

FedEFM: Federated Endovascular Foundation Model with Unseen Data

Tuong Do^{1,2}, Nghia Vu², Tudor Jianu¹, Baoru Huang¹, Minh Vu³, Jionglong Su⁴, Erman Tjiputra²,
Quang D. Tran², Te-Chuan Chiu⁵, Anh Nguyen¹

Abstract—In endovascular surgery, the precise identification of catheters and guidewires in X-ray images is essential for reducing intervention risks. However, accurately segmenting catheter and guidewire structures is challenging due to the limited availability of labeled data. Foundation models offer a promising solution by enabling the collection of similar-domain data to train models whose weights can be fine-tuned for downstream tasks. Nonetheless, large-scale data collection for training is constrained by the necessity of maintaining patient privacy. This paper proposes a new method to train a foundation model in a decentralized federated learning setting for endovascular intervention. To ensure the feasibility of the training, we tackle the unseen data issue using differentiable Earth Mover’s Distance within a knowledge distillation framework. Once trained, our foundation model’s weights provide valuable initialization for downstream tasks, thereby enhancing task-specific performance. Intensive experiments show that our approach achieves new state-of-the-art results, contributing to advancements in endovascular intervention and robotic-assisted endovascular surgery, while addressing the critical issue of data sharing in the medical domain.

I. INTRODUCTION

Endovascular surgery is now usually a minimally invasive procedure that diagnoses and treats vascular diseases with several advantages such as reduced trauma and quick recovery time [1]. During endovascular surgery, surgeons use catheters and guidewires to access arteries. However, this procedure also entails risks such as potential vessel wall damage [2]. Precise identification of catheters and guidewires within X-ray images is crucial for patient safety [3]. The rise of deep learning has played a vital role in improving surgical precision and enhancing patient safety in endovascular intervention [4]. However, accurately segmenting intricate catheters and guidewires in X-ray images remains challenging due to the limited quantity of data [1].

Recently, vision language models have received attention from researchers from various domains [5], [6]. For example, CLIP [7] and ALIGN [8] demonstrate proficiency in cross-modal alignment and zero-shot learning tasks. In the medical domain, EndoFM [5] is developed as a foundation model for endoscopy video analysis. The LVM-Med model [9] is introduced as a foundation model for medical images across multiple modalities. Although these models show promising results on downstream tasks, *most assume that the data can be collected and trained centrally*, which is usually challenging in the medical domain.

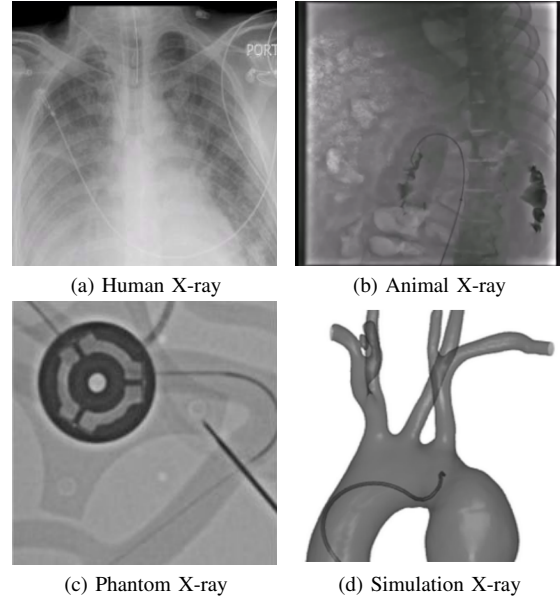


Fig. 1: Different types of endovascular X-ray data.

In practice, collecting large-scale data in the medical domain is not a trivial task due to data privacy [10], [11]. To overcome this limitation, federated learning is emerging as a candidate, enabling the training process to occur between hospital silos without collecting patient data. Despite the advantages of federated learning, current challenges include ensuring convergent training across different silos [12] and heterogeneous data [13]. In endovascular intervention, these challenges primarily stem from data gathered from various sources, hence leading to the domain gap between X-ray data. Fig. 1 shows an example of X-ray images from different endovascular datasets. We observe that due to privacy, endovascular datasets with real human X-ray images are usually small, compared to data collected with animal, silicon phantom models, or from simulation environments [14].

In this paper, our goal is to train a foundation model using diverse endovascular datasets with federated learning. Since we aim to use all possible endovascular data (i.e., from humans, animals, phantoms, etc.), there is an unseen data problem between silos (Fig. 2). To tackle this problem, we propose the Federated Endovascular Foundation Model (FedEFM), a new distillation algorithm using differentiable Earth Mover’s Distance (EMD). Once trained, FedEFM provides crucial initializations for downstream tasks, thereby enhancing task-specific performance. Our approach outperforms existing methods and holds significant potential for application in robotic-assisted endovascular surgery, while effectively maintaining data privacy.

¹Department of Computer Science, University of Liverpool, UK

²AIOZ Ltd., Singapore

³Automation & Control Institute, TU Wien, Austria

⁴Xi’an Jiaotong-Liverpool University, China

⁵National Tsing Hua University, Taiwan

Our contribution can be summarized as below:

- We propose a new method to train a federated endovascular foundation model with unseen data using a multishot distillation technique.
- We collect new datasets for training endovascular foundation models. Our proposed model is verified under several downstream tasks. Our code will be released.

II. LITERATURE REVIEW

Endovascular Intervention. Endovascular intervention has significantly advanced the treatment of various vascular diseases, such as aneurysms and embolisms under X-ray fluoroscopy [15]–[17]. However, these procedures face several challenges due to poor contrast [18], the complexity of anatomical structures [19], and the limited availability of expert-labeled data [20], [21]. Recent research has focused on improving these aspects through advanced imaging technologies and machine learning approaches [22]–[24]. Specifically, the authors in [25] proposed an improved U-Net-based method for guidewire endpoint localization in X-ray images. Recently, FW-Net [26] is proposed to enhance catheter segmentation by leveraging frame-to-frame temporal continuity. While several works focus on traditional tasks, few develop foundation models for endovascular intervention [5], [9]. The main reason is that patient data must be kept private, which becomes a major barrier preventing foundation models from being trained [27]–[29].

