/64)
Accuracy	0.844	0.918	0.762	0.816	0.767	0.804	0.757	0.788
	(238/282)	(259/282)	(215/282)	(230/282)	(145/189)	(152/189)	(143/189)	(149/189)
PPV	0.944	0.945	0.894	0.904	0.893	0.873	0.876	0.897
	(151/160)	(173/183)	(135/151)	(150/166)	(92/103)	(103/118)	(92/105)	(96/107)
NPV	0.713	0.869	0.611	0.690	0.616	0.690	0.607	0.646
	(87/122)	(86/99)	(80/131)	(80/116)	(53/86)	(49/71)	(51/84)	(53/82)
Cabaut		Traini	ng set 5			Testin	g set 5	
Cohort 5	SVM	RF	LR	ELM	SVM	RF	LR	ELM
AUC	0.886	0.949	0.837	0.924	0.893	0.805	0.831	0.893
(95%CI)	(0.793-0.980)	(0.894-1.000)	(0.724-0.951)	(0.850-0.998)	(0.780-1.000)	(0.656-0.955)	(0.679-0.983)	(0.770-1.000)
Sensitivity	0.833	0.875	0.708	0.958	0.882	0.941	0.882	0.824
	(20/24)	(21/24)	(17/24)	(23/24)	(15/17)	(16/17)	(15/17)	(14/17)
Specificity	0.864	0.909	0.864	0.727	0.750	0.500	0.750	0.813
	(19/22)	(20/22)	(19/22)	(16/22)	(12/16)	(8/16)	(12/16)	(13/16)
Accuracy	0.848	0.891	0.783	0.848	0.818	0.727	0.818	0.818
	(39/46)	(41/46)	(36/46)	(39/46)	(27/33)	(24/33)	(27/33)	(27/33)
PPV	0.870	0.913	0.850	0.793	0.789	0.667	0.789	0.824
	(20/23)	(21/23)	(17/20)	(23/29)	(15/19)	(16/24)	(15/19)	(14/17)
NPV	0.826 (19/23)	0.870 (20/23)	0.731 (19/26)	0.941 (16/17)	0.857 (12/14)	0.889 (8/9)	0.857 (12/14)	0.813 (13/16)

removed, the performance across centers showed little variation but remained at a relatively low level, indicating that the absence of meta-update limited the global model's ability to provide a good initialization. When the personalization step was removed, some centers achieved relatively good results, while others showed significantly lower AUC values, leading to large performance discrepancies across sites. This suggests that without personalization, relying solely on the global model cannot effectively adapt to heterogeneous data distributions. In contrast, the complete PFLS achieved both the best overall performance and balanced results across centers, further demonstrating the complementary roles of the Reptile step and the personalization step in enhancing model generalization and adaptability.

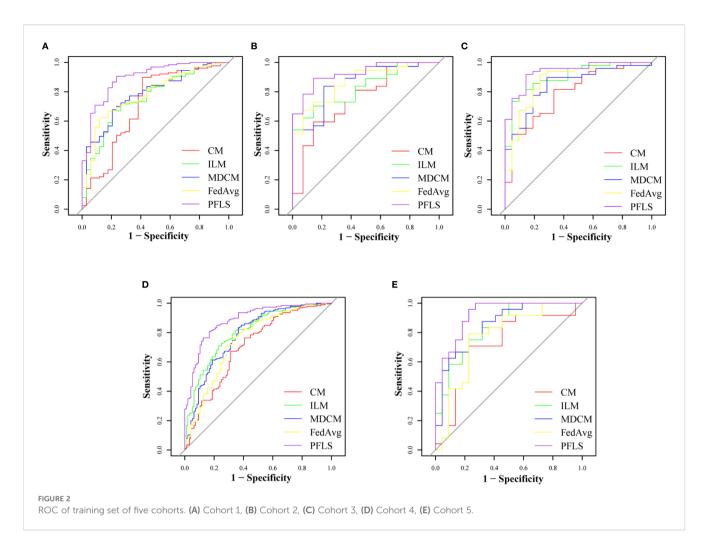
Discussion

The analysis of lung CT images to accurately differentiate between patients with TBG and LAC in a non-invasive manner holds considerable clinical value. In this study, a PFLS is used to collaboratively use CT image data of TBG and LAC from five cohorts while maintaining patient privacy, enabling the training of robust models for each medical institution. The proposed method shows superior predictive performance than the compared methods

on the data from each medical institution for distinguishing between TBG and LAC.

Research shows that gender, age, morphology, and spiculation are significantly different between patients with TBG and LAC (27). As malignant tumors predominantly grow in lung parenchyma, nodules of LAC patients are more likely to exhibit irregular margins, lobulation, and spiculation (28). Although lobulation is a distinct feature of malignant lung nodules, several studies have shown that 25% of benign nodules also exhibit lobulation. Spiculation presents a more significant correlation with LAC (29). Pathologically, spiculation is attributed to fibrous tissue proliferation induced by interstitial thickening and peripheral vascular occlusion. However, identification based on morphologic features of nodules is highly subjective, and the criteria for determination vary among radiologists, limiting the utility of shape features in distinguishing between benign and malignant nodules (30). Therefore, the performance of the CM we constructed based on the above features was not satisfactory.

