# FedMix: Mixed Supervised Federated Learning for Medical Image Segmentation

Jeffry Wicaksana, Zengqiang Yan, Dong Zhang, Xijie Huang, Huimin Wu, Xin Yang, and Kwang-Ting Cheng [*][†]

May 5, 2022

## Abstract

The purpose of federated learning is to enable multiple clients to jointly train a machine learning model without sharing data. However, the existing methods for training an image segmentation model have been based on an unrealistic assumption that the training set for each local client is annotated in a similar fashion and thus follows the same image supervision level. To relax this assumption, in this work, we propose a label-agnostic unified federated learning framework, named FedMix, for medical image segmentation based on mixed image labels. In FedMix, each client updates the federated model by integrating and effectively making use of all available labeled data ranging from strong pixel-level labels, weak bounding box labels, to weakest image-level class labels. Based on these local models, we further propose an adaptive weight assignment procedure across local clients, where each client learns an aggregation weight during the global model update. Compared to the existing methods, FedMix not only breaks through the constraint of a single level of image supervision, but also can dynamically adjust the aggregation weight of each local client, achieving rich yet discriminative feature representations. To evaluate its effectiveness, experiments have been carried out on two challenging medical image segmentation tasks, *i.e.*, breast tumor segmentation and skin lesion segmentation. The results validate that our proposed FedMix outperforms the state-of-the-art methods by a large margin[1].

**Keywords**: Federated learning, mixed supervisions, medical image segmentation, pseudo labeling, adaptive weight aggregation

[*]J. Wicaksana, D. Zhang, X. Huang, H. Wu and K. -T. Cheng are with the Department of Electronic and Computer Engineering, The Hong Kong University of Science and Technology, Kowloon, Hong Kong. E-mail: {jwicaksana, dongz, xhuangbs, hwubl, timcheng}@ust.hk

[†]Z. Yan and X. Yang are with the School of Electronic Information and Communications, Huazhong University of Science and Technology, Wuhan, China. E-mail: {z_yan ,xinyang2014}@hust.edu.cn

[1]The code is available at: https://github.com/Jwicaksana/FedMix

(a) Input Image     (b) Image-Level Class Labels     (c) Bounding Box Labels     (d) Pixel-Level Labels

Weak Supervisions     (more labour inputs)     Strong Supervisions
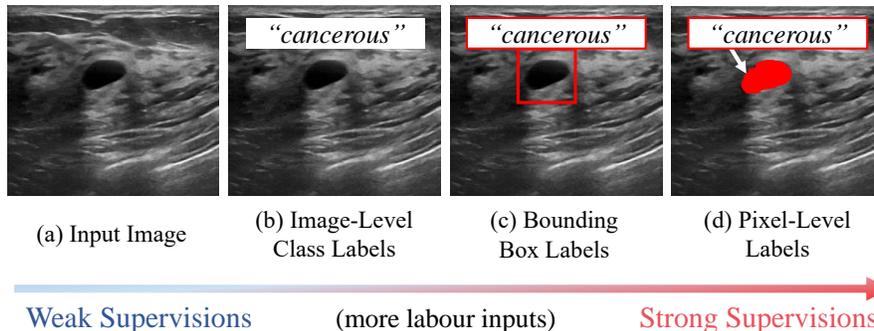
**Figure 1:** Examples of different levels of medical image labels, where the image-level class labels in (b) contain only the lesion category. The bounding box labels in (c) contain not only the lesion category, but also a coarse location. The pixel-level labels in (d) contain both the lesion category and location information of each pixel, which is a strong image supervision. Although strong image supervisions are more informative, they are very expensive to obtain. The utilization of some easily obtained image supervisions is beneficial in practice.

## 1 Introduction

Medical image segmentation is a representative task for image content analysis supporting computer aided diagnosis, which can not only recognize the lesion category, but also locate the specific areas [1]. In the past few years, this task has been extensively studied and applied in a wide range of underlying scenarios, *e.g.*, lung nodule segmentation [2], skin lesion boundary detection [4], and COVID-19 lesion segmentation [3].

The optimization of deep learning models usually relies on a vast amount of training data [5]. For example, for a fully-supervised semantic segmentation model, the ideal scenario is that we can collect the pixel-level annotated images as much as possible from diverse sources. However, this scenario is almost infeasible due to the following two reasons: *i*) the strict sharing protocol of sensitive patient information between medical institutions and *ii*) the exceedingly high pixel-level annotation cost. As the expert knowledge usually required for annotating medical images is much more demanding and difficult to obtain, various medical institutions have very limited strong pixel-level annotated images and most available images are unlabeled or weakly-annotated [3, 20, 21]. Therefore, a realistic clinical mechanism which utilizes every available supervision for cross-institutional collaboration without data sharing is highly desirable.

Thanks to the timely emergence of Federated Learning (FL), which aims to enable multiple clients to jointly train a machine learning model without sharing data, the problem of data privacy being breached can be alleviated [11]. FL has gained significant attention in the medical imaging community [12, 17], due to the obvious reason that medical images often contain some personal information. During the training process of a standard FL model, each local client
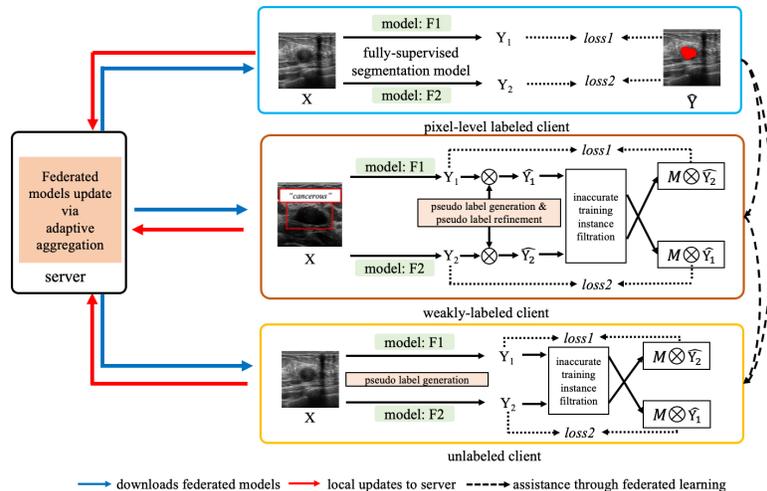
**Figure 2:** An illustration of our proposed Mixed Supervised Federated Learning (FedMix) framework. The local client update utilizes every available supervision for training. Based on which, an adaptive weight aggregation procedure is used for the global federated model update. Compared to the existing methods, FedMix not only breaks through the constraint of a single level of image supervision, but also can dynamically adjust the aggregation weight of each local client, achieving a rich yet discriminative feature representation.

first downloads the federated model from a server and updates the model locally. Then, the locally-trained model parameters of each client are sent back to the server. Finally, all clients' model parameters are aggregated to update the global federated model. Most of the existing FL frameworks [13,18] require that the data used for training by each local client needs to follow the same level of labels, *e.g.*, pixel-level labels (as shown in Fig. 1 (d)) for an image semantic segmentation model, which limits the model learning ability. Although, some semi-supervised federated learning methods [31,33] attempt to utilize the unlabeled data in addition to pixel-level labeled images in training, they do not make any use of the weakly-labeled images (*e.g.*, image-level class labels in Fig. 1 (b) and bounding box labels in Fig. 1 (c)), which are invaluable.