Federated Learning. Federated learning has emerged as a promising solution for training machine learning models on decentralized data while preserving data privacy [30]. This approach is particularly beneficial in medical domains, where data sensitivity and privacy concerns are paramount [31]. Although various studies have explored the application of federated learning to train foundation models in the medical field [32], [33], the privacy issue can be handled but not fully resolved [34], [35]. The inherent heterogeneity and non-IID nature of medical data across different institutions present significant challenges [36], [37]. Additionally, the unseen data issue, where certain types of data are present in some datasets but absent in others, complicates the training process and model generalization [38], [39].

Knowledge Distillation with Earth Mover’s Distance. Knowledge distillation involves transferring knowledge from a large, complex model (the teacher) to a smaller, more efficient model (the student) [40]. In the context of federated learning, distillation can be used to enable local models to learn from aggregated global models without sharing raw data [41]. The Earth Mover’s Distance (EMD), also known as the Wasserstein distance, measures the dissimilarity between two probability distributions and is particularly useful for comparing distributions that do not have overlapping support. By leveraging the differentiable EMD, it is possible to align distributions of labels across different models, facilitating better model convergence and knowledge transfer [42], [43]. In this paper, we leverage EMD within a distillation training process to address the unseen label data issue when training endovascular foundation models in federated scenarios.

III. FEDERATED ENDOVASCULAR FOUNDATION MODEL

A. Motivation

We aim to train a federated foundation model for endovascular intervention with all possible types of X-ray data. In practice, each silo (hospital) retains certain data sources that may not be available at other hospitals. The issue arises from the dissimilarity in data corpora across hospitals, i.e., some data are available in one hospital but not in others. Fig. 2 shows an illustration of this problem. Consequently, this leads to the *unseen data* issue that needs to be addressed to ensure the feasibility of the federated training process.

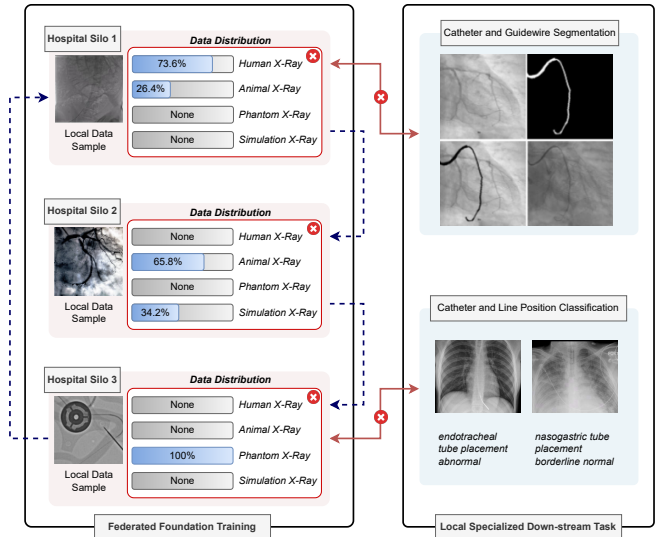


Fig. 2: Unseen data issue visualization. Red lines with crosses indicate insufficient data for training. Blue dotted lines between data silos indicate transferable weights.

According to [44], while federated learning upholds data privacy by prohibiting direct data sharing, it permits the transmission of model weights between connected hospital silos. To take advantage of this characteristic, we propose the Federated Endovascular Foundation Model (FedEFM), a multishot foundation federated distillation algorithm using EMD to ensure the feasibility of learning. Specifically, our approach enables a local silo model to learn from its neighbors’ data and subsequently integrate the acquired knowledge back into the original silo through a distillation mechanism. Unlike other approaches that require a similar label set in both local and global models trained on contributed silos, devices, or servers [44]–[46], our method allows a smooth federated training procedure where hospitals do not need to share their data corpora, thus further improving data privacy. Moreover, once trained, the foundational model’s weights serve as valuable initialization for downstream tasks.

B. Federated Distillation with EMD

We propose Algorithm 1 for training a foundation model within a decentralized federated learning process, effectively addressing the issue of the unseen data problem. Specifically, in the initial round, local model weights θ_i of each i -th

hospital silo is trained using their respective local data ξ_i . Note that N is the maximum number of hospital silos. Within the next communication round, we first perform overseas training where local model weights θ_i of each i -th silo is transmitted to each of their j -th neighbor hospital silo; $j \in \mathcal{N}(i)$ denotes a list of neighbors of i -th silo. This process aims to let local weights θ_i learn knowledge from the data ξ_j of its j -th neighbor silo.

We consider $\theta_{i \rightarrow j}$, the so-called overseas expert, to denote the weight of silo i being transmitted to silo j to learn external knowledge. In $(k+1)$ specific communication round, each transferred weight $\theta_{i \rightarrow j}$ is optimized in j -th silo using the Equation 1.

$$\theta_{i \rightarrow j}(k+1) = \theta_{i \rightarrow j}(k) - \alpha_k \nabla \mathcal{L}_c(\theta_{i \rightarrow j}(k), \xi_j(k)) \quad (1)$$

where ξ is the data in a mini-batch, α is the learning rate, and \mathcal{L}_c is the Cross-Entropy loss used for training a typical foundation classification model [47].

Algorithm 1: Federated knowledge distillation with Earth Mover’s Distance.

Input: Initial weight $\theta_i(0)$ for each silo i ; Maximum training round K .

1 **for** $k = 0$ **to** $K - 1$ **do**

// The loop below is parallel

2 **foreach** *silo* i **do**

$\mathcal{N}(i) \leftarrow$ List of i -th neighbour nodes.

$\xi_i(k) \leftarrow$ Sampling data from local silo i

3 **foreach** *silo* $j \in \mathcal{N}(i)$ **do**

$\xi_j(k) \leftarrow$ Sampling data from the j -th neighbor of silo i

$\theta_{i \rightarrow j} \leftarrow$ Train overseas expert model at j -th silo using Equation 1.