The convolutional neural networks (CNNs) are capable of automatically extracting features from images and generating features at various levels of abstraction. Among CNNs, Low-level layers produce details like edges, textures, and corners, while high-level layers produce globally abstracted features. Deep CNNs are extensively applied in medical imaging and achieving commendable

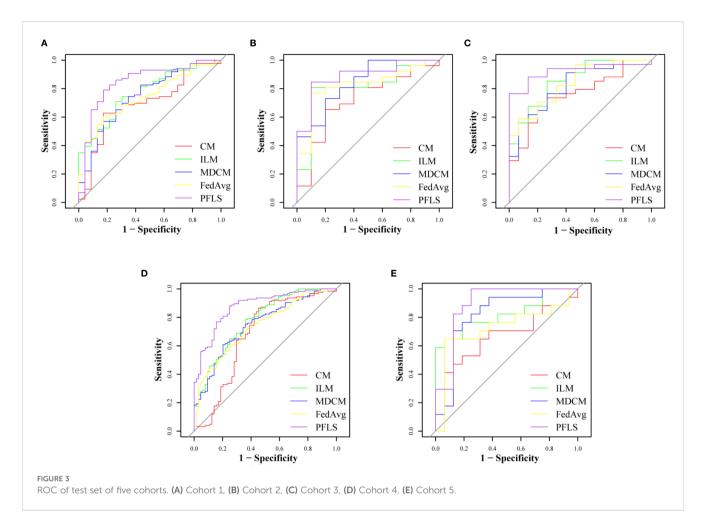


results (31–33). Nevertheless, deep CNNs are susceptible to overfitting, particularly when training with a limited number of samples. Therefore, when training the ILM, we used the transfer learning strategy to pre-train the model with Imagnet data, and then fine-tune the model with the CT data of this study to alleviate the problem of overfitting. However, perhaps due to the smaller number of training samples, the performance of ILM on the test sets of various medical institutions is generally worse than that of PFLS.

Federated learning is a distributed machine learning approach that enables collaborative training of machine learning models using data from multiple hospitals, while eliminating data leakage. It efficiently addresses the data island issue and mitigates model overfitting due to insufficient training samples from a single medical institution. However, due to differences of CT equipment between hospitals, imaging results from different hospitals are heterogeneous, meaning that using the same machine learning model to infer CT images from different hospitals does not ensure satisfactory performance for all hospitals. Except for the slightly better performance of fedag on the test set of institution 2 compared to ILM, the performance of other institutions was worse than that of ILM. This is primarily due to the apparent heterogeneity of the data across hospitals, resulting in apparent

discrepancies between the global model and the actual optimal model of a local hospital, and then leading to poorer predictive performance of the global model in certain hospitals. In addition, we also compared our framework with a centralized model (MDCM) trained by directly merging data from all centers. Although this setting represents an ideal scenario without privacy constraints, the MDCM exhibited unstable performance and large variations across sites. This inconsistency can be attributed to the substantial heterogeneity in imaging protocols, scanner vendors, and patient populations among different hospitals. When data are simply pooled, the model may become biased toward dominant site-specific features, thereby limiting its generalizability. In contrast, PFLS achieved more stable and consistent performance across centers, further underscoring the necessity of federated and personalized federated strategies for real-world multicenter applications.

The proposed PFLS employs the Reptile algorithm by treating the local models of federated learning as the different task of meta learning to collaboratively train the global model, and then fine-tunes global models for personalizing local model of each medical institution. The meta-learning-based federated learning frameworks of the proposed method effectively alleviates overfitting. And the robust personalized local models, which are fine-tuned global models by local hospitals,



enhance the generalization ability of the local models while effectively mitigating the performance decrease caused by data heterogeneity among different hospitals. In terms of classification performance, the proposed method shows remarkable advantages on all datasets, except

for slightly lower AUC values compare to ILM, FedAvg, and CM. NRI and IDI show that, except for medical institutions 2 and 5 with relatively less test data, the PFLS performance of other institutions has significantly improved compared with ILM, FedAvg, and CM.

TABLE 3 Performance of the of models on the each cohort.

Cabaut 1		1	Training set	1		Testing set 1							
Cohort 1	СМ	ILM	MDCM	Fedavg	PFLS	СМ	ILM	MDCM	Fedavg	PFLS			
AUC (95%CI)	0.712 (0.601- 0.824)	0.759 (0.669- 0.849)	0.784 (0.702- 0.866)	0.782 (0.698- 0.866)	0.895 (0.833- 0.956)	0.692 (0.569- 0.815)	0.773 (0.672- 0.873)	0.752 (0.630- 0.857)	0.740 (0.637- 0.843)	0.846 (0.748- 0.944)			
Sensitivity	0.898	0.669	0.677	0.622	0.906	0.744	0.709	0.588	0.477	0.826			
	(114/127)	(85/127)	(86/127)	(79/127)	(115/127)	(64/86)	(61/86)	(50/86)	(41/86)	(71/86)			
Specificity	0.588	0.794	0.794	0.882	0.765	0.391	0.739	0.750	0.870	0.739			
	(20/34)	(27/34)	(27/34)	(30/34)	(26/34)	(9/23)	(17/23)	(18/23)	(20/23)	(17/23)			
Accuracy	0.832	0.696	0.702	0.677	0.876	0.670	0.716	0.624	0.560	0.807			
	(134/161)	(112/161)	(113/161)	(109/161)	(141/161)	(73/109)	(78/109)	(68/109)	(61/109)	(88/109)			
PPV	0.891	0.924	0.925	0.952	0.935	0.821	0.910	0.893	0.932	0.922			
	(114/128)	(85/92)	(86/93)	(79/83)	(115/123)	(64/78)	(61/67)	(50/56)	(41/44)	(71/77)			
NPV	0.606	0.391	0.397	0.385	0.684	0.290	0.405	0.340	0.308	0.531			
	(20/33)	(27/69)	(27/68)	(30/78)	(26/38)	(9/31)	(17/42)	(18/53)	(20/65)	(17/32)			