Clients participating in FL may have different labeling budgets. Therefore, there may be a wide range of inter-client variations in label availability. Weak labels are easier to acquire and thus more broadly available compared to pixel-level ones. In practice, there is a wide range of weak labels with varying strengths and acquisition costs. While an image-level label indicating whether a breast ultrasound image is cancerous or not is easier to acquire compared to a bounding box label pointing out the specific location of the cancerous region, it is also less informative. Therefore, effectively utilizing the information from these weakly-labeled data with varying levels of label strengths as well as unlabeled data, especially for clients without pixel-level labeled data would be highly ben-

eficial for improving the federated model's robustness while preventing training instability.

In this work, as illustrated in Fig. 2, we propose a label-agnostic Mixed Supervised Federated Learning (FedMix) framework, which is a unified FL model making use of data labeled in any form for medical image segmentation. Specifically, in the absence of pixel-level labels, FedMix first effectively utilizes unlabeled images as well as useful information contained in the weakly-labeled images (*i.e.*, image-level class labels and bounding box labels) for producing and selecting high-quality pseudo labels. Through an iterative process, the accuracy of selected pseudo labels which are then used for local training on the client sides improves, leading to better model performance. To further improve the model robustness, FedMix takes into account the variability of local clients' available labels through an adaptive aggregation procedure for updating the global federated model. Compared to the existing methods, FedMix not only breaks through the constraint of a single type of labels, but also can dynamically assign an optimized aggregation weight to each local client. Experimental results on two challenging segmentation tasks demonstrate the superior performance of FedMix on learning from mixed supervisions, which is valuable in the clinical setting. Our contributions are summarized as follows:

- The mixed supervised FL framework targeting multi-source medical image segmentation through an iterative pseudo label generator followed by a label refinement operation, based on the information derived from weakly-labeled data, to target high-quality pseudo labels for training.

- An adaptive weight assignment across clients, where each client can learn an aggregation weight. Adaptive weight assignment is essential to handle inter-client variations in supervision availability.

- Extensive experiments on the challenging breast tumor segmentation and skin lesion segmentation. FedMix outperforms the state-of-the-art methods by a large margin.

The rest of this paper is organized as follows: Existing and related work are summarized and discussed in Section 2. The details of FedMix are introduced in Section 3. In Section 4, we present thorough evaluation of FedMix compared with the existing methods. We provide ablation studies as well as analysis in Section 5, and conclude the paper in Section 6.

## 2   Related Work

### 2.1   Federated Learning

Federated learning (FL) is a distributed learning framework, which is designed to allow different clients, institutions, and edge devices to jointly train a machine learning model without sharing the raw data [11], which plays a big role in protecting data privacy. In the past years, FL has drawn great attention from the

medical image communities [18, 46] and has been validated for multi-site functional magnetic resonance imaging classification [13], health tracking through wearables [52], COVID-19 screening and lesion detection [47], and brain tumor segmentation [12, 17]. In clinical practice, different clients may have great variations in data quality, quantity, and supervision availability. Improper use of these data may lead to significant performance degradation among different clients. To reduce the inter-client variations, FL has been combined with domain adaptation [16, 53, 56], contrastive learning [54] and knowledge distillation [55] to learn a more generalizable federated model. However, existing works do not consider the variation in supervision availability (*i.e.*, different clients have different levels of image labels), which is often observed in clinical practice. In our work, we use all the available image label information including image-level class labels, bounding box labels, and pixel-level labels to train a medical image segmentation model and propose a mixed supervised FL framework.

## 2.2   Semi-supervised Federated Learning

In a standard federated learning setting, not every local client has access to pixel-level supervision for image segmentation to facilitate model learning with weakly-labeled and unlabeled training data. To this end, some semi-supervised federated learning approaches require clients to share supplementary information, *e.g.*, client-specific disease relationship [32], extracted features from raw data [34], metadata of the training data [35], and ensemble predictions from different clients' locally-updated models besides their parameters [33]. Additional information sharing beyond the locally-updated model parameters may leak privacy-sensitive information [45] about clients' data. Yang *et al.* [31] proposed to avoid additional information sharing by first training a fully-supervised federated learning model only on clients with available pixel-level supervision for several training rounds and then using the model to generate pseudo labels for local clients based on the unlabeled data. Those confident pseudo labels are used to supervise the local model updates on unlabeled clients for subsequent rounds. In this work, we design a unified federated learning framework that utilizes various weakly supervised data in addition to fully-supervised and unlabeled data for training while limiting the information sharing between clients to only locally-updated model parameters for privacy preservation.

## 2.3   Medical Image Segmentation

The deep learning-based image recognition technology has been used for various medical image segmentation tasks, *e.g.*, optic disc segmentation [24], lung nodules segmentation [2], skin lesion boundary detection [4], and COVID-19 lesion segmentation [3]. However, training a fully-supervised deep model for image semantic segmentation often requires access to a mass of pixel-level supervisions, which are expensive to acquire [21]. In particular, the problem of the expensive pixel-level supervision is much more obstructive for medical image segmentation [26]. To this end, efforts have been made to explore the use of some easily

5

obtained image supervisions (*e.g.*, scribbles [43], image-level classes [6], bounding boxes [7], points [8], and even unlabeled image [36]) to train a pixel-level image segmentation model. However, most of the existing works are based on only one or two types of image supervisions, which greatly limits the model learning efficiency. In most cases, access to some pixel-level annotated data is required to facilitate model training, which may not always be available for each participating client. In our work, we carefully use image-level class labels, bounding box labels, and pixel-level labels to train local clients and propose an adaptive weight assignment procedure across clients for medical image segmentation.

# 3   Our Approach

In this section, we first introduce the notation and experimental setting of the proposed unified federated learning framework, *i.e.*, Fedmix, in Section 3.1. Then, we provide a framework overview in Section 3.2. Finally, we present implementation details including pseudo label generation, selection, and federated model update of the proposed FedMix in Section 3.3 and Section 3.4.

## 3.1   Preliminaries

### 3.1.1   Experimental Settings

To emulate the real scenario setting, we focus on deep learning from multi-source datasets, where each client's data is collected from different medical sources. We focus on exploring variations in cross-client supervisions and thus limit each client to a single level of labels.

### 3.1.2   Training Notations

In this paper, we denote $\overline{D} = [D_1, ..., D_N]$ as the collection of $N$ clients' training data. Given client $i$, $D_i^L = [X, Y_{gt}]$, $D_i^U = [X]$, $D_i^{img} = [X, Y_{img}]$, and $D_i^{bbox} = [X, Y_{bbox}]$ represent the training data that is pixel-level labeled, unlabeled, image-level class labeled, and bounding box-level labeled, respectively. $X$ and $Y$ represent the sets of the training images and the available labels.

To integrate various levels of image labels, in our work, we modify the bounding box labels and image-level class labels to pixel-level labels. Specifically, the bounding box point representation is converted into pixel-level label where the foreground class falls inside the bounding box and the background class falls outside the bounding box. For image-level class labels, we constrain the pixel-level label to the corresponding image class. Consequently, $Y_{gt}$, $Y_{img}$, and $Y_{bbox}$ has the same dimension, *e.g.*, $Y \in \mathbb{R}^{(C+1) \times H \times W}$, $C$ indicates the total number of foreground classes while $W$ and $H$ indicates the weight and height of the respective image data.