$\hat{\theta}_{i \rightarrow j} \leftarrow \theta_{i \rightarrow j}$ // Collect overseas expert weights from j -th neighbor back to i -th silo.

$\text{EMD}(\theta_i, \hat{\theta}_{i \rightarrow j}) \leftarrow$ Compute Earth Mover’s Distance using Equation 3.

4 $\theta_i(k+1) \leftarrow$ Compute $\mathcal{L}_{\text{MD}}^i$ with Equation 11 and train i -th local model using Equation 2.

Then, we perform knowledge transfer where each learned overseas expert $\theta_{i \rightarrow j}$ from the previous step is transferred back to the i -th silo. Successfully transferred weights is denoted as $\hat{\theta}_{i \rightarrow j}$ which shares values with $\theta_{i \rightarrow j}$.

In the local silo i , the local weight is updated based on both the original weight θ_i and the transferred weights $\hat{\theta}_{i \rightarrow j}$ that is learned from the neighbour silo j . In particular, we aim to find regions that share similarities between two weights using the Earth Mover’s Distance $\text{EMD}(\theta_i, \hat{\theta}_{i \rightarrow j})$. In this way, the distance measures the contribution of transferred weights during distillation, enabling the local silo to learn from its neighbors while avoiding divergence when weight convergence goals differ significantly. Local weights θ_i is

then optimized using:

$$\begin{aligned} \theta_i(k+1) &= \theta_i(k) - \\ \alpha_k \sum_{j \in \mathcal{N}(i)} \text{EMD}(\theta_i, \hat{\theta}_{i \rightarrow j}, k) \nabla \mathcal{L}_{\text{MD}}^i(\theta_i(k), \hat{\theta}_{i \rightarrow j}(k), \xi_i(k)) \end{aligned} \quad (2)$$

where \mathcal{L}_{MD} is the distillation loss (Equation 11), and $\mathcal{N}(i)$ indicates in-neighbors of silo i .

Differentiable Earth Mover’s Distance. Assume that the input sample ξ_i from i -th local silo passes through the foundation architecture θ_i to generate the dense representation $\mathbf{U} \in \mathbb{R}^{H \times W \times C}$, where H and W denote the spatial size of the feature map and C is the feature dimension. In a parallel manner, $\mathbf{V} \in \mathbb{R}^{H \times W \times C}$ also denotes the dense representation when ξ_i passes through $\hat{\theta}_{i \rightarrow j}$.

Under Earth Mover circumstance, \mathbf{V} represents suppliers transporting goods to demanders \mathbf{U} . Then, EMD between two feature sets $\mathbf{U} = \{u_1, u_2, \dots, u_{HW}\}$ and $\mathbf{V} = \{v_1, v_2, \dots, v_{HW}\}$ can be computed as:

$$\text{EMD}(\theta_i, \hat{\theta}_{i \rightarrow j}) = \text{EMD}(\mathbf{U}, \mathbf{V}) = \sum_{p=1}^{HW} \sum_{q=1}^{HW} (1 - c_{pq}) \tilde{x}_{pq}. \quad (3)$$

where \tilde{x} is conducted from optimal matching flow $\tilde{X} = \{x_1, x_2, \dots, x_{pq}\}$ for each sample pair of two sets \mathbf{U} and \mathbf{V} ; c_{pq} is the cost per unit transported from supplier to demander and is obtained by computing the pairwise distance between embedding nodes $u_p \in \mathbf{U}$ and $v_q \in \mathbf{V}$.

The cost per unit c_{pq} is computed as below and also plays a virtual role in computing the optimal matching flow:

$$c_{pq} = 1 - \frac{u_p^T v_q}{\|u_p\| \|v_q\|} \quad (4)$$

where nodes with similar representations tend to generate small matching costs between each other. Then, the optimal matching flow \tilde{X} is conducted by optimizing \tilde{x} as below:

$$\begin{aligned} &\underset{x}{\text{minimize}} && \sum_{p=1}^{HW} \sum_{q=1}^{HW} c_{pq} x_{pq} \\ &\text{subject to} && x_{pq} > 0, \quad p = 1, \dots, HW, \quad q = 1, \dots, HW \\ &&& \sum_{p=1}^{HW} x_{pq} = v_q, \quad q = 1, \dots, HW \\ &&& \sum_{q=1}^{HW} x_{pq} = u_p, \quad p = 1, \dots, HW \end{aligned} \quad (5)$$

Here, EMD seeks an optimal matching \tilde{X} between suppliers and demanders such that the overall matching cost is minimized. The global optimal matching flows \tilde{X} can be achieved by solving a Linear Programming problem (LP). For the sake of completeness, we transform the optimization in Equation 5 to a compact matrix form:

$$\begin{aligned} &\underset{x}{\text{minimize}} && c(\theta)^T x \\ &\text{subject to} && G(\theta)x \leq h(\theta), \\ &&& A(\theta)x = b(\theta). \end{aligned} \quad (6)$$

Here $x \in \mathbb{R}^{HW \times HW}$ is our optimization variable. $Ax = b$ represents the equality constraint and $Gx \leq h$ denotes

the inequality constraint in Equation 5. Accordingly, the Lagrangian of the LP problem in Equation 6 is given by:

$$L(\theta, x, \nu, \lambda) = c^T x + \lambda^T (Gx - h) + \nu^T (Ax - b), \quad (7)$$

where ν denotes the dual variables on the equality constraints and $\lambda \geq 0$ denotes the dual variables on the inequality constraints. Following the KKT conditions, we obtain the optimum $(\tilde{x}, \tilde{\nu}, \tilde{\lambda})$ of the objective function by solving $g(\theta, \tilde{x}, \tilde{\nu}, \tilde{\lambda}) = 0$ with primal-dual interior point methods, where

$$g(\theta, x, \nu, \lambda) = \begin{bmatrix} \nabla_{\theta} L(\theta, x, \nu, \lambda) \\ \mathbf{diag}(\lambda)(G(\theta)x - h(\theta)) \\ A(\theta)x - b(\theta) \end{bmatrix}. \quad (8)$$

Then, with the theorem below, we can derive the gradients of the LP parameters.