TABLE 3 Continued

		Т	raining set	2		Testing set 2							
Cohort 2	СМ	ILM	MDCM	FedAvg	PFLS	СМ	ILM	MDCM	FedAvg	PFLS			
AUC (95%CI)	0.749 (0.594- 0.905)	0.828 (0.715- 0.941)	0.855 (0.742- 0.968)	0.869 (0.767- 0.971)	0.925 (0.853- 0.997)	0.719 (0.523- 0.915)	0.808 (0.640- 0.975)	0.842 (0.698- 0.987)	0.823 (0.672- 0.974)	0.889 (0.771- 1.000)			
Sensitivity	0.595	0.703	0.838	0.676	0.892	0.654	0.808	0.885	0.615	0.923			
	(22/37)	(26/37)	(31/37)	(25/37)	(33/37)	(17/26)	(21/26)	(23/26)	(16/26)	(24/26)			
Specificity	0.857	0.857	0.786	0.857	0.786	0.700	0.900	0.600	0.900	0.700			
	(12/14)	(12/14)	(11/14)	(12/14)	(11/14)	(7/10)	(9/10)	(6/10)	(9/10)	(7/10)			
Accuracy	0.667	0.745	0.824	0.726	0.863	0.667	0.833	0.806	0.694	0.861			
	(34/51)	(38/51)	(42/51)	(37/51)	(44/51)	(24/36)	(30/36)	(29/36)	(25/36)	(31/36)			
PPV	0.917	0.929	0.912	0.926	0.917	0.850	0.955	0.852	0.941	0.889			
	(22/24)	(26/28)	(31/34)	(25/27)	(33/36)	(17/20)	(21/22)	(23/27)	(16/17)	(24/27)			
NPV	0.444	0.522	0.647	0.500	0.733	0.438	0.643	0.667	0.474	0.778			
	(12/27)	(12/23)	(11/17)	(12/24)	(11/15)	(7/16)	(9/14)	(6/9)	(9/19)	(7/9)			
Cohort 3		Т	raining set	3			T	esting set	3				
Conort 3	СМ	ILM	MDCM	FedAvg	PFLS	СМ	ILM	MDCM	FedAvg	PFLS			
AUC (95%CI)	0.794 (0.680- 0.908)	0.903 (0.829- 0.977)	0.847 (0.752- 0.943)	0.880 (0.786- 0.973)	0.939 (0.884- 0.994)	0.759 (0.618- 0.900)	0.857 (0.745- 0.969)	0.825 (0.701- 0.950)	0.839 (0.724- 0.955)	0.922 (0.845- 0.999)			
Sensitivity	0.816	0.735	0.898	0.918	0.918	0.882	0.735	0.824	0.971	0.882			
	(40/49)	(36/49)	(44/49)	(45/49)	(45/49)	(30/34)	(25/34)	(28/34)	(33/34)	(30/34)			
Specificity	0.667	0.905	0.714	0.714	0.857	0.333	0.733	0.600	0.533	0.867			
	(14/21)	(19/21)	(15/21)	(15/21)	(18/21)	(5/15)	(11/15)	(9/15)	(8/15)	(13/15)			
Accuracy	0.771	0.786	0.843	0.857	0.900	0.714	0.735	0.755	0.837	0.878			
	(54/70)	(55/70)	(59/70)	(60/70)	(63/70)	(35/49)	(36/49)	(37/49)	(41/49)	(43/49)			
PPV	0.851	0.947	0.880	0.882	0.938	0.750	0.862	0.824	0.825	0.938			
	(40/47)	(36/38)	(44/50)	(45/51)	(45/48)	(30/40)	(25/29)	(28/34)	(33/40)	(30/32)			
NPV	0.609	0.594	0.750	0.790	0.818	0.556	0.550	0.600	0.889	0.765			
	(14/23)	(19/32)	(15/20)	(15/19)	(18/22)	(5/9)	(11/20)	(9/15)	(8/9)	(13/17)			
Cohort 4		Т	raining set	4			Т	esting set	4				
COHOIC 4	СМ	ILM	MDCM	FedAvg	PFLS	СМ	ILM	MDCM	FedAvg	PFLS			
AUC (95%CI)	0.706 (0.639- 0.772)	0.805 (0.751- 0.858)	0.785 (0.727- 0.844)	0.748 (0.685- 0.812)	0.886 (0.845- 0.926)	0.679 (0.587- 0.770)	0.778 (0.707- 0.848)	0.746 (0.673- 0.818)	0.737 (0.664- 0.810)	0.876 (0.825- 0.927)			
Sensitivity	0.672	0.710	0.833	0.801	0.807	0.664	0.720	0.768	0.752	0.768			
	(125/186)	(132/186)	(155/186)	(149/186)	(150/186)	(83/125)	(90/125)	(96/125)	(94/125)	(96/125)			
Specificity	0.688	0.771	0.635	0.615	0.833	0.641	0.656	0.594	0.594	0.828			
	(66/96)	(74/96)	(61/96)	(59/96)	(80/96)	(41/64)	(42/64)	(38/64)	(38/64)	(53/64)			
Accuracy	0.677	0.731	0.766	0.738	0.816	0.656	0.698	0.709	0.698	0.788			
	(191/282)	(206/282)	(216/282)	(208/282)	(230/282)	(124/189)	(132/189)	(134/189)	(132/189)	(149/189)			
PPV	0.806	0.857	0.816	0.801	0.904	0.783	0.804	0.787	0.783	0.897			
	(125/155)	(132/154)	(155/190)	(149/186)	(150/166)	(83/106)	(90/112)	(96/122)	(94/120)	(96/107)			
NPV	0.520	0.578	0.663	0.615	0.690	0.494	0.546	0.567	0.551	0.646			
	(66/127)	(74/128)	(61/92)	(59/96)	(80/116)	(41/83)	(42/77)	(38/67)	(38/69)	(53/82)			