6

---
**Algorithm 1:** Pseudocode of FedMix
---
**input**     : $\overline{D}$
**parameter:** $\beta$, $\lambda$: hyperparameters for adaptive aggregation
            $T$: maximum training rounds
            $\epsilon$: threshold for dynamic sample selection
**output**    : $\theta_{\xi 1}$: parameters of $F_1$
            $\theta_{\xi 2}$: parameters of $F_2$

$\theta_{\xi 1}$, $\theta_{\xi 2} \leftarrow$ initialize()
**for** $t = 1 : T$ **do**
    $\mathcal{L}^t = \{\}$, $\theta_{\xi 1}^t = \{\}$, $\theta_{\xi 2}^t = \{\}$
    **for** $i = 1 : N$ **do**
        $f_1^i, f_2^i \leftarrow$ download($\theta_{\xi 1}$, $\theta_{\xi 2}$)
        $X, Y \leftarrow D_i$
        $Y_1, Y_2 \leftarrow F_1^i(X), F_2^i(X)$
        $M_i \leftarrow$ **sample**($Y_1, Y_2, \epsilon$)
        $\hat{Y}_1, \hat{Y}_2 \leftarrow$ **refine**($Y_1, Y_2, Y$)
        $d_i \leftarrow M_i * D_i$
        $\Delta \theta_{i1}^t, \Delta \theta_{i2}^t, \mathcal{L}_i^t \leftarrow$ update($F_1^i, F_2^i; d_i$)
        $\theta_{\xi 1}^t$.add($\Delta \theta_{i1}^t$), $\theta_{\xi 2}^t$.add($\Delta \theta_{i2}^t$), $\mathcal{L}^t$.add($\mathcal{L}_i^t$)
    **end**
    $\theta_{\xi 1}$, $\theta_{\xi 2} \leftarrow$ **aggregate**($\theta_{\xi 1}^t$, $\theta_{\xi 2}^t$, $\mathcal{L}^t$; $\beta$, $\lambda$)
**end**
return $\theta_{\xi 1}$ and $\theta_{\xi 2}$
---

## 3.2 Overview

As illustrated in Fig. 2, to fully utilize *every level of labels* at various clients, the pseudo-code of FedMix is presented in Algorithm 1 and FedMix has two main components:

1. **Pseudo Label Generation and Selection.** In the mixed supervised setting, clients without access to pixel-level label rely on the pseudo labels for training. To improve the pseudo labels' accuracy, we design a unified refinement process using *every level of labels* and dynamically select high-quality pseudo labels for training.

2. **Adaptive Aggregation for Federated Model Update.** FedMix uses an adaptive aggregation operation where the weight of each client is determined by not only its data quantity but also the quality of its pseudo labels. Our aim is to learn a federated model for tumor segmentation, the local model updates without access to pixel-level labels have to be integrated with care. In this way, the reliable clients will be assigned higher aggregation weights, leading to a better federated model.

### 3.3 Pseudo Label Generation and Selection

#### 3.3.1 Pseudo Label Generation

Based on the cross-pseudo supervisions [36], we train two differently initialized models, $F_1(.)$ and $F_2(.)$ to co-supervise each other with pseudo labels when no pixel-level label is available. The training image $X$ is fed to the two models $F_1$ and $F_2$ to generate pseudo labels $Y_1$ and $Y_2$, respectively. The pseudo labels are then refined, denoted as $\hat{Y}_1$ and $\hat{Y}_2$, and used for training the model of each local client. Details of the corresponding refinement strategies for each type of label are introduced as follows:

1. **Pixel-level labels**: Under this kind of supervision, we do refine the pseudo labels, which can be expressed as $\hat{Y}_1 = \hat{Y}_2 = Y_{gt}$.

2. **Bounding box labels**: Each of the predictions $Y_1 = F_1(X_1)$ and $Y_2 = F_2(X_2)$ is refined according to the corresponding bounding box label, *i.e.*, $\hat{Y}_1 = Y_1 * Y_{bbox}$ and $\hat{Y}_2 = Y_2 * Y_{bbox}$.

3. **Image-level class labels**: We do not apply pseudo label refinement, which can be formulated as $\hat{Y}_1 = Y_1$, and $\hat{Y}_2 = Y_2$.

4. **No labels** (*i.e.*, without supervisions): We do not refine the pseudo labels, which is formulated as $\hat{Y}_1 = Y_1$, and $\hat{Y}_2 = Y_2$.

A specific client $i$ is trained by minimizing:

$$\mathcal{L}_i = \mathcal{L}_{dice}(Y_1, \hat{Y}_2) + \mathcal{L}_{dice}(Y_2, \hat{Y}_1), \tag{1}$$

where $\mathcal{L}_{dice}$ is the Dice loss function.

#### 3.3.2 Dynamic Sample Selection

Despite the effectiveness of the above pseudo label generation and refinement processes, the pseudo labels may be incorrect. Therefore, we propose a dynamic sample selection approach to select high-quality data and pseudo labels. Specifically, given client $i$ and its training data $D_i$, we generate a mask $M_i = \{m_1, ..., m_{|D_i|} | m_i \in [0, 1]\}$ to select reliable training samples according to Eq. 2. We measure the consistency between pseudo labels before refinement, *i.e.*, $Y_1$ and $Y_2$. Higher prediction consistency between $Y_1$ and $Y_2$ indicates a higher likelihood that the pseudo labels are closer to ground truth. The above process is expressed as:

$$m_i = \begin{cases} 1 & \text{if } dice(Y_1, Y_2) >= \epsilon \\ 0 & \text{o.w.,} \end{cases} \tag{2}$$

where $\epsilon \in [0, 1]$ is a threshold which is inversely proportional to the amount of selected training samples. For pixel-level labels, $m_i = 1$ for all training samples

as $\hat{Y}_1 = \hat{Y}_2 = Y_{gt}$. As training progresses, the models are more capable of generating more accurate pseudo labels. Consequently, $\sum_{i=1}^{i=|M_i|} m_i$ progressively increases to $|D_i|$, allowing the model to learn from a growing set of training data. More discussions of dynamic sample selection are provided in Section 5.1.

## 3.4 Federated Model Update

At each training round, every local client $i$ first receives the federated model's parameters $\theta_\xi^t$ from the server at time or iteration $t$. Then, every client updates the model locally with its training data $D_i$. Finally, the gradient update from each local client $\Delta\theta_i^{t+1}$ will be sent to the server to update the federated model's parameters according to Eq. 3.

$$\theta_\xi^{t+1} \leftarrow \theta_\xi^t + \sum_{i=1}^{N} w_i \Delta\theta_i^{t+1}. \tag{3}$$

In FedAvg [11], the aggregation weight of each client, $w_i$, is defined as $|D_i| / \sum_{i=1}^{i=|\overline{D}|} |D_i|$. In the mixed supervised setting, relying only on data quantity for weight assignment is sub-optimal. Thus, supervision availability of each client should also be taken into account during the aggregation process. To this end, we propose to utilize the client-specific training loss to infer the data quality. Each client's training loss not only provides a more objective measurement of its importance during FedMix optimization but also prevents the federated model from relying on the over-fitting clients. The proposed adaptive aggregation function is defined by

$$c_i \leftarrow \frac{|D_i|}{\sum_{i=1}^{i=|\overline{D}|} |D_i|}, d_i \leftarrow \frac{\Delta\mathcal{L}_i^\beta}{\sum_{i=1}^{i=|\overline{D}|} \Delta\mathcal{L}_i^\beta}, \tag{4}$$

and

$$w_i \leftarrow \frac{c_i + \lambda \cdot d_i}{\sum_{i=1}^{i=|\overline{D}|} c_i + \lambda \cdot d_i}, \tag{5}$$

where $\lambda$ and $\beta$ are hyper-parameters to tune, impacting the degree of reliance towards different clients. More discussions of adaptive aggregation can be found in Section 5.2.