Suppose $g(\theta, \tilde{x}, \tilde{\nu}, \tilde{\lambda}) = 0$. Then, when all derivatives exist, the partial Jacobian of \tilde{x} with respect to θ at the optimal solution $(\tilde{\lambda}, \tilde{\nu}, \tilde{x})$, namely $J_{\theta}\tilde{x}$, can be obtained by satisfying:

$$J_{\theta}\tilde{x} = - \left(J_x g(\theta, \tilde{\lambda}, \tilde{\nu}, \tilde{x}) \right)^{-1} J_{\theta} g(\theta, \tilde{x}, \tilde{\nu}, \tilde{\lambda}). \quad (9)$$

Then, applying to the KKT conditions, the (partial) Jacobian with respect to θ can be defined as:

$$J_{\theta} g(\theta, \tilde{\lambda}, \tilde{\nu}, \tilde{x}) = \begin{bmatrix} J_{\theta} \nabla_x L(\theta, \tilde{x}, \tilde{\nu}, \tilde{\lambda}) \\ \mathbf{diag}(\tilde{\lambda}) J_{\theta} (G(\theta)x - h(\theta)) \\ J_{\theta} (A(\theta)\tilde{x} - b(\theta)) \end{bmatrix}. \quad (10)$$

After obtaining the optimal \tilde{x} , we can derive a closed-form gradient for θ , enabling efficient backpropagation without altering the optimization path.

C. Distillation Loss

Assume that each $\theta_{i \rightarrow j}$ is a teacher transmitted from j -th neighbor silo. The distillation loss of i -th silo \mathcal{L}_{MD}^i based on student model loss is designed as:

$$\mathcal{L}_{MD}^i = \beta T^2 \sum_{j=1}^{\mathcal{N}(i)} (\mathcal{L}_c(Q_S^{\tau}, Q_T^{\tau})) + (1 - \beta) \mathcal{L}_c(Q_S, y_{true}^i) \quad (11)$$

where Q_S is the standard softmax output of the local student; y_{true}^i is the ground-truth labels; β is a hyper-parameter for controlling the importance of each loss component; Q_S^{τ}, Q_T^{τ} are the softened outputs of the i -th local student and the j -th overseas teachers using the same temperature parameter T [40], which are computed as follows:

$$Q_k^{\tau} = \frac{\exp(l_k/T)}{\sum_k \exp(l_k/T)} \quad (12)$$

where the logit l is outputted from the pre-final layers for both teacher and student models. Besides, as stated in Equation 2, the objective function computed for each j -th contributed transferrable weights is controlled by the corresponding EMD to ensure the learning convergence.

When the training in all silos is completed in each communication round, local model weights in all silos are aggregated to obtain global weights $\Theta = \sum_{i=0}^{N-1} \vartheta_i \theta_i$, which are further utilized for downstream fine-tuning. Note that $\vartheta \in \{0, 1\}$ indicates accumulation status.

IV. EXPERIMENTS

A. Data Preparation

Robotic Setup. To collect large-scale X-ray images, we employ a robotic platform and a full-size silicon phantom. A surgeon uses a master device joystick to control a follower robot for cannulating three arteries: the left subclavian (LSA), left common carotid (LCCA), and right common carotid (RCCA). Fig. 3 shows an overview of our robotics setup. During each catheterization procedure, the surgeon activates the X-ray fluoroscopy using a pedal in the operating room. The experiments are conducted using the Epsilon X-ray Generator. We develop a real-time image grabber to transmit the video feed of the surgical scene to a workstation, a computer-based device equipped with an 8-Core ARM v8.2 64-bit CPU. Overall, we collect and label 4,700 new X-ray images to create our EIPhantom dataset.

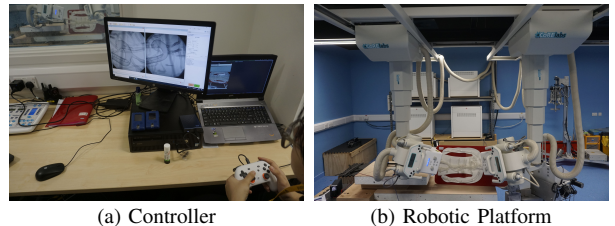


Fig. 3: Data collection with endovascular robot.

Simulation Data. Apart from X-ray images collected from our real robot, we also collect an EISimulation dataset from the CathSim simulator [2] for simulated X-ray images. We manually label both data from the robot and CathSim simulator to use them in downstream tasks. We note that the datasets used to train the foundation model are not being used in downstream endovascular understanding tasks.

Phase	Dataset	#Frames
Federated Foundation Training	CathAction [14]	500,000
	VESSEL12 [48]	12,892
	Drive [49]	8,028
	SenNet [50]	7,436
	Medical Decathlon [51]	442
Downstream Fine-tuning	EISimulation (ours)	1,683
	EIPhantom (ours)	4,710
	RANZCR [52]	33,664
	CathAnimal [53]	25,000

TABLE I: X-ray datasets used in our experiments.

Dataset Summary. Table I summarises datasets related to endovascular intervention [48]–[51], [53] we use in this paper. All datasets cover different endovascular procedures with X-ray images as the main modality. The data are collected from diverse sources, including human/animal studies, human phantoms, and simulated environments.

Learning Scenario	Method	Accuracy	Avg. Cycle Time (ms)	
CLL	CLIP [7]	67.5	-	
	SAM [55]	72.4	-	
	LVM-Med [9]	98.8	-	
CFL	FedAvg [54]	80.9	57.7	
	MOON [56]	85.2	69.2	
	STAR [57]	82.4	63.8	
	FedEFM (ours)	w/o EMD	84.7	42.5
		w EMD	98.2	61.1
DFL	MATCHA [58]	42.4	43.4	
	RING [44]	52.2	73.2	
	CDL [45]	78.5	59.9	
	FedEFM (ours)	w/o EMD	72.4	47.3
		w EMD	97.5	62.1

TABLE II: Foundation model performance comparison.