TABLE 3 Continued

			Training set	: 5		Testing set 5							
Cohort 5	СМ	ILM	MDCM	FedAvg	PFLS	СМ	ILM	MDCM	FedAvg	PFLS			
AUC (95%CI)	0.714 (0.555- 0.873)	0.835 (0.720- 0.951)	0.854 (0.744- 0.964)	0.769 (0.623- 0.915)	0.924 (0.850- 0.998)	0.665 (0.471- 0.860)	0.824 (0.676- 0.972)	0.820 (0.661- 0.979)	0.728 (0.537- 0.919)	0.893 (0.770- 1.000)			
Sensitivity	0.708	0.750	0.792	0.792	0.958	0.647	0.882	0.412	0.647	0.824			
	(17/24)	(18/24)	(19/24)	(19/24)	(23/24)	(11/17)	(15/17)	(7/17)	(11/17)	(14/17)			
Specificity	0.773	0.727	0.773	0.773	0.727	0.625	0.375	0.875	0.688	0.813			
	(17/22)	(16/22)	(17/22)	(17/22)	(16/22)	(10/16)	(6/16)	(14/16)	(11/16)	(13/16)			
Accuracy	0.739	0.739	0.783	0.783	0.848	0.636	0.636	0.636	0.667	0.818			
	(34/46)	(34/46)	(36/46)	(36/46)	(39/46)	(21/33)	(21/33)	(21/33)	(22/33)	(27/33)			
PPV	0.773	0.750	0.792	0.792	0.793	0.647	0.600	0.778	0.688	0.824			
	(17/22)	(18/24)	(19/24)	(19/24)	(23/29)	(11/17)	(15/25)	(7/9)	(11/16)	(14/17)			
NPV	0.708	0.727	0.773	0.773	0.941	0.625	0.750	0.583	0.647	0.813			
	(17/24)	(16/22)	(17/22)	(17/22)	(16/17)	(10/16)	(6/8)	(14/24)	(11/17)	(13/16)			

Compared with traditional machine learning models, Bayesian ELM has significant advantages in alleviating overfitting and balancing model complexity with generalization ability, thereby providing more robust and clinically valuable predictive performance. From a theoretical perspective, extreme learning machine (ELM) randomly generates hidden layer weights and analytically solves the output weights, avoiding the complex

gradient-based iterative process in traditional neural networks, which grants it faster training speed and stronger representation capacity. Building on this, the introduction of the Bayesian framework not only enables probabilistic modeling of model parameters and provides uncertainty estimation but also effectively suppresses overfitting through prior and posterior constraints. This combination allows Bayesian ELM to maintain

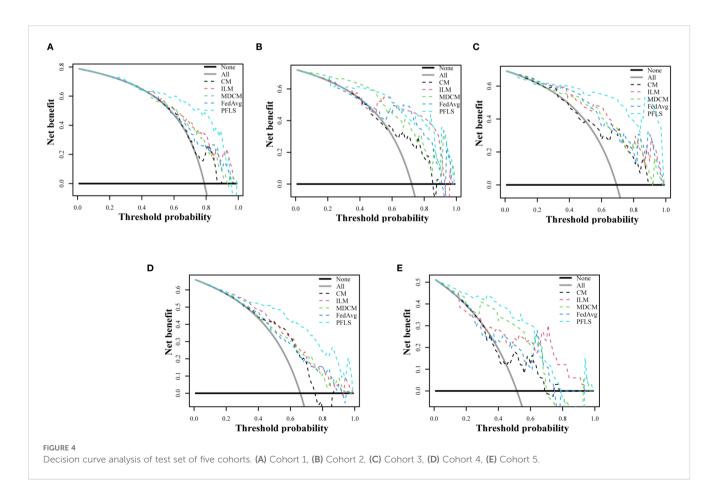


TABLE 4 Performance comparison of models trained under different FL algorithms.