# 4 Experiments

## 4.1 Datasets and Evaluation Metrics

**Dataset.** In our work, experiments are carried out on two challenging medical image segmentation tasks:

- **Breast tumor segmentation.** In this task, three public breast ultrasound datasets, namely BUS [37], BUSIS [38], and UDIAT [39], are used

**Table 1:** Statistics of the breast ultrasound dataset

| Site | # Patients | # Images | # Healthy / # Cancerous |
|------|-----------|----------|-------------------------|
| BUS | 600 | 780 | $\dfrac{133}{647}$ |
| BUSIS | 562 | 562 | $\dfrac{0}{562}$ |
| UDIAT | 163 | 163 | $\dfrac{0}{163}$ |

**Table 2:** Statistics of the HAM10K dataset

| Site | Source | # Patients | # Images |
|------|--------|-----------|----------|
| Rosendahl | rosendahl | 1552 | 2259 |
| Vidir | modern | 1695 | 3363 |
| | old | 278 | 439 |
| | molemax | 3954 | 3954 |

for evaluation and each of them is regarded as a separate client. More details of this dataset are introduced in Table 1.

- **Skin tumor segmentation.** HAM10K [40] consists of four different sources. Each source acts as a client in FL. The statistics of HAM10K are presented in Table 2.

Following the standard practice, the training data is randomly and patient-wisely split into 80% for training and 20% for testing. All the breast ultrasound and skin dermoscopy images are resized to 256×256 pixels and then randomly flipped and cropped to 224×224 pixels for training.

**Evaluation metrics.** In this work, Dice coefficient (DC) is used for the evaluation of the two segmentation tasks. Considering the two-model architecture of FedMix, the predictions or outputs of $F_1$ are used for evaluation.

## 4.2   Implementation Details

**Network architectures.** UNet [41] combined with the group norm [42] is selected as the baseline segmentation model.

**Supervision types.** The following types of labels are included in our experiments: 1) pixel-level labels (denoted as $L$), 2) bounding box labels (denoted as $B$), 3) image-level class labels (denoted as $I$), and 4) unlabeled (denoted as $U$), *e.g.*, training with only the raw images.

**Comparison methods.** The following four prevailing frameworks are included for comparison:

- Local learning (LL): Each client trains a deep learning network based on its pixel-level labeled data.

**Table 3:** Quantitative results of local learning (LL) and FedAvg under the fully-supervised setting for breast tumor segmentation.

| Frameworks | C1 | C2 | C3 | Avg. |
|:---:|:---:|:---:|:---:|:---:|
| | L | L | L | |
| LL | 66.96 | 87.37 | 87.23 | 80.52 |
| FedAvg | **77.46** | **92.44** | **87.12** | **85.67** |

**Table 4:** Quantitative results of different learning frameworks under the semi-supervised setting for breast tumor segmentation.

| Frameworks | C1 | C2 | C3 | Avg. |
|:---:|:---:|:---:|:---:|:---:|
| | U | U | L | |
| LL (trained on C3) | 64.72 | 83.40 | 87.23 | 78.45 |
| FedST [31] | 64.83 | 85.66 | 86.38 | 78.95 |
| FedMix | **68.17** | **89.19** | **87.97** | **81.77** |

- Federated Averaging (FedAvg): All clients, owning pixel-level labels, collaboratively train a federated model.

- Semi-supervised federated learning via self-training [31] (FedST): FedST is proposed to utilize both pixel-level labeled and unlabeled data for federated training. FedST is selected as it does not require additional information sharing beyond the locally-updated model parameters.

- Our proposed Federated learning with mixed supervisions (FedMix): FedMix integrates various levels of labels.

The performance of FedAvg under the fully-supervised setting is regarded as an upper bound of federated learning techniques. We evaluate the performance of FedMix under the semi-supervised setting by comparing FedMix with FedST. We also evaluate the performance of FedMix under various settings to show how additional weak labels improve the federated model's performance.

**Training details.** All the networks are trained using the Adam optimizer with an initial learning rate of 1e-3 and a batch size of 16. All methods are implemented within the PyTorch framework and trained on Nvidia GeForce Titan RTX GPUs for 300 rounds. The federated training is performed synchronously and the federated model parameters are updated every training round. We set $\epsilon = 0.9$, $\lambda = 10$, and $\beta = 1.5$ and $\beta = 3$ for adaptive aggregation on breast tumor and skin lesion segmentation respectively.

## 4.3 Results on Breast Tumor Segmentation

**Experiment settings.** Data from BUS, BUSIS, and UDIAT are represented by C1, C2, and C3 respectively. To better demonstrate the value of weak labels, C3, owning the least amount of data, is selected as the client with pixel-level

**Table 5:** Quantitative results of FedMix under various weakly-supervised settings for breast tumor segmentation.

| supervision [C1, C2, C3] | C1 | C2 | C3 | Avg. |
|---|---|---|---|---|
| $[U, U, L]$ | 68.17 | 89.19 | 87.97 | 81.77 |
| $[I, U, L]$ | 68.37 | **89.47** | 88.56 | 82.13 |
| $[B, B, L]$ | **71.26** | 89.36 | **89.41** | **83.34** |

**Table 6:** Quantitative results of federated learning under the fully-supervised setting with various aggregation functions for breast tumor segmentation. AdaptAgg is the proposed aggregation function.

| Aggregation Function | C1 $L$ | C2 $L$ | C3 $L$ | Avg. |
|---|---|---|---|---|
| FedAvg | 77.46 | **92.44** | 87.12 | 85.67 |
| AdaptAgg | **79.02** | 93.08 | **88.27** | **86.79** |

labels. The levels of the labels on C1 and C2 are adjusted accordingly for different cases. As only C1 contains both healthy and cancerous images, it is regarded as the client with image-level labels when needed.

**Quantitative evaluation.** According to Table 3, *e.g.*, in the fully-supervised setting, the LL model of C1 has the lowest DC of 66.96%, indicating the large intra-client variations among its data. C2 and C3 performs better than C1, *i.e.*, 87.37% and 87.23% respectively. With FedAvg, every client benefits from the federation, especially C1 with an increase of 10.50% in DC. Besides, training a federated model with more data from different clients is useful to learn more generalizable features, leading to an average increase of 5.15% in DC.

Quantitative results of FedMix and FedST under the semi-supervised setting are provided in Table. 4. For LL, the results of C1 and C2 are produced using the model trained on C3. Compared to the locally-learned models under the fully-supervised setting in Table 3, there exists slight performance degradation on C1 and C2, *i.e.*, 2.24% and 3.97% decrease in DC respectively, indicating the limitation of the model trained only on C3. By utilizing the unlabeled data on C1 and C2 for training, FedST and FedMix are able to train better federated models compared to LL. The overall improvements of FedST are quite limited with an average increase of 0.50% in DC while the segmentation results on C3 are badly affected. Comparatively, FedMix consistently improves the results of all the three clients, leading to an average increase of 3.32% and 2.82% in DC for LL and FedST respectively.