B. Federated Endovascular Foundation Model Validation

Setup. We first validate our proposed method (FedEFM) and compare it with different foundation models in different learning scenarios. In particular, we consider three scenarios, including Centralized Local Learning (CLL), Client-server Federated Learning (CFL) [54], and Decentralized Federated Learning (DFL) [44]. We note that CLL is the traditional training scenario (i.e., no federated learning) where data are merged for local training. Multiple algorithms have been conducted for the comparison purpose, including CLIP [7], SAM [55], LVM-Med [9], FedAvg [54], MOON [56], STAR [57], MATCHA [58], RING [44], and CDL [45]. We use ViT [59] backbone in all benchmarking algorithms and train on datasets for the training phase in Table I. Note that our default setup is maintained at 100% unseen label corpus.

Results. Table II shows the comparison with different algorithms on multiple learning scenarios. When we train ViT in CFL and DFL setup using FedAvg and MATCHA, the accuracy is only 80.9% and 42.4%, respectively, reflecting the inherent challenges in federated learning. Applying our proposed FedEFM method resulted in a substantial accuracy improvement to 98.2% and 97.5%. These results show that our proposed method can obtain competitive results even compared with the centralized training that can gather all data and only has a minor cycle time trade-off compared with most of the federated learning methods.

C. Fine-tuning Results

Setup. We use ViT backbone and fine-tune it using our FedEFM and different foundation models, including, CLIP [7], SAM [55], and LVM-Med [9]. Note that, all models are evaluated under segmentation and classification tasks in endovascular intervention.

Evaluation Metrics. We use the metrics in [9], [53] to evaluate the performance of the trained foundation model in downstream tasks. Specifically, we use Accuracy (%) for the classification task; 2D Dice score, mIoU, and Jaccard metric are used for the segmentation task. For the segmentation task, we compare on our collected EIPhantom, EISimulation dataset, and CathAnimal [53]. In the classification task, we benchmark using the RANZCR dataset [52].

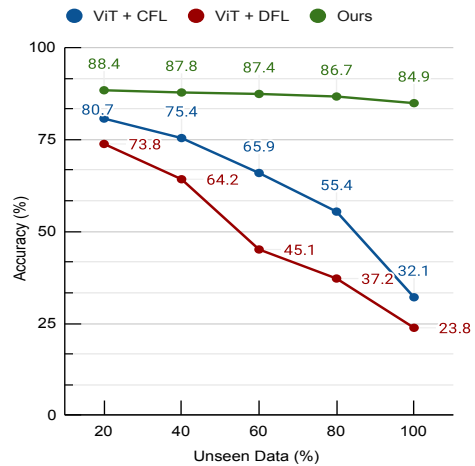


Fig. 4: Results with different unseen data proportions.

Results. Table III shows the comparison between our method and other foundation models. This table shows that the ViT backbone under our proposed algorithm outperforms other models with a clear margin. Furthermore, models trained on medical data such as LVM-Med [9] and our FedEFM archive better results compared with models trained on non-medical data such as CLIP [7] and SAM [55]. This shows that developing a domain-specific foundation model is important in the medical domain.

D. Ablation Study

Unseen Data Proportion Analysis. Fig. 4 presents an analysis of our method under different percentages of unseen data. In this experiment, we assume that each silo (hospital) only keeps an amount of data (e.g., human/animal/simulated X-ray) where their data corpus only shares the similarity in a given percentage. A 100% unseen data corpus means that the data of each hospital silo have no similarity in their data types compared to others. As the percentage of unseen data types increases, we observe a notable decline in the accuracy of the baseline on CFL and DFL scenarios. However, our proposed approach demonstrates remarkable resilience to unseen data, maintaining high accuracy even when confronted with a higher percentage of unfamiliar semantic data. In specific instances, when all data labels are unseen (100%), ViT under CFL and DFL scenarios exhibit significantly lower accuracies at 32.1% and 23.8%, respectively. In contrast, our approach achieves an accuracy of 84.9%, showcasing its effectiveness in handling unseen data.

Backbones Analysis. We verify the stability of our method on different networks, including UNet [60], TransUNet [61], and SwinUnet [62], and ViT [59] under federated learning scenario. Table IV shows the performance of the different backbones when we fine-tune them using our FedEFM. Table IV demonstrates that using our foundation model to initialize the weights of those backbones significantly improves the results. These results validate the effectiveness of our training process in addressing the unseen data problem, and our FedEFM is useful for different backbones in endovascular downstream tasks.

Models	Segmentation									Classification
	EIPhantom			CathAnimal [53]			EISimulation			RANZCR [52]
	Dice	mIoU	Jaccard	Dice	mIoU	Jaccard	Dice	mIoU	Jaccard	Accuracy
CLIP [7]	46.7	23.8	43.5	59.1	43.5	52.1	52.4	37.3	32.0	60.4
SAM [55]	47.3	29.9	50.7	62.2	41.1	58.8	77.9	30.5	51.1	55.4
LVM-Med [9]	56.2	31.8	51.5	66.6	52.5	70.7	70.9	49.1	61.2	62.3
FedEFM (ours)	63.1	35.5	57.1	67.2	50.1	71.8	82.9	63.2	81.2	67.9

TABLE III: Fine-tuning results of different foundation models on endovascular segmentation and classification tasks.