0.14		Trainir	ng set 1		Testing set 1							
Cohort 1	Fedprox	FedBN	Moon	PFLS	Fedprox	FedBN	Moon	PFLS				
AUC	0.790	0.770	0.761	0.895	0.764	0.741	0.752	0.846				
(95%CI)	(0.716-0.863)	(0.677-0.863)	(0.671-0.851)	(0.833-0.956)	(0.646-0.883)	(0.628-0.853)	(0.639-0.865)	(0.748-0.944)				
Sensitivity	0.567	0.835	0.512	0.906	0.512	0.651	0.523	0.826				
	(72/127)	(106/127)	(65/127)	(115/127)	(44/86)	(56/86)	(45/86)	(71/86)				
Specificity	0.912	0.618	0.912	0.765	0.870	0.739	0.783	0.739				
	(31/34)	(21/34)	(31/34)	(26/34)	(20/23)	(17/23)	(18/23)	(17/23)				
Accuracy	0.640	0.789	0.596	0.876	0.587	0.670	0.578	0.807				
	(103/161)	(127/161)	(96/161)	(141/161)	(64/109)	(73/109)	(63/109)	(88/109)				
PPV	0.960	0.891	0.956	0.935	0.936	0.903	0.900	0.922				
	(72/75)	(106/119)	(65/68)	(115/123)	(44/47)	(56/62)	(45/50)	(71/77)				
NPV	0.360	0.500	0.333	0.684	0.323	0.362	0.305	0.531				
	(31/86)	(21/42)	(31/93)	(26/38)	(20/62)	(17/47)	(18/59)	(17/32)				
Cohort 2		Trainir	ig set 2		Testing set 2							
Conort 2	FedProx	FedBN	Moon	PFLS	FedProx	FedBN	Moon	PFLS				
AUC	0.844	0.849	0.842	0.925	0.800	0.808	0.813	0.889				
(95%CI)	(0.729-0.959)	(0.735-0.965)	(0.717-0.970)	(0.853-0.997)	(0.610-0.990)	(0.620-0.991)	(0.648-0.967)	(0.771-1.000)				
Sensitivity	0.784	0.757	0.946	0.892	0.654	0.962	0.885	0.923				
	(29/37)	(28/37)	(35/37)	(33/37)	(17/26)	(25/26)	(23/26)	(24/26)				
Specificity	0.857	0.857	0.643	0.786	0.800	0.500	0.600	0.700				
	(12/14)	(12/14)	(9/14)	(11/14)	(8/10)	(5/10)	(6/10)	(7/10)				
Accuracy	0.804	0.784	0.863	0.863	0.694	0.833	0.806	0.861				
	(41/51)	(40/51)	(44/51)	(44/51)	(25/36)	(30/36)	(29/36)	(31/36)				
PPV	0.935	0.933	0.875	0.917	0.895	0.833	0.852	0.889				
	(29/31)	(28/30)	(35/40)	(33/36)	(17/19)	(25/30)	(23/27)	(24/27)				
NPV	0.600	0.571	0.818	0.733	0.471	0.833	0.667	0.778				
	(12/20)	(12/21)	(9/11)	(11/15)	(8/17)	(5/6)	(6/9)	(7/9)				
Cohort 3		Trainir	ıg set 3		Testing set 3							
Conort 3	FedProx	FedBN	Moon	PFLS	FedProx	FedBN	Moon	PFLS				
AUC	0.875	0.847	0.869	0.939	0.831	0.812	0.829	0.922				
(95%CI)	(0.772-0.977)	(0.753-0.945)	(0.758-0.980)	(0.884-0.994)	(0.701-0.961)	(0.664-0.960)	(0.665-0.994)	(0.845-0.999)				
Sensitivity	0.776	0.796	0.776	0.918	0.824	0.853	0.794	0.882				
	(38/49)	(39/49)	(38/49)	(45/49)	(28/34)	(29/34)	(27/34)	(30/34)				
Specificity	0.905	0.857	0.905	0.857	0.733	0.667	0.800	0.867				
	(19/21)	(18/21)	(19/21)	(18/21)	(11/15)	(10/15)	(12/15)	(13/15)				
Accuracy	0.814	0.814	0.814	0.900	0.796	0.796	0.796	0.878				
	(57/70)	(57/70)	(57/70)	(63/70)	(39/49)	(39/49)	(39/49)	(43/49)				
PPV	0.950	0.929	0.950	0.938	0.875	0.853	0.900	0.938				
	(38/40)	(39/42)	(38/40)	(45/48)	(28/32)	(29/34)	(27/30)	(30/32)				
NPV	0.633	0.643	0.633	0.818	0.647	0.667	0.632	0.765				
	(19/30)	(18/28)	(19/30)	(18/22)	(11/17)	(10/15)	(12/19)	(13/17)				
Cobort 1		Trainin	g set 4		Testing set 4							
Cohort 4	FedProx	FedBN	Moon	PFLS	FedProx	FedBN	Moon	PFLS				
AUC	0.779	0.754	0.802	0.886	0.748	0.730	0.755	0.876				
(95%CI)	(0.723-0.834)	(0.696-0.813)	(0.746-0.857)	(0.845-0.926)	(0.674-0.822)	(0.651-0.809)	(0.683-0.827)	(0.825-0.927)				

10.3389/fonc.2025.1666937 Chen et al.