One interesting observation is that FedMix in semi-supervised learning outperforms LL with full supervisions, demonstrating the effectiveness of FedMix in exploring hidden information in unlabeled data. Quantitative results of FedMix under different settings are presented in Table 5. When C1 owns image-level labels, not only C1 but also C2 and C3 would benefit from the federation, shown
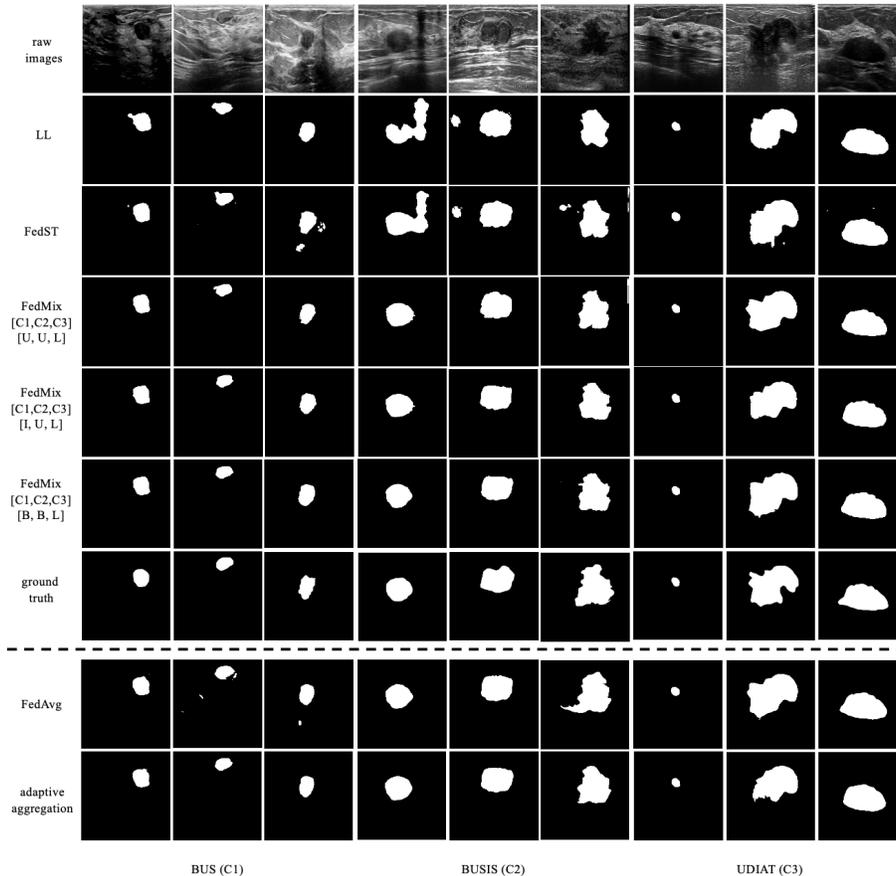
12

**Figure 3:** Exemplar qualitative results of different learning frameworks for breast tumor segmentation. **The upper part** (Rows 1 to 7): the raw images, the segmentation maps produced by local learning (LL), FedST and FedMix under semi-supervision (*i.e.*,[C1, C2, C3] = [$U$, $U$, $L$]), the segmentation maps of FedMix under mixed supervision (*i.e.*,[C1, C2, C3] = [$I$, $U$, $L$] and [C1, C2, C3] = [$B$, $B$, $L$]), and the manual annotations by experts respectively. **The lower part** (Rows 8 to 9): the segmentation maps obtained by federated learning under full pixel-level supervision using FedAvg and the proposed adaptive aggregation function respectively.

by performance improvements across clients, *i.e.*, an average of 0.36% increase in DC. When C1 and C2 have access to bounding box labels, the DC scores of C1 and C3 are further improved, with an average increase of 1.57% and 1.11% compared to FedMix with weaker supervisions. To validate the effectiveness of adaptive aggregation, we compare FedAvg and adaptive aggregation under the fully-supervised setting. The results are presented in Table 6. Putting more emphasis on more reliable clients via adaptive aggregation effectively improves the DC by 1.12%.

**Table 7:** Quantitative results of local learning (LL) and FedAvg under the fully-supervised setting for skin lesion segmentation.

| Frameworks | C1 | C2 | C3 | C4 | Avg. |
| --- | --- | --- | --- | --- | --- |
| | $L$ | $L$ | $L$ | $L$ | |
| LL | 88.98 | 93.21 | 94.33 | 94.93 | 92.86 |
| FedAvg | **90.39** | **93.57** | **95.88** | **95.44** | **93.82** |

**Table 8:** Quantitative results of different learning frameworks under the semi-supervised setting for skin lesion segmentation.

| Frameworks | C1 | C2 | C3 | C4 | Avg. |
| --- | --- | --- | --- | --- | --- |
| | $U$ | $U$ | $L$ | $U$ | |
| LL (trained on C3) | 74.55 | 72.85 | 94.33 | 91.21 | 83.23 |
| FedST [31] | 75.08 | 74.08 | 93.78 | **92.24** | 83.79 |
| FedMix | **80.55** | **81.72** | **94.54** | 90.92 | **86.93** |

**Qualitative evaluation.** According to Fig. 3, LL on C3 produces quite a few false positives on C2, indicating poor generalization capability due to limited training data. Under the semi-supervised setting, though the unlabeled data of C1 and C2 is used for training, the segmentation results of FedST are close to those of LL as learning from incorrect pseudo labels is not helpful and may be detrimental. Comparatively, FedMix can utilize the useful information in unlabeled data and the model generates predictions close to the experts' annotations. The introduction of stronger supervision signals (*i.e.*, from $U$ to $I$ and $B$) to FedMix would further reduce false positives and improve the shape preservation of tumor regions. The utilization of adaptive aggregation in federated learning is beneficial even under the fully-supervised setting. Adaptively aggregated federated model can better capture the boundaries and shapes of the tumor regions and contain fewer false positives compared to the model learned using FedAvg.

## 4.4 Results on Skin Lesion Segmentation

**Experiment setting.** Images from Rosendahl, Vidir-modern, Vidir-old, and Vidir-molemax are represented by C1, C2, C3, and C4 respectively, and C3, owning the least amount of data, is selected as the client with pixel-level labels. The levels of the labels on C1, C2, and C4 are adjusted accordingly for different cases.

**Quantitative results.** From Table 7, under the fully-supervised setting, FedAvg improves the performance of the locally-learned models by an average of 0.96% in DC, indicating that cross-client collaboration is beneficial.