Backbones	Initialize	Segmentation									Classification
		EIPhantom			CathAnimal [53]			EISimulation			RANZCR [52]
		Dice	mIoU	Jaccard	Dice	mIoU	Jaccard	Dice	mIoU	Jaccard	Accuracy
U-Net [60]	From-scratch	48.1	20.2	50.2	52.5	42.7	59.4	51.1	22.5	66.6	49.4
	Fine-tuned	52.1	30.5	51.7	66.9	48.3	65.4	56.4	27.9	72.9	56.0
TransUnet [61]	From-scratch	46.7	30.1	49.9	51.2	44.4	59.5	62.2	19.7	68.3	52.9
	Fine-tuned	58.9	34.0	55.9	54.3	46.2	64.4	80.2	22.3	72.2	58.3
SwinUnet [62]	From-scratch	47.3	32.2	51.7	50.6	43.4	58.5	60.8	19.1	67.2	55.7
	Fine-tuned	58.5	34.3	56.0	66.2	48.4	65.5	76.8	19.0	68.9	62.5
ViT [59]	From-scratch	50.9	30.2	50.8	59.4	44.7	60.0	72.1	61.4	74.5	60.6
	Fine-tuned	63.1	35.5	57.1	67.2	50.1	71.8	82.9	63.2	81.2	67.9

TABLE IV: Performance of different backbones when using our FedEFM for fine-tuning.

Qualitative Results. Fig. 5 illustrates the catheter and guidewire segmentation results of fine-tuning ViT on our method and different foundation models. The visualization portrays that our method excels in accurately delineating the catheter and guidewire structures, showcasing superior segmentation performance compared to other approaches. This figure further confirms that we can successfully train a federated endovascular foundation model without collecting users’ data and the trained foundation model is useful for the downstream segmentation task.

E. Limitations

While our proposed approach demonstrates significant potential, it is subject to certain limitations that warrant further investigation. Firstly, the requirement for additional weight exchange among silos extends the overall training time. However, this limitation is mitigated to some extent by the higher convergence speed of our method compared to other approaches. Additionally, our method is designed for deployment in silos with strong GPU computing resources, but the varying hardware capabilities present in many real-world federated learning networks necessitate further examination. Overcoming these limitations will open new research in federated foundation learning for endovascular interventions and other medical applications. Furthermore, addressing the challenges of managing heterogeneous data distributions and ensuring robust data privacy remains a critical focus. Moving forward, we plan to extend our approach to robotic-assisted endovascular surgery and other areas, such as pathology, to further investigate the application of federated foundation models in medical imaging and robotic systems.

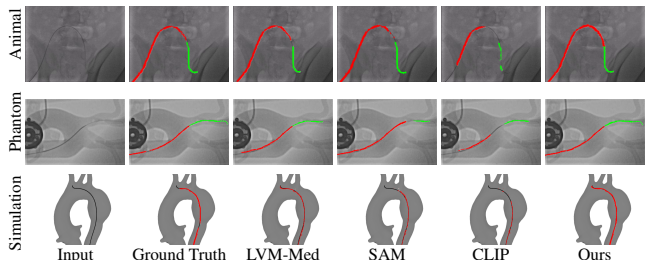


Fig. 5: Catheter and guidewire segmentation between methods. Red lines are catheters and green ones are guidewires.

V. CONCLUSION

We present a new approach to train an endovascular foundation model in a federated learning setting, leveraging differentiable Earth Mover’s Distance and knowledge distillation to handle the unseen data issue. Our method ensures that once the foundational model is trained, its weights can be effectively fine-tuned for downstream tasks, thereby enhancing performance. Our approach achieves state-of-the-art results and contributes to the field of endovascular intervention, particularly by addressing the critical issue of data sharing in the medical domain. By enabling weight exchange among local silos and fostering knowledge transfer, our method improves model generalization while preserving data privacy. Experimental results across various endovascular imaging tasks validate the efficacy of our approach, demonstrating its potential for application in privacy-sensitive medical domains. We will release our implementation and trained models to facilitate reproducibility and further research.