TABLE 4 Continued

		Trainiı	ng set 4		Testing set 4						
Cohort 4	FedProx	FedBN	Moon	PFLS	FedProx	FedBN	Moon	PFLS			
Sensitivity	0.672	0.672	0.688	0.807	0.760	0.680	0.648	0.768			
	(125/186)	(125/186)	(128/186)	(150/186)	(95/125)	(85/125)	(81/125)	(96/125)			
Specificity	0.792	0.729	0.813	0.833	0.641	0.672	0.734	0.828			
	(76/96)	(70/96)	(78/96)	(80/96)	(41/64)	(43/64)	(47/64)	(53/64)			
Accuracy	0.713	0.691	0.730	0.816	0.720	0.677	0.677	0.788			
	(201/282)	(195/282)	(206/282)	(230/282)	(136/189)	(128/189)	(128/189)	(149/189)			
PPV	0.862	0.828	0.877	0.904	0.805	0.802	0.827	0.897			
	(125/145)	(125/151)	(128/146)	(150/166)	(95/118)	(85/106)	(81/98)	(96/107)			
NPV	0.555	0.534	0.574	0.690	0.577	0.518	0.516	0.646			
	(76/137)	(70/131)	(78/136)	(80/116)	(41/71)	(43/83)	(47/91)	(53/82)			
Calaant F		Traini	ng set 5		Testing set 5						
Cohort 5	FedProx	FedBN	Moon	PFLS	FedProx	FedBN	Moon	PFLS			
AUC	0.845	0.758	0.820	0.924	0.798	0.746	0.809	0.893			
(95%CI)	(0.731-0.958)	(0.614-0.901)	(0.698-0.941)	(0.850-0.998)	(0.644-0.951)	(0.571-0.921)	(0.652-0.962)	(0.770-1.000)			
Sensitivity	0.667	0.667	0.583	0.958	0.471	0.588	0.412	0.824			
	(16/24)	(16/24)	(14/24)	(23/24)	(8/17)	(10/17)	(7/17)	(14/17)			
Specificity	0.909	0.818	0.955	0.727	0.875	0.688	0.875	0.813			
	(20/22)	(18/22)	(21/22)	(16/22)	(14/16)	(11/16)	(14/16)	(13/16)			
Accuracy	0.783	0.739	0.761	0.848	0.667	0.636	0.636	0.818			
	(36/46)	(34/46)	(35/46)	(39/46)	(22/33)	(21/33)	(21/33)	(27/33)			
PPV	0.889	0.800	0.933	0.793	0.800	0.667	0.778	0.824			
	(16/18)	(16/20)	(14/15)	(23/29)	(8/10)	(10/15)	(7/9)	(14/17)			
NPV	0.714	0.692	0.677	0.941	0.609	0.611	0.583	0.813			
	(20/28)	(18/26)	(21/31)	(16/17)	(14/23)	(11/18)	(14/24)	(13/16)			

efficient training while better balancing model complexity and generalization, thereby demonstrating stronger robustness and stability when applied to heterogeneous multi-center medical data.

Despite the promising results, our study has some limitations. First, the implementation of strict inclusion and exclusion criteria for samples could introduce bias in sample selection, potentially affecting model training. Second, the study includes only TBG, a

specific type of benign nodule, and misses other benign nodules such as inflammatory pseudotumors, hamartomas, and fibromas. Third, the aggregation of local models in the Reptile stages dose not account for differences in data distribution across hospitals, potentially leading to suboptimal global models in the Reptile stage and affecting the performance of the robust local models trained in the next stage.

TABLE 5 AUC results of ablation experiments.

		Personalization stage	AUC (95%CI)											
Fedavg	Fedavg Reptile stage		Cohort 1		Cohort 2		Cohort 3		Cohort 4		Cohort 5			
stage	stage		Train set	Test set										
✓		1	0.781 (0.698- 0.865)	0.753 (0.650- 0.857)	0.851 (0.743- 0.960)	0.812 (0.671- 0.953)	0.828 (0.718- 0.938)	0.812 (0.691- 0.933)	0.811 (0.756- 0.867)	0.768 (0.696- 0.842)	0.845 (0.735- 0.955)	0.816 (0.653- 0.980)		
1	/		0.770 (0.685- 0.856)	0.740 (0.641- 0.839)	0.826 (0.693- 0.960)	0.804 (0.620- 0.988)	0.832 (0.717- 0.947)	0.808 (0.657- 0.959)	0.738 (0.677- 0.800)	0.698 (0.622- 0.775)	0.782 (0.647- 0.918)	0.735 (0.562- 0.909)		
1	1	1	0.895 (0.833- 0.956)	0.846 (0.748- 0.944)	0.925 (0.853- 0.997)	0.889 (0.771- 1.000)	0.939 (0.884- 0.994)	0.922 (0.845- 0.999)	0.886 (0.845- 0.926)	0.876 (0.825- 0.927)	0.924 (0.850- 0.998)	0.893 (0.770- 1.000)		

Conclusion

The PFLS proposed in this study facilitates collaborative training across multiple hospitals while maintaining the data privacy of each hospital. It effectively mitigates the model overfitting caused by insufficient samples from a single hospital. Moreover, the personalizing process of local model address the heterogeneity of data across hospitals, which cannot be adequately performed by a single global model. The resulting robust local models show excellent discrimination between LAC and TBG, providing invaluable assistance to clinicians in improving diagnostic accuracy.

Data availability statement

Due to institutional policies and patient privacy regulations, the data are not publicly available. Access to the dataset is restricted and can only be granted upon reasonable request and with approval from the corresponding ethics committees of the participating institutions. Requests to access the datasets should be directed to XMC: 3897001254@qq.com.

Ethics statement

The study was approved by the Ethics Committee of Jiangmen Central Hospital, and ethical approvals were obtained from all participating institutions. The studies were conducted in accordance with the local legislation and institutional requirements. Written informed consent for participation was not required from the participants or the participants' legal guardians/next of kin in accordance with the national legislation and institutional requirements. Written informed consent was obtained from the individual(s) for the publication of any potentially identifiable images or data included in this article.

Author contributions

CY: Supervision, Writing – original draft, Data curation, Project administration, Resources, Writing – review & editing. LL: Formal analysis, Writing – review & editing, Conceptualization, Investigation, Writing – original draft. BF: Methodology, Writing – review & editing, Funding acquisition, Supervision. YC: Project administration, Writing – review & editing, Data curation, Validation. JX: Visualization, Formal analysis, Methodology, Writing – review & editing, HL: Writing – review & editing, Resources, Data curation. KL: Writing – review & editing, Supervision, Resources, Data curation. XDC: Writing – review & editing, Data curation, Project administration. YK: Resources, Writing – review & editing, Supervision, Data curation. HZ:

Writing – review & editing, Supervision, Project administration, Methodology. QH: Writing – review & editing, Data curation, Resources. QJ: Conceptualization, Visualization, Writing – review & editing. WL: Formal analysis, Visualization, Supervision, Writing – review & editing, Methodology. QL: Writing – review & editing, Project administration, Data curation, Resources. XMC: Writing – original draft, Resources, Investigation, Funding acquisition, Data curation.

