








CXR-FL: Deep Learning-Based Chest X-ray Image Analysis Using Federated Learning

Filip Ślęzyk^{1,2} , Przemysław Jabłeczki^{1,2} , Aneta Lisowska¹ ,
Maciej Malawski^{1,2} , and Szymon Płotka^{1,3} 

¹ Sano Centre for Computational Medicine, Krakow, Poland
s.plotka@sanoscience.org

² AGH University of Science and Technology, Krakow, Poland

³ Warsaw University of Technology, Warsaw, Poland

Abstract. Federated learning enables building a shared model from multicentre data while storing the training data locally for privacy. In this paper, we present an evaluation (called CXR-FL) of deep learning-based models for chest X-ray image analysis using the federated learning method. We examine the impact of federated learning parameters on the performance of central models. Additionally, we show that classification models perform worse if trained on a region of interest reduced to segmentation of the lung compared to the full image. However, focusing training of the classification model on the lung area may result in improved pathology interpretability during inference. We also find that federated learning helps maintain model generalizability. The pre-trained weights and code are publicly available at (<https://github.com/SanoScience/CXR-FL>).

Keywords: Deep learning · Federated learning · Medical imaging

1 Introduction

Federated Learning (FL) is an effective privacy-preserving machine learning technique used to train models across multiple decentralized devices. It enables using a large amount of labeled data in a secure and privacy-preserving process [12] to improve the generalizability of the model [2]. Recent work on the application of federated learning in medical imaging shows promising results in dermoscopic diagnosis [3], volumetric segmentation [4] and chest X-ray image analysis [5]. In this paper, we evaluate the application of deep learning-based models to medical image analysis using the FL method. To gain insight into the impact of FL-related parameters on the global model, we conduct experiments with a variable number of clients and local training epochs. We explore utilisation of cascading approach, where medical image segmentation is performed prior to classification, for increased pathology classification interpretability. We compare our results with [1] in terms of explainable AI (XAI) and classification performance. We find faster convergence of the learning process for a greater fraction

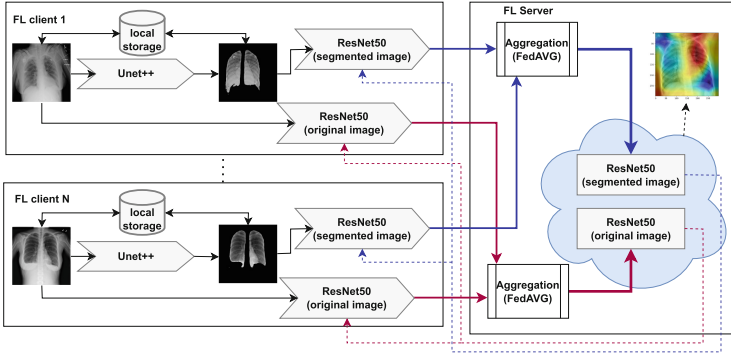


Fig. 1. Methodology: combining segmentation and classification in FL setting

of selected clients and a greater number of local epochs in the segmentation task. We show that federated learning improves the generalizability of the model and helps avoid overfitting in the classification task. We show that Grad-CAM explanations for classification models trained on segmented images may be more focused on the lung area than those trained on full images.

2 Method: FL for Segmentation and Classification

Our method consists of federated training of segmentation and classification models. First, we train segmentation models in a federated manner. For this purpose, we utilize the UNet++ model (with an EfficientNet-B4 backbone) that is later used to prepare the input for classification models. At the classification stage, we use the best segmentation model in terms of the chosen performance metric, and preprocess CXR images (from the training and testing set) to extract lung regions and reduce the impact of the background noise on the prediction. We subsequently train one model on full images and the second on segmented ones independently, all in a federated fashion. During each round of federated training, clients download the global model and fine-tune it with the use of locally stored data. Once all models are fine-tuned in the given round, the server aggregates weights and the next round begins. After the training phase, both types of models pass through the visual explanation step using GradCAM, as in [1]. We test two architectures: ResNet50 and DenseNet121, both commonly used in medical image data classification [10]. An overview of the proposed method for classification stage is depicted in Fig 1.

3 Experiments and Results

3.1 Datasets

Chest X-Ray Dataset: To train the UNet++ model in a federated manner, we use this data set, which is a union of two other data sets known as Chest X-Ray Images (Pneumonia) [8,9]. The dataset consists of 6380 CXR images.

RSNA 2018 Dataset : To evaluate our method, we use an open-source RSNA Pneumonia Detection Challenge 2018 chest X-ray data [7]. In total, the dataset consists of 26684 CXR images in the DICOM format. There are 3 classes in the dataset: “Normal” (8525 - train/326 - test), “No Lung Opacity/Not Normal” (11500 - train/321 - test) and “Lung Opacity” (5659 - train/353 - test).