The key for semi-supervised federated learning is to extract and use accurate information from the unlabeled data. Under the semi-supervised setting, where only C3 has access to annotation (*i.e.*, $L$), we present the results in Table 8. The

**Table 9:** Quantitative results of FedMix under various mixed supervised settings for skin lesion segmentation.

| Supervision [C1, C2, C3, C4] | C1 | C2 | C3 | C4 | Avg. |
|---|---|---|---|---|---|
| $[U, U, L, U]$ | 80.55 | 81.72 | 94.54 | 90.92 | 86.93 |
| $[B, B, L, B]$ | **88.80** | **93.11** | **95.82** | **94.41** | **93.04** |

**Table 10:** Quantitative results under the fully-supervised setting with various aggregation functions for skin lesion segmentation. AdaptAgg is the proposed adaptive aggregation.

| Aggregation Function | C1 $L$ | C2 $L$ | C3 $L$ | C4 $L$ | Avg. |
|---|---|---|---|---|---|
| FedAvg | 90.39 | 93.57 | 95.88 | 95.44 | 93.82 |
| AdaptAgg | **90.91** | **93.73** | **96.78** | **95.51** | **94.23** |

locally-learned (LL) model on C3 does not perform well on C1 and C2, observed through the significant performance degradation which indicates severe inter-client variations between {C3, C4} and {C1, C2}. As a result, the pseudo labels on {C1, C2} generated by the model trained on C3 may be inaccurate, utilizing which for training would be harmful. Instead of using all the pseudo labels, FedST makes use of only confident predictions. While the model learned through FedST has an average of 0.56% increase in DC compared to LL, it performs worse on C3, *i.e.*, 0.55% decrease in DC. The performance drop may disincentive C3 to participate in the federation thus hindering the deployment of FedST. With dynamic sample selection and adaptive aggregation, FedMix manages to select high-quality unlabeled data and more accurate pseudo labels for training, thus improving the segmentation performance on C3. Additionally, compared to LL, both C1 and C2 obtain significant performance improvements with an average increase of 6.00% and 8.87% in DC respectively. In general, FedMix achieves better overall performance, resulting in an average increase of 3.14% in DC compared to FedST.

Quantitative results of FedMix under various settings are presented in Table 9. Incorporating bounding box labels for training improves the pseudo labels' accuracy. Consequently, the segmentation performance of FedMix is further improved by 6.11%, approaching the performance of FedAvg under the fully-supervised setting. Bounding box labels are much easier to obtain than pixel-level labels, making FedMix more valuable in clinical scenarios. We further conduct a comparison between FedAvg and adaptive aggregation under the fully-supervised setting, presented in Table 10. The proposed adaptive aggregation function can better utilize the high-quality data and balance the weights among clients, leading to better convergence and segmentation performance.

**Qualitative results.** Qualitative results of skin lesion segmentation are shown in Fig. 4. Consistent with the quantitative results, the segmentation maps on
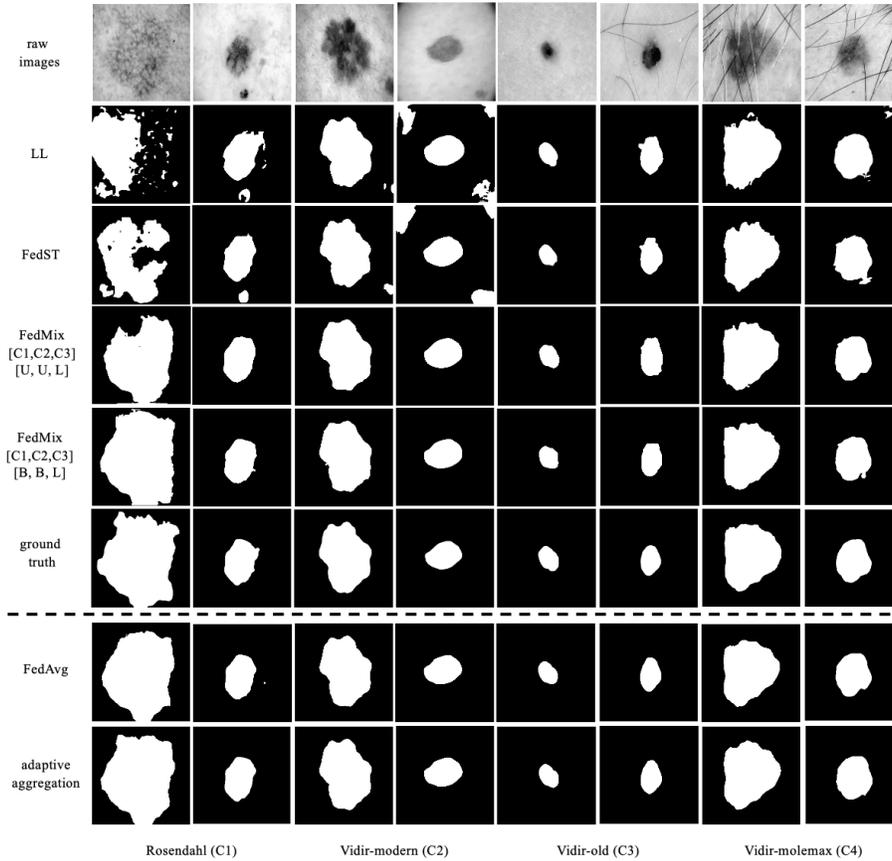
**Figure 4:** Qualitative results of different learning frameworks for skin lesion segmentation. **The upper part** (Rows 1 to 6): the raw images, the segmentation maps produced by local learning (LL), FedST, FedMix under semi-supervision (*i.e.*,[C1, C2, C3, C4] = [U, U, L, U]), FedMix under mixed supervision (*i.e.*,[C1, C2, C3, C4] = [B, B, L, B]), and the expert annotations respectively. **The lower part** (Rows 7 to 8): the segmentation maps obtained by federated learning under the fully-supervised setting with FedAvg and the proposed adaptive aggregation function respectively.

C1 and C2, produced by the locally-learned model on C3, are inaccurate, due to large inter-client variations between {C1, C2} and {C3, C4}. While the segmentation maps produced by FedST are slightly more accurate compared to LL, learning from confident pseudo labels is insufficient to train a generalizable model, shown through the inaccurate segmentation maps produced by FedST on C1 and C2. Under the same supervision setting, FedMix produces more accurate segmentation maps by dynamically selecting the high-quality pseudo labels for training. Given stronger supervisions, *e.g.*, bounding box labels, FedMix improves the segmentation quality, especially on tumor shape preservation.

16

**Table 11:** Quantitative results of FedMix with and without dynamic sample selection for breast tumor and skin lesion segmentation.

| Sample Selection | C1 $U$ | C2 $U$ | C3 $L$ | C4 $U$ | Avg. |
|:---:|:---:|:---:|:---:|:---:|:---:|
| Breast tumor segmentation | | | | | |
| $\times$ | 34.92 | 47.69 | 30.41 | N/A | 37.67 |
| $\checkmark$ | **66.92** | **88.49** | **86.95** | N/A | **80.78** |
| Skin lesion segmentation | | | | | |
| $\times$ | 45.38 | 33.10 | 55.11 | 41.28 | 43.27 |
| $\checkmark$ | **81.30** | **78.10** | **94.43** | **91.11** | **86.24** |

Through the comparison under the fully-supervised setting, we observe that the segmentation maps produced by adaptive aggregation contain fewer false negatives and have better shape consistencies with manual annotations compared to FedAvg.

# 5 Ablation Studies

## 5.1 Effectiveness of Dynamic Sample Selection

We remove the label refinement step in FedMix and utilize FedAvg for comparison. Quantitative results are presented in Table 11. We can observe that without dynamic sample selection, the model may learn from incorrect pseudo labels which is detrimental for convergence. Dynamic sample selection is based on the intuition where the prediction consistencies between the two models given the same input image are positively correlated with the accuracy of the pseudo labels. We perform separate evaluations on the three datasets for breast tumor segmentation, (*i.e.*, BUS (C1), BUSIS (C2), and UDIAT (C3)). For each client, we train two differently initialized models, $F_1$ and $F_2$, locally on 80% of the data for 20 training rounds.