REFERENCES

- [1] A. Ramadani, M. Bui, T. Wendler, H. Schunkert, P. Ewert, and N. Navab, "A survey of catheter tracking concepts and methodologies," *Medical Image Analysis*, 2022.
- [2] T. Jianu, B. Huang, M. N. Vu, M. E. Abdelaziz, S. Fichera, C.-Y. Lee, P. Berthet-Rayne, F. R. y Baena, and A. Nguyen, "Cathsim: an open-source simulator for endovascular intervention," *IEEE Transactions on Medical Robotics and Bionics*, 2024.
- [3] G. M. Pereira Junior, A. Souza Alvarenga, C. R. Almeida Felipe, A. Vale Monteiro, L. R. Rezende, and M. G. Moreira Guimarães Penido, "Use of ultrasound to confirm guidewire position in hemodialysis catheter implantation," *Nephrology*, 2022.
- [4] B. Huang, Y. Hu, A. Nguyen, S. Giannarou, and D. S. Elson, "Detecting the sensing area of a laparoscopic probe in minimally invasive cancer surgery," in *MICCAI*, 2023.
- [5] Z. Wang, C. Liu, S. Zhang, and Q. Dou, "Foundation model for endoscopy video analysis via large-scale self-supervised pre-train," in *MICCAI*, 2023.
- [6] Y. Zhang, J. Gao, M. Zhou, *et al.*, "Text-guided foundation model adaptation for pathological image classification," in *MICCAI*, 2023.
- [7] A. Radford, J. W. Kim, C. Hallacy, *et al.*, "Learning transferable visual models from natural language supervision," in *ICML*, 2021.
- [8] C. Jia, Y. Yang, Y. Xia, Y.-T. Chen, Z. Parekh, H. Pham, Q. Le, Y.-H. Sung, Z. Li, and T. Duerig, "Scaling up visual and vision-language representation learning with noisy text supervision," in *ICML*, 2021.
- [9] D. M. Nguyen, H. Nguyen, N. T. Diep, *et al.*, "Lvm-med: Learning large-scale self-supervised vision models for medical imaging via second-order graph matching," *arXiv*, 2023.
- [10] G. A. Kaissis, M. R. Makowski, D. Rückert, and R. F. Braren, "Secure, privacy-preserving and federated machine learning in medical imaging," *Nature Machine Intelligence*, 2020.
- [11] B. Huang, A. Nguyen, S. Wang, Z. Wang, E. Mayer, D. Tuch, K. Vyas, S. Giannarou, and D. S. Elson, "Simultaneous depth estimation and surgical tool segmentation in laparoscopic images," *IEEE transactions on medical robotics and bionics*, 2022.
- [12] M. Jiang, Z. Wang, and Q. Dou, "Harmoff: Harmonizing local and global drifts in federated learning on heterogeneous medical images," in *AAAI*, 2022.
- [13] Q. Zhou and G. Zheng, "Fedcontrast-gpa: Heterogeneous federated optimization via local contrastive learning and global process-aware aggregation," in *MICCAI*, 2023.
- [14] B. Huang, T. Vo, C. Kongtongvattana, G. Dagnino, D. Kundrat, W. Chi, M. Abdelaziz, T. Kwok, T. Jianu, T. Do, *et al.*, "Cathaction: A benchmark for endovascular intervention understanding," *arXiv preprint arXiv:2408.13126*, 2024.
- [15] S. A. Baert, M. A. Viergever, and W. J. Niessen, "Guide-wire tracking during endovascular interventions," *IEEE Transactions on Medical Imaging*, vol. 22, no. 8, pp. 965–972, 2003.
- [16] Y.-J. Zhou, X.-L. Xie, X.-H. Zhou, S.-Q. Liu, G.-B. Bian, and Z.-G. Hou, "A real-time multifunctional framework for guidewire morphological and positional analysis in interventional x-ray fluoroscopy," *IEEE Transactions on Cognitive and Developmental Systems*, vol. 13, no. 3, pp. 657–667, 2020.
- [17] M. Gherardini, E. Mazomenos, A. Menciassi, and D. Stoyanov, "Catheter segmentation in x-ray fluoroscopy using synthetic data and transfer learning with light u-nets," *Computer Methods and Programs in Biomedicine*, vol. 192, p. 105420, 2020.
- [18] P. Moore, "Mri-guided congenital cardiac catheterization and intervention: The future?" *Catheterization and cardiovascular interventions*, vol. 66, no. 1, pp. 1–8, 2005.
- [19] F. O. Efthymiou, S. K. Kakkos, V. I. Metaxas, C. P. Dimitroukas, K. G. Moulakakis, S. I. Papadoulas, N. K. Kouri, A. L. Tsimpoukis, K. M. Nikolakopoulos, C. P. Papageorgopoulou, *et al.*, "Factors influencing fluoroscopy time in endovascular treatment of abdominal aneurysms: a retrospective study," *Radiation Protection Dosimetry*, vol. 199, no. 5, pp. 443–452, 2023.
- [20] K. Breininger, T. Würfl, T. Kurzdorfer, S. Albarqouni, M. Pfister, M. Kowarschik, N. Navab, and A. Maier, "Multiple device segmentation for fluoroscopic imaging using multi-task learning," in *Intravascular Imaging and Computer Assisted Stenting and Large-Scale Annotation of Biomedical Data and Expert Label Synthesis: 7th Joint International Workshop, CVII-STENT 2018 and Third International Workshop, LABELS 2018, Held in Conjunction with MICCAI 2018, Granada, Spain, September 16, 2018, Proceedings 3*. Springer, 2018, pp. 19–27.
- [21] K. Breininger, S. Albarqouni, T. Kurzdorfer, M. Pfister, M. Kowarschik, and A. Maier, "Intraoperative stent segmentation in x-ray fluoroscopy for endovascular aortic repair," *International journal of computer assisted radiology and surgery*, vol. 13, pp. 1221–1231, 2018.
- [22] Y. Ma, D. Zhou, L. Ye, R. J. Housden, A. Fazili, and K. S. Rhode, "A tensor-based catheter and wire detection and tracking framework and its clinical applications," *IEEE Transactions on Biomedical Engineering*, vol. 69, no. 2, pp. 635–644, 2021.
- [23] A. Ranne, Y. Velikova, N. Navab, *et al.*, "Aiaresseg: Catheter detection and segmentation in interventional ultrasound using transformers," *arXiv preprint arXiv:2309.14492*, 2023.
- [24] Z. Mei, H. Wang, S. Pan, H. Chen, J. Wei, Q. Zhang, J. Mao, G. Liu, and Y. Zhao, "Real time detection and tracking of guide wire/catheter for interventional embolization robot based on deep learning," in *2023 IEEE International Conference on Mechatronics and Automation (ICMA)*. IEEE, 2023, pp. 778–783.
- [25] W. Du, G. Yi, O. M. Omisore, W. Duan, X. Chen, T. Akinyemi, J. Liu, B.-G. Lee, and L. Wang, "Guidewire endpoint detection based on pixel-adjacent relation during robot-assisted intravascular catheterization: In vivo mammalian models," *Advanced Intelligent Systems*, vol. 6, no. 4, p. 2300687, 2024.
- [26] A. Nguyen, D. Kundrat, G. Dagnino, W. Chi, E. Abdelaziz, Y. Guo, Y. Ma, T. Kwok, C. Riga, and G.-Z. Yang, "End-to-end real-time catheter segmentation with optical flow-guided warping during endovascular intervention," in *ICRA*, 2020.
- [27] Z. Wang, Z. Wu, D. Agarwal, and J. Sun, "Medclip: Contrastive learning from unpaired medical images and text," *arXiv preprint arXiv:2210.10163*, 2022.
- [28] J. Ma, Y. He, F. Li, L. Han, C. You, and B. Wang, "Segment anything in medical images," *Nature Communications*, vol. 15, pp. 1–9, 2024.
- [29] T. Koleilat, H. Asgariandehkordi, H. Rivaz, and Y. Xiao, "Medclip-sam: Bridging text and image towards universal medical image segmentation," *arXiv preprint arXiv:2403.20253*, 2024.
- [30] P. M. Mammen, "Federated learning: Opportunities and challenges," *arXiv preprint arXiv:2101.05428*, 2021.
- [31] M. Jiang, Y. Zhong, A. Le, X. Li, and Q. Dou, "Client-level differential privacy via adaptive intermediary in federated medical imaging," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2023, pp. 500–510.
- [32] H. Chen, Y. Zhang, D. Krompass, J. Gu, and V. Tresp, "Feddat: An approach for foundation model finetuning in multi-modal heterogeneous federated learning," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 10, 2024, pp. 11 285–11 293.
- [33] Y. Liu, G. Luo, and Y. Zhu, "Fedfms: Exploring federated foundation models for medical image segmentation," *arXiv preprint arXiv:2403.05408*, 2024.
- [34] R. Xue, K. Xue, B. Zhu, X. Luo, T. Zhang, Q. Sun, and J. Lu, "Differentially private federated learning with an adaptive noise mechanism," *IEEE Transactions on Information Forensics and Security*, 2023.
- [35] S. D. Okegbile, J. Cai, H. Zheng, J. Chen, and C. Yi, "Differentially private federated multi-task learning framework for enhancing human-to-virtual connectivity in human digital twin," *IEEE Journal on Selected Areas in Communications*, 2023.
- [36] X. Lin, J. Wu, J. Li, C. Sang, S. Hu, and M. J. Deen, "Heterogeneous differentially-private federated learning: Trading privacy for utility truthfully," *IEEE Transactions on Dependable and Secure Computing*, vol. 20, no. 6, pp. 5113–5129, 2023.
- [37] J. Wang, X. Yang, S. Cui, L. Che, L. Lyu, D. D. Xu, and F. Ma, "Towards personalized federated learning via heterogeneous model reassembly," *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [38] Q. Liu, C. Chen, J. Qin, Q. Dou, and P.-A. Heng, "Feddg: Federated domain generalization on medical image segmentation via episodic learning in continuous frequency space," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 1013–1023.
- [39] M. Tölle, F. Navarro, S. Eble, I. Wolf, B. Menze, and S. Engelhardt, "Funavg: Federated uncertainty weighted averaging for datasets with diverse labels," *arXiv preprint arXiv:2407.07488*, 2024.
- [40] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," *arXiv*, 2015.
- [41] E. Jeong, S. Oh, H. Kim, J. Park, M. Bennis, and S.-L. Kim, "Communication-efficient on-device machine learning: Federated dis-