3.2 Implementation Details

We implement our models in Python 3.8 with the PyTorch v.1.10.1 and Flower v.0.17.0 frameworks, based on our previous experience [6]. We train our models on 4 nodes of a cluster with $1 \times$ NVIDIA v100 GPU each.

For the **segmentation task**, we use UNet++ with EfficientNet-B4 backbone pretrained on ImageNet. Adagrad is utilised as an optimizer for clients. We use a batch size of 2 and set learning rate and weight decay to $lr = 1 \times 10^{-3}$, $wd = 0$ respectively. We assess Jaccard score and BCE-Dice loss on a test set on the central server. The data set used to train the segmentation model was split into a training set and a test set with a 9:1 ratio, maintaining IID distribution of samples. Images are rescaled to 1024×1024 px and augmented with random flip and random affine transformations. The central model is evaluated on a server-side test set after each training round. For the **classification task**, we use Adam optimizer with learning rate $lr = 1 \times 10^{-4}$ and weight decay $wd = 1 \times 10^{-5}$, and set batch size to 8. Images are rescaled to 224×224 px and augmented with random flip and random affine transformations. We evaluate accuracy and CE loss on the test sets (segmented/non-segmented) on the central server. In both tasks, the models are pretrained on the ImageNet dataset. Such pretrained models are downloaded by clients during the first round of the process. We use the FedAvg [11] aggregation strategy and split data in the IID manner among FL clients both in segmentation and classification.

3.3 Segmentation Results

In order to find the optimal central segmentation model, we evaluate several configurations of parameters typical for FL such as the number of local epochs performed by each client during every training round and the fraction of clients selected by the server during each round. The process of training each model consists of 15 rounds. The Jaccard score and loss obtained by each model are presented in Fig 2. For each configuration, we check the number of rounds required to achieve a Jaccard score of 0.92 twice. Results are presented in Table 1. We identify that for a fixed number of local epochs, a greater fraction of selected clients results in a smaller number of rounds needed to exceed the score of 0.92, similarly to the trend observed in [11]. The highest score (0.924) is achieved by the model trained with 3 local epochs and 3 selected clients in the 15th round of training. This model is later used to generate masks for classification.

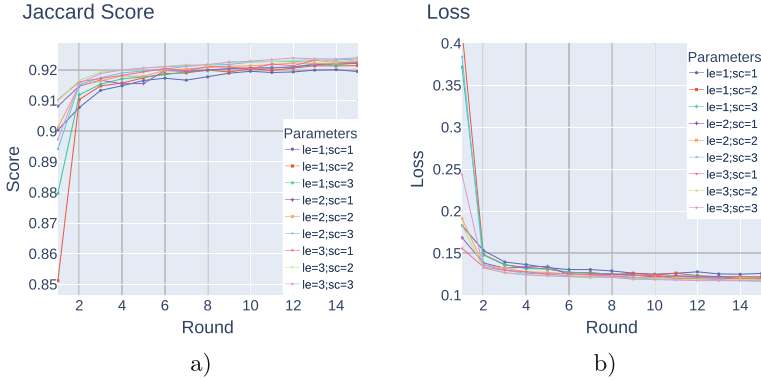


Fig. 2. (a) Jaccard score for the test dataset, achieved by segmentation models, and (b) loss of segmentation models for the test dataset, in successive rounds of training. “sc” - the number of clients selected by the server in each round, “le” - the number of local epochs performed by each client per round.

Table 1. Number of rounds needed by the segmentation model exceeded a Jaccard Score of 0.92 for the serverside test dataset. “sc” - the number of the clients selected by the server in each round, “le” - the number of local epochs performed by each client per round.

Configuration	Experiment 1	Experiment 2
le = 1 & sc = 1	13	14
le = 1 & sc = 2	11	10
le = 1 & sc = 3	9	9
le = 2 & sc = 1	9	9
le = 2 & sc = 2	7	7
le = 2 & sc = 3	6	6
le = 3 & sc = 1	6	6
le = 3 & sc = 2	5	5
le = 3 & sc = 3	5	5

3.4 Classification Results

In the case of the classification task, we evaluate how splitting the same amount of training data between 1, 2 and 3 clients impacts global model quality. Additionally, we assess differences between results obtained with ResNet50 and DenseNet121 architectures on full and segmented images. The accuracy score and loss for 10 rounds of training are presented in Fig. 4. It can be noted that the training process overfits in the case of 1 client and DenseNet121 model, both for segmented and full images, which is represented by a high loss value in the two last rounds for those configurations. The degradation of the global model quality can be also observed for DenseNet121 trained with full images on 2 and

3 clients. The lowest and most stable loss values are obtained for the ResNet50 model trained with 2 and 3 clients for full images and 1 to 3 clients for segmented images. Table 2 presents maximum accuracy and minimum loss values for each configuration of model architecture and dataset type.