The prediction consistencies between the two models, measured in DC (%), are used to select the evaluation set from the remaining 20% of the data according to the consistency threshold $\epsilon$. With a smaller $\epsilon$, more samples with lower prediction consistencies are included for evaluation. With the increase of $\epsilon$, as only the samples with high prediction consistencies are selected, the overall DC accuracy is higher. The findings in Table 12 validate our assumption and demonstrate the value of dynamic sample selection in filtering inaccurate pseudo labels during training.

## 5.2 Effectiveness of Adaptive Aggregation

We compare adaptive aggregation with FedAvg and present the results in Table 13. For breast tumor segmentation, adaptive aggregation consistently improves performance across clients, with an average of 1.00% increase in DC

**Table 12:** The effect of the threshold $\epsilon$ to the quantitative results (DC %) on each client for breast tumor segmentation.

| $\epsilon$ | BUS (C1) | BUSIS (C2) | UDIAT (C3) |
|---|---|---|---|
| 0.1 | 12.6 | 13.5 | 18.8 |
| 0.2 | 22.0 | 22.9 | 21.3 |
| 0.3 | 25.3 | 39.8 | 26.5 |
| 0.4 | 56.1 | 40.2 | 45.5 |
| 0.5 | 55.9 | 45.2 | 40.7 |
| 0.6 | 66.1 | 60.9 | 52.9 |
| 0.7 | 66.2 | 73.6 | 64.5 |
| 0.8 | 72.3 | 77.0 | 64.9 |
| 0.9 | **86.07** | **89.1** | **79.6** |

**Table 13:** Quantitative results of FedMix with and without adaptive aggregation for breast tumor and skin lesion segmentation.

| Adaptive Aggregation | C1 | C2 | C3 | C4 | Avg. |
|---|---|---|---|---|---|
| | $U$ | $U$ | $L$ | $U$ | |
| Breast tumor segmentation | | | | | |
| $\times$ | 66.92 | 88.49 | 86.95 | N/A | 80.78 |
| $\checkmark$ | **68.17** | **89.19** | **87.97** | N/A | **81.78** |
| Skin lesion segmentation | | | | | |
| $\times$ | **81.30** | 78.10 | 94.43 | **91.11** | 86.24 |
| $\checkmark$ | 80.55 | **81.72** | **94.54** | 90.92 | **86.93** |

compared to FedAvg. For skin lesion segmentation, due to the inter-client variations between {C1, C2} and {C3, C4}, adaptive aggregation focuses more on minimizing the training losses on C1 and C2. As a result, the average DC increase of {C1, C2} is 1.44% while the corresponding increase on C4 is limited to 0.19%. Overall, adaptive aggregation outperforms FedAvg. Till now, aggregation weight optimization in federated learning is still an open problem and should be further explored in the future.

## 6 Conclusion

FedMix is the first federated learning framework that makes effective use of different levels of labels on each client for medical image segmentation. In FedMix, we first generate pseudo labels from clients and use supervision-specific refinement strategies to improve the accuracy and quality of pseudo labels. Then the high-quality data of each client is selected through dynamic sample selection for local model updates. To better update the federated model, FedMix utilizes an adaptive aggregation function to adjust the weights of clients according to both data quantity and data quality. Experimental results on two segmentation

tasks demonstrate the effectiveness of FedMix on learning from various supervisions, which is valuable to reduce the annotation burden of medical experts. In the semi-supervised federated setting, FedMix outperforms the state-of-the-art approach FedST. Compared to FedAvg, the proposed adaptive aggregation function achieves consistent performance improvements on the two tasks under the fully-supervised setting. We believe the methods proposed in FedMix are widely-applicable in FL for medical image analysis beyond mixed supervisions.

# References

[1] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proc. CVPR*, 2015.

[2] D. Jin, Z. Xu, Y. Tang, A. P. Harrison, and D. J. Mollura, "CT-realistic lung nodule simulation from 3D conditionoal generative adversarial networks for robust lung segmentation," in *Proc. MICCAI*, 2018, pp. 732-740.

[3] Y. Li, L. Luo, H. Lin, H. Chen, and P-A. Heng, "Dual-consistency semi-supervised learning with uncertainty quantification for COVID-19 lesion segmentation from CT images," in *Proc. MICCAI*, 2021, pp. 199-209.

[4] N. C. Codella, *et al.*, "Skin lesion analysis toward melanoma detection: A challenge at the 2017 international symposium on biomedical imaging," in *Proc. ISBI*, 2018.

[5] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. CVPR*, 2016, pp. 770-778.

[6] J. Ahn and S. Kwak, "Learning pixel-level semantic affinity with image-level supervision for weakly supervised semantic segmentation," in *Proc. CVPR*, 2018, pp. 4981-4990.

[7] J. Dai, K. He, and J. Sun, "Boxsup: Exploiting bounding boxes to supervise convolutional networks for semantic segmentation," in *Proc. ICCV*, 2015, pp. 1635-1643.

[8] A. Bearman, O. Russakovsky, V. ferrari, and L. Fei-fei, "What's the point: Semantic segmentation with point supervision," in *Proc. ECCV*, 2016, pp. 549-565.

[9] D. Lin, J. Dai, J. Jia, K. He, and J. Sun, "Scribblesup: Scribble-supervised convolutional networks for semantic segmentation," in *Proc. CVPR*, 2016, pp. 3159-3167.

[10] M. Cordts *et al.*, "The cityscapes dataset for semantic urban scene understanding," in *Proc. ICCV*, 2015, pp. 3213-3223.

[11] B. Mcmahan, E. Moore, D. Ramage, S. Hampson, and B. A. y. Arcas, "Communication-efficient learning of deep networks from decentralized data," 2016, *arXiv:1602.05629*.

[12] W. Li *et al.*, "Privacy-preserving federated brain tumour segmentation," in *Proc. MICCAI BrainLes. Workshop*, 2019, pp.92-104.

[13] X. Li *et al.*, "Multi-site fMRI analysis using privacy-preserving federated learning and domain adaptation: ABIDE results," *Med. Image Anal.*, vol. 65, no. 101765, 2020.

[14] Y. Wang, J. Zhang, M. Kan, S. Shan, and X. Chen, "Self-supervised equivariant attention mechanism for weakly supervised semantic segmentation," in *Proc. CVPR*, 2020, pp. 12275-12284.

[15] J. Lee, J. Yi, C. Shin, and S. Yoon, "BBAM: Bounding box attribution map for weakly supervised semantic and instance segmentation," in *Proc. CVPR*, 2021, pp. 2643-2652.

[16] Z. Yan, J. Wicaksana, Z. Wang, X. Yang, and K. T. Cheng, "Variation-aware federated learning with multi-source decentralized medical data," *IEEE J. Biomed. Health Informat.*, vol. 25, no. 7, pp. 2615-2628, 2021.

[17] M. J. Sheller, G. A. Reina, B. Edwards, J. Martin, and S. Bakas, "Multi-institutional deep learning model without sharing patient data: A feasibility study on brain tumor segmentation," in *Proc. MICCAI BrainLes. Workshop.*, 2018, pp. 92-104.

[18] Q. Dou *et al.*, "Federated deep learning for detecting COVID-19 lung abnormalities in CT: A privacy-preserving multinational validation study," *NPJ Digit. Med.*, vol. 4, no. 1, pp. 1-11, 2021.