Funding

The author(s) declare financial support was received for the research and/or publication of this article. This work was supported by the National Natural Science Foundation of China (82460361, 62176104, 12261027), GUAT Special Research Project on the Strategic Development of Distinctive Interdisciplinary Fields (TS2024231) and The Bagui Youth Top Talent Training Program.

Conflict of interest

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Generative Al statement

The author(s) declare that no Generative AI was used in the creation of this manuscript.

Any alternative text (alt text) provided alongside figures in this article has been generated by Frontiers with the support of artificial intelligence and reasonable efforts have been made to ensure accuracy, including review by the authors wherever possible. If you identify any issues, please contact us.

Publisher's note

All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

Supplementary material

The Supplementary Material for this article can be found online at: https://www.frontiersin.org/articles/10.3389/fonc.2025.1666937/full#supplementary-material

References

- 1. Henschke CI, Salvatore M, Cham M, Powell CA, DiFabrizio L, Flores R, et al. Baseline and annual repeat rounds of screening: implications for optimal regimens of screening. *Eur Radiol.* (2017) 28:1085–94. doi: 10.1007/s00330-017-5029-z
- Vishal K, Sagar K, David P, William D, Jeremy A, Richard R, et al. A practical algorithmic approach to the diagnosis and management of solitary pulmonary nodules: part 1: radiologic characteristics and imaging modalities. *Chest.* (2013) 143:825–39. doi: 10.1378/chest.12-0960
- 3. Vishal K, Sagar K, David P, William D, Jeremy A, Richard R, et al. A practical algorithmic approach to the diagnosis and management of solitary pulmonary nodules: part 2: pretest probability and algorithm. *Chest.* (2013) 143:840–6. doi: 10.1378/chest.12-1487
- 4. Henschke CI, Salvatore M, Cham M, Powell CA, DiFabrizio L, Flores R, et al. Current trends and emerging diagnostic techniques for lung cancer. *BioMed Pharmacother*. (2018) 106:1586–99. doi: 10.1016/j.biopha.2018.07.145
- 5. Tiberi S, du Plessis N, Walzl G, Vjecha MJ, Rao M, Ntoumi F, et al. Tuberculosis: progress and advances in development of new drugs, treatment regimens, and host-directed therapies. *Lancet Infect Dis.* (2018) 18:e183–98. doi: 10.1016/S1473-3099(18) 30110-5
- 6. Dennie C, Bayanati H, Souza CA, Peterson R, Shamji FM. Role of the thoracic radiologist in the evaluation and management of solid and subsolid lung nodules. *Thorac Surg Clin.* (2021) 31:283–92. doi: 10.1016/j.thorsurg.2021.04.004
- 7. Groheux D, Quere G, Blanc E, Lemarignier C, Vercellino L, de Margerie-Mellon C, et al. FDG PET-CT for solitary pulmonary nodule and lung cancer: Literature review. *Diagn Interv Imaging*. (2016) 97:1003–17. doi: 10.1016/j.diii.2016.06.020
- 8. Ng YL, Patsios D, Roberts H, Walsham A, Paul NS, Chung T, et al. CT-guided percutaneous fine-needle aspiration biopsy of pulmonary nodules measuring 10mm or less. Clin Radiol. (2007) 63:272–7. doi: 10.1016/j.crad.2007.09.003
- 9. Dominguez-Konicki L, Karam AR, Furman MS, Grand DJ. CT-guided biopsy of pulmonary nodules ≤10 mm: Diagnostic yield based on nodules' lobar and segmental distribution. *Clin Imag.* (2020) 66:7–9. doi: 10.1016/j.clinimag. 2020.04.040
- 10. Yu H, Yang LT, Zhang Q, Armstrong D, Deen MJ. Convolutional neural networks for medical image analysis: State-of-the-art, comparisons, improvement and perspectives. *Neurocomputing*. (2021) 444:92–110. doi: 10.1016/j.neucom.2020.04.157
- 11. Litjens G, Kooi T, Bejnordi BE, Setio AAA, Ciompi F, Ghafoorian M, et al. A survey on deep learning in medical image analysis. *Med Image Anal.* (2017) 42:60–88. doi: 10.1016/j.media.2017.07.005
- 12. Kalra S, Wen J, Cresswell JC, Volkovs M, Tizhoosh HR. Decentralized federated learning through proxy model sharing. *Nat Commun.* (2023) 14:2899. doi: 10.1038/s41467-023-38569-4
- 13. Henschke CI, Salvatore M, Cham M, Powell CA, DiFabrizio L, Flores R, et al. Federated learning and differential privacy for medical image analysis. *Sci Rep.* (2022) 12:1953. doi: 10.1038/s41598-022-05539-7
- 14. Henschke CI, Salvatore M, Cham M, Powell CA, DiFabrizio L, Flores R, et al. Federated learning for multi-center imaging diagnostics: a simulation study in cardiovascular disease. *Sci Rep.* (2022) 12:3551. doi: 10.1038/s41598-022-07186-4
- 15. Shen T, Zhang J, Jia X, Zhang F, Lv Z, Kuang K, et al. Federated mutual learning: a collaborative machine learning method for heterogeneous data, models, and objectives. *Front Inform Technol Electron Eng.* (2023) 24:1390–402. doi: 10.1631/FITEE.2300098
- Henschke CI, Salvatore M, Cham M, Powell CA, DiFabrizio L, Flores R, et al. HarmoFL: harmonizing local and global drifts in federated learning on heterogeneous medical images. Proc Conf AAAI Artif Intell. (2022) 36:1087–95. doi: 10.1609/ aaai.y36i1.19993