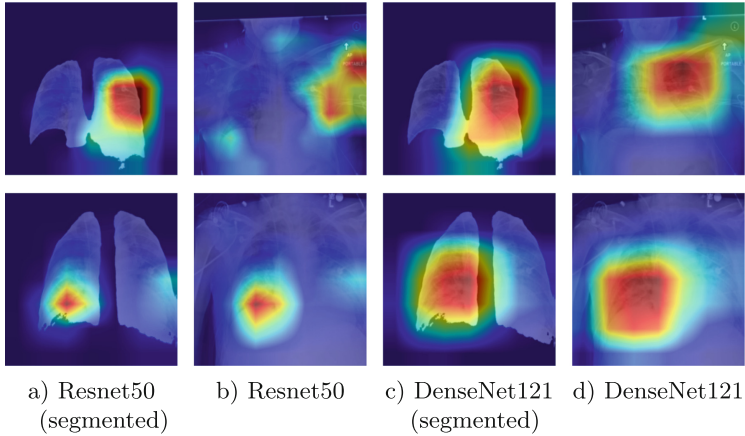


Fig. 3. Grad-CAM visualisations of Lung Opacity samples. In some instances, segmentation resulted in activations focused more on the lung area (upper sample). However, for a majority of cases, visualisation was comparable for segmented and full images (lower sample).

The best accuracy, 0.757, is achieved for ResNet50 model trained on two clients. The worst-performing model is DenseNet121 trained on full images on a single client. In general, the evaluation shows that training on a single client results in overall worse accuracy compared to training with 2 and 3 clients, which is reflected in Fig. 2. This leads to the conclusion that in this case, splitting the data among distinct clients and training the model in the FL manner helps maintain generalizability and avoid overfitting. We observe that models trained on segmented images perform consistently worse than models trained on full images, as is the case for [1]. There is one exception: for the DenseNet121 model the best accuracy is achieved for segmented images (0.742).

To understand qualitative differences in the classification of segmented and full images, we perform Grad-CAM visualisation for ResNet50 and DenseNet121 models. We identify samples that show that the use of segmented images leads to activations more focused on the lung area (as presented in the upper sample in Fig. 3), which is beneficial for model interpretability. However, it can be observed that samples in which the activations are already focused on regions with pathological lung changes, for both full and segmented images, are prevalent. We believe that the small difference in the quality of the models trained on full and segmented images can be explained by the common presence of that similarity.

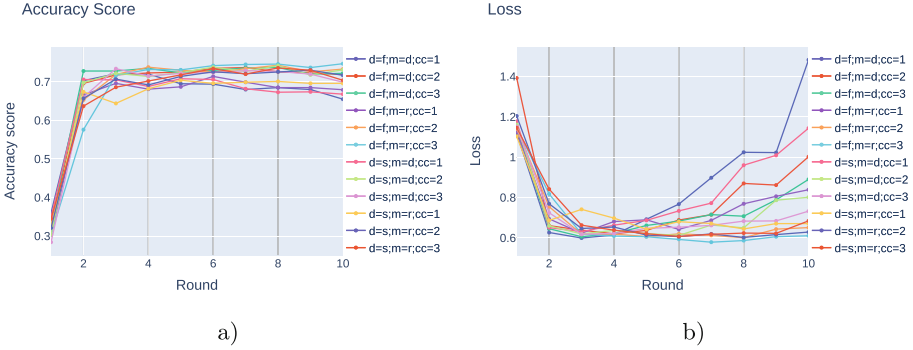


Fig. 4. (a) Accuracy score achieved for the test dataset by classification models, and (b) loss of classification models on test dataset, in successive rounds of training. “d” defines dataset type (f - full/s - segmented), “m” defines model (d - DenseNet121, r - ResNet50), “cc” indicates number of clients participating in training.

Table 2. Maximum accuracy and minimum loss values obtained for each classification model on the test set. “d” defines dataset type (f - full/s - segmented), “m” defines model (d - DenseNet121, r - ResNet50), “cc” indicates number of clients participating in training. Values listed in boldface correspond to extremes in each *model/dataset kind* subset.