- tillation and augmentation under non-iid private data,” *arXiv preprint arXiv:1811.11479*, 2018.
- [42] C. Zhang, Y. Cai, G. Lin, and C. Shen, “Deepemd: Few-shot image classification with differentiable earth mover’s distance and structured classifiers,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 12 203–12 213.
- [43] Q. Zhao, Z. Yang, and H. Tao, “Differential earth mover’s distance with its applications to visual tracking,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 32, no. 2, pp. 274–287, 2008.
- [44] O. Marfoq, C. Xu, G. Neglia, and R. Vidal, “Throughput-optimal topology design for cross-silo federated learning,” *NIPS*, 2020.
- [45] T. Do, B. X. Nguyen, H. Nguyen, E. Tjiputra, Q. D. Tran, and A. Nguyen, “Addressing non-iid problem in federated autonomous driving with contrastive divergence loss,” in *ICRA*, 2024.
- [46] A. Nguyen, N. Nguyen, K. Tran, E. Tjiputra, and Q. Tran, “Autonomous navigation in complex environments with deep multimodal fusion network,” in *IROS*, 2020.
- [47] G. E. Hinton, P. Dayan, B. J. Frey, and R. M. Neal, “The “wake-sleep” algorithm for unsupervised neural networks,” *Science*, 1995.
- [48] R. D. Rudyanto, S. Kerkstra, E. M. Van Rikxoort, *et al.*, “Comparing algorithms for automated vessel segmentation in computed tomography scans of the lung: the vessel12 study,” *Medical image analysis*, 2014.
- [49] J. Staal, M. D. Abràmoff, M. Niemeijer, M. A. Viergever, and B. Van Ginneken, “Ridge-based vessel segmentation in color images of the retina,” *TMI*, 2004.
- [50] C. Walsh, P. Tafforeau, W. Wagner, *et al.*, “Imaging intact human organs with local resolution of cellular structures using hierarchical phase-contrast tomography,” *Nature methods*, 2021.
- [51] M. Antonelli, A. Reinke, *et al.*, “The medical segmentation decathlon,” *Nature communications*, 2022.
- [52] L. Hansen, M. Sieren, M. Hobe, *et al.*, “Radiographic assessment of cvc malpositioning: How can ai best support clinicians?” in *Medical Imaging with Deep Learning*, 2021.
- [53] C. Kongtongvattana, B. Huang, J. Kang, H. Nguyen, O. Olufemi, and A. Nguyen, “Shape-sensitive loss for catheter and guidewire segmentation,” *arXiv*, 2023.
- [54] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, “Communication-efficient learning of deep networks from decentralized data,” in *Artificial intelligence and statistics*, 2017.
- [55] A. Kirillov, E. Mintun, N. Ravi, *et al.*, “Segment anything,” *arXiv*, 2023.
- [56] Q. Li, B. He, and D. Song, “Model-contrastive federated learning,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 10 713–10 722.
- [57] U. Brandes, “On variants of shortest-path betweenness centrality and their generic computation,” *Social networks*, 2008.
- [58] J. Wang, A. K. Sahu, Z. Yang, G. Joshi, and S. Kar, “Matcha: Speeding up decentralized sgd via matching decomposition sampling,” in *2019 Sixth Indian Control Conference (ICC)*. IEEE, 2019, pp. 299–300.
- [59] A. Dosovitskiy, L. Beyer, A. Kolesnikov, *et al.*, “An image is worth 16x16 words: Transformers for image recognition at scale,” *ICLR*, 2021.
- [60] O. Ronneberger, P. Fischer, and T. Brox, “U-net: Convolutional networks for biomedical image segmentation,” in *MICCAI*, 2015.
- [61] J. Chen, Y. Lu, Q. Yu, *et al.*, “Transunet: Transformers make strong encoders for medical image segmentation,” *arXiv*, 2021.
- [62] H. Cao, Y. Wang, J. Chen, D. Jiang, X. Zhang, Q. Tian, and M. Wang, “Swin-unet: Unet-like pure transformer for medical image segmentation,” in *ECCV*, 2022.