[19] I. Dayan *et al.*, "Federated learning for predicting clinical outcomes in patients with COVID-19," *Nat. Med.*, vol. 27, no. 10, pp. 1735-1743, 2021.

[20] Y. -X. Zhao, Y. -M. Zhang, M. Song, C. -L. Liu, "Multi-view semi-supervised 3D whole brain segmentation with a self-ensemble network," in *Proc. MICCAI*, 2019, pp. 256-265.

[21] S. Reib, C. Seibold, A. Freytag, E. Rodner, and R. Stiefelhagen, "Every annotation counts: Multi-label deep supervision for medical image segmentation," in *Proc. CVPR*, 2021, pp. 9532-9542.

[22] P. Bilic *et al.*, "The liver tumor segmentation benchmark (LITS)," in *Proc. MICCAI*, 2017.

[23] B. H. Menze *et al.*, "The multimodal brain tumor image segmentation benchmark (BRATS)," *IEEE Trans. Med. Imag.*, vol. 34, no. 10, pp. 1993-2024, 2014.

[24] H. Fu *et al.*, "Joint optic disc and cup segmentation based on multi-label deep network and polar transformation," *IEEE Trans. Med. Imag.*, vol. 37, no. 7, pp. 1597-1605, 2018.

[25] J. Liu *et al.*, "Active cell appearance model induced generative adversarial networks for annotation-efficient cell segmentation and identification on adaptive optics retinal images," *IEEE Trans. Med. Imag.*, vol. 40, no. 10, pp. 2820-2831, 2021.

[26] X. Li, L. Yu, H. Chen, L. Xing, and P. -A. Heng, "Transformation-consistent self-ensembling model for semisupervised medical image segmentation," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 32, no. 2, pp. 523-534, 2021.

[27] Y. Ouali, C. Hudelot, and M. Tami, "Semi-supervised semantic segmentation with cross-consistency training," in *Proc. CVPR*, 2020, pp. 12674-12684.

[28] Z. Ke, D. Qiu, K. Li, Q.Yan, and R. W. H. Lau, "Guided collaborative training for pixel-wise semi-supervised learning," in *Proc. ECCV*, 2020, pp. 429-445.

[29] S. Li, C. Zhang, and X. He, "Shape-aware semi-supervised 3D semantic segmentation for medical images," in *Proc. MICCAI*, 2020, pp. 552-561.

[30] D. Berthelot *et al.*, "Mixmatch: A holistic approach to semi-supervised learning," in *Proc. NeurIPS*, 2019.

[31] D. Yang *et al.*, "Federated semi-supervised learning for COVID region segmentation in chest CT using multi-national data from China, Italy, Japan," *Med. Image Anal.*, vol. 70, no. 101992, 2021.

[32] Q. Liu, H. Yang, Q. Dou, and P. -A. Heng, "Federated semi-supervised medical image classification via inter-client relation matching," in *Proc. MICCAI*, 2021, pp. 325-335.

[33] T. Bdair, N. Navab, and S. Albarqouni, "FedPerl: Semi-supervised peer learning for skin lesion classification," in *Proc. MICCAI*, 2021, pp. 336-346.

[34] Y. Wu, D. Zeng, Z. Wang, Y. Shi, and J. Hu, "Federated contrastive learning for volumetric medical image segmentation," in *Proc. MICCAI*, 2021, pp. 367-377.

[35] N. Dong and I. Voiculescu, "Federated contrastive learning for decentralized unlabeled medical images," in *Proc. MICCAI*, 2021, pp. 378-387.

[36] X. Chen, Y. Yuan, G. Zeng, and J. Wang, "Semi-supervised semantic segmentation with cross pseudo supervision," in *Proc. CVPR*, 2021, pp. 2613-2622.

[37] W. Al-Dhabyani, M. Gooma, H. Khaled, and A. Fahmy, "Dataset of breast ultrasound images," *Data Brief*, vol. 28, no. 104863, 2020.

[38] Y. Zhang *et al.*, "BUSIS: A benchmark for breast ultrasound image segmentation," 2021, *arXiv:1801.03182*.

[39] M. H. Yap *et al.*, "Automated breast ultrasound lesions detection using convolutional neural networks," *IEEE J. Biomed. Health Informat.*, vol. 22, no. 4, pp. 1218-1226, 2018.

[40] P. Tschandl *et al.*, "The HAM10000 dataset, a large collection of multi-source dermatoscopic images of common pigmented skin lesions," *Sci. Data*, vol. 5, no. 180161, 2018.

[41] O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional networks for biomedical image segmentation," in *Proc. MICCAI*, 2015, pp. 234-241.

[42] Y. Wu and K. He, "Group normalization," in *Proc. ECCV*, 2018, pp. 3-19.

[43] Y. B. Can *et al.*,"Learning to segment medical images with scribble-supervision alone," in *Proc. MICCAI DLMIA, ML-CDS Workshop.*, 2018, pp.236-244.

[44] M. Izadyyazdanabadi *et al.*, "Weakly-supervised learning based-feature localization in confocal laser endomicroscopy glioma images," in *Proc. MICCAI*, 2018, pp. 300-308.

[45] J. Sun *et al.*, "Soteria: Provable defense against privacy leakage in federated learning from representation perspective," in *Proc. CVPR*, 2021, pp. 9311-9319.

[46] G. A. Kaissis, M. R. Makowski, D. Ruckert, and R. F. Braren, "Secure, privacy-preserving and federated machine learning in medical imaging," *Nat. Mach. Intell.*, vol. 2, no. 6, pp. 305-311, 2020.

[47] M. Flores *et al.*, "Federated learning used for predicting outcomes in SARS-COV-2 patients," *Res Sq.*, 2021.

[48] M. Y. Lu *et al.*, "Federated learning for computational pathology on gigapixel whole slide images," *Med. Image Anal.*, vol. 76, no. 102298, 2022.

[49] I. Feki, S. Ammar, Y. Kessentini, and K. Muhammad, "Federated learning for COVID-19 screening from chest x-ray images," *Appl. Soft Comput.*, vol. 106, no. 107330, 2021.

[50] N. Rieke *et al.*, "The future of digital health with federated learning," *NPJ Digit. Med.*, vol. 3, no. 1, pp. 1-7, 2020.

[51] P. Guo, P. Wang, J. ZHou, S. Jiang, and V. M. Patel, "Multi-institutional collaborations fori mproving deep learning-based magnetic resonance image reconstruction using federated learning," in *Proc. CVPR*, 2021, pp. 2423-2432.

[52] Y. Chen, X. Qin, J. Wang, C. Yu, and W. Gao, "FedHealth: A federated transfer learning framework for wearable healthcare," *IEEE. Intel. Sys.*, vol. 35, no. 4, pp. 83-93, 2020.

[53] C. Ju *et al.*, "Federated transfer learning for EEG signal classification," in *Proc. EMBC*, 2020, pp. 3040-3045.

[54] Q. Li, B. He, and D. Song, "Model contrastive federated learning," in *Proc. CVPR*, 2021, pp. 10713-10722.

[55] D. Li and J. Wang, "FedMD: Heterogeneous federated learning via model distillation," 2019, *arXiv:1910.03581*.

[56] 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 *Proc. CVPR*, 2021, pp. 1013-1023.