Configuration	Max. accuracy	Min. loss
d = f & m = d & cc = 1	0.721	0.599
d = f & m = d & cc = 2	0.737	0.620
d = f & m = d & cc = 3	0.737	0.606
d = f & m = r & cc = 1	0.714	0.623
d = f & m = r & cc = 2	0.757	0.601
d = f & m = r & cc = 3	0.747	0.579
d = s & m = d & cc = 1	0.708	0.631
d = s & m = d & cc = 2	0.742	0.612
d = s & m = d & cc = 3	0.734	0.618
d = s & m = r & cc = 1	0.704	0.643
d = s & m = r & cc = 2	0.730	0.602
d = s & m = r & cc = 3	0.736	0.607

4 Conclusions

In this paper, we evaluated deep learning-based models in the context of CXR image analysis. We conducted experiments in a FL environment to understand the impact of FL-related parameters on the global model performance in segmentation and classification tasks. We also prepared Grad-CAM visualisations for classification models. We found that in the segmentation task, when the number

of local epochs is fixed, the model reaches the desired quality faster with a greater fraction of selected clients. In addition, setting a greater number of local epochs for each client also leads to the same behaviour, which may contribute to lower network traffic in FL processes. Moreover, we conclude that splitting the same dataset among distinct FL clients may lead to improvements in classification for the tested models. We observed a higher accuracy score for full images compared to segmented images in the classification task. However, models trained on segmented images may be characterized by improved interpretability.

Acknowledgements. This publication is partly supported by the EU H2020 grant Sano (No. 857533) and the IRAP Plus programme of the Foundation for Polish Science. This research was supported in part by the PL-Grid Infrastructure. We would like to thank Piotr Nowakowski for his assistance with proofreading the manuscript.

References

1. Teixeira, L.O., et al.: Impact of lung segmentation on the diagnosis and explanation of COVID-19 in chest X-ray images. *Sensors* **21**, 7116 (2021). <https://doi.org/10.3390/s21217116>
2. Sheller, M.J., et al.: Federated learning in medicine: facilitating multi-institutional collaborations without sharing patient data. *Sci. Rep.* **10**, 12598 (2020). <https://doi.org/10.1038/s41598-020-69250-1>
3. Chen, Z., Zhu, M., Yang, C., Yuan, Y.: Personalized retrogress-resilient framework for real-world medical federated learning. In: de Bruijne, M., et al. (eds.) *MICCAI 2021. LNCS*, vol. 12903, pp. 347–356. Springer, Cham (2021). https://doi.org/10.1007/978-3-030-87199-4_33
4. Wu, Y., Zeng, D., Wang, Z., Shi, Y., Hu, J.: Federated contrastive learning for volumetric medical image segmentation. In: de Bruijne, M., et al. (eds.) *MICCAI 2021. LNCS*, vol. 12903, pp. 367–377. Springer, Cham (2021). https://doi.org/10.1007/978-3-030-87199-4_35
5. Dong, N., Voiculescu, I.: Federated contrastive learning for decentralized unlabeled medical images. In: de Bruijne, M., et al. (eds.) *MICCAI 2021. LNCS*, vol. 12903, pp. 378–387. Springer, Cham (2021). https://doi.org/10.1007/978-3-030-87199-4_36
6. Jabłęcki, P., Ślęzyk, F., Malawski, M.: Federated learning in the cloud for analysis of medical images - experience with open source frameworks. In: Oyarzun Laura, C., et al. (eds.) *DCL/PPML/LL-COVID19/CLIP -2021. LNCS*, vol. 12969, pp. 111–119. Springer, Cham (2021). https://doi.org/10.1007/978-3-030-90874-4_11
7. Shih, G., et al.: Augmenting the national institutes of health chest radiograph dataset with expert annotations of possible pneumonia. *Radiol. Artif. Intell.* **1**, e180041 (2019). <https://doi.org/10.1148/ryai.2019180041>
8. Kermany, D., et al.: Labeled optical coherence tomography (OCT) and chest x-ray images for classification. *Mendeley Data* **2** (2018). <https://doi.org/10.17632/rscbjbr9sj.2>
9. Cohen, J.P., et al.: COVID-19 image data collection: prospective predictions are the future. [arXiv:2006.11988](https://arxiv.org/abs/2006.11988) [cs, eess, q-bio] (2020)
10. Tang, Y.-X., et al.: Automated abnormality classification of chest radiographs using deep convolutional neural networks. *npj Digit. Med.* **3**, 70 (2020). <https://doi.org/10.1038/s41746-020-0273-z>

11. McMahan, H.B., et al.: Communication-efficient learning of deep networks from decentralized data (2016). <https://doi.org/10.48550/ARXIV.1602.05629>
12. Kaissis, G.A., et al.: Secure, privacy-preserving and federated machine learning in medical imaging. *Nat. Mach. Intell.* **2**(6), 305–311 (2020)