- 17. Henschke CI, Salvatore M, Cham M, Powell CA, DiFabrizio L, Flores R, et al. Federated optimization in heterogeneous networks. *Proc Mach Learn Syst.* (2020) 2:429–50. doi: 10.48550/arXiv.1812.06127
- 18. Karimireddy SP, Kale S, Mohri M, Reddi SJ, Stich SU, Suresh AT, et al. Scaffold: Stochastic controlled averaging for on-device federated learning. *arXiv preprint arXiv*. (2019) 2:1910.06378. doi: 10.48550/arXiv.1910.06378
- 19. Marx A, Chan JK, Coindre JM, Detterbeck F, Girard N, Harris NL, et al. The 2015 world health organization classification of tumors of the thymus: continuity and changes. *J Thorac Oncol.* (2015) 10:1383–95. doi: 10.1097/JTO.00000000000000654
- 20. McKee BJ, Regis SM, McKee AB, Flacke S, Wald C. Performance of ACR lung-RADS in a clinical CT lung screening program. J Am Coll Radiol. (2016) 13:R25–9. doi: 10.1016/j.jacr.2015.12.009
- 21. Henschke CI, Salvatore M, Cham M, Powell CA, DiFabrizio L, Flores R, et al. Sparse Bayesian extreme learning committee machine for engine simultaneous fault diagnosis. *Neurocomputing*. (2016) 174:331–43. doi: 10.1016/j.neucom.2015.02.097
- 22. McMahan B, Moore E, Ramage D, Hampson S, Arcas BA. Communication-efficient learning of deep networks from decentralized data. In: *Artificial intelligence and statistics*. Fort Lauderdale, FL, USA: PMLR (2017). p. 1273–82.
- 23. Henschke CI, Salvatore M, Cham M, Powell CA, DiFabrizio L, Flores R, et al. Federated optimization in heterogeneous networks. *Proc Mach Learn Syst.* (2020) 2:429–50. doi: 10.48550/arXiv.1812.06127
- 24. Henschke CI, Salvatore M, Cham M, Powell CA, DiFabrizio L, Flores R, et al. Fedbn: Federated learning on non-iid features via local batch normalization. *arXiv* preprint arXiv. (2021) 2102:07623. doi: 10.48550/arXiv.2102.07623
- 25. Henschke CI, Salvatore M, Cham M, Powell CA, DiFabrizio L, Flores R, et al. Model-contrastive federated learning. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. Piscataway, NJ, USA: IEEE (2021). p. 10713–22.
- 26. Jiang Y, Konecny J, Rush K, Kannan S. Improving federated learning personalization via model agnostic meta learning. *arXiv.* (2019). doi: 10.48550/arXiv.1909.12488
- 27. Feng B, Chen X, Chen Y, Lu S, Liu K, Li K, et al. Solitary solid pulmonary nodules: a CT-based deep learning nomogram helps differentiate tuberculosis granulomas from lung adenocarcinomas. *Eur Radiol.* (2020) 30:6497–507. doi: 10.1007/s00330-020-07024-z
- 28. Yanagawa M, Johkoh T, Noguchi M, Morii E, Shintani Y, Okumura M, et al. Radiological prediction of tumor invasiveness of lung adenocarcinoma on thin-section CT. *Med (Baltimore).* (2017) 96:e6331. doi: 10.1097/MD.0000000000006331
- 29. Hawkins S, Wang H, Liu Y, Garcia A, Stringfield O, Krewer H, et al. Predicting Malignant nodules from screening CT scans. *J Thorac Oncol.* (2016) 11:2120–8. doi: 10.1016/j.jtho.2016.07.002
- 30. Xu DM, van Klaveren RJ, de Bock GH, Leusveld A, Zhao Y, Wang Y, et al. Limited value of shape, margin and CT density in the discrimination between benign and Malignant screen detected solid pulmonary nodules of the NELSON trial. *Eur J Radiol.* (2008) 68:347–52. doi: 10.1016/j.ejrad.2007.08.027
- 31. Ashraf SF, Yin K, Meng CX, Wang Q, Wang Q, Pu J, et al. Predicting benign, preinvasive, and invasive lung nodules on computed tomography scans using machine learning. *J Thorac Cardiovasc Surg.* (2022) 163:1496–1505.e10. doi: 10.1016/j.jtcvs.2021.02.010
- 32. Wang S, Shi J, Ye Z, Dong D, Yu D, Zhou M, et al. Predicting EGFR mutation status in lung adenocarcinoma on computed tomography image using deep learning. Eur Respir J. (2019) 53:1800986. doi: 10.1183/13993003.00986-2018
- 33. Zhao W, Yang J, Ni B, Bi D, Sun Y, Xu M, et al. Toward automatic prediction of EGFR mutation status in pulmonary adenocarcinoma with 3D deep learning. *Cancer Med.* (2019) 8:3532–43. doi: 10.1002/cam4.2233