

SegViz: A Federated Learning Framework for Medical Image Segmentation from Distributed Datasets with Different and Incomplete Annotations

Adway U. Kanhere^{1,2}

AKANHERE@SOM.UMARYLAND.EDU

Pranav Kulkarni¹

PKULKARNI@SOM.UMARYLAND.EDU

Paul H. Yi¹

PYI@SOM.UMARYLAND.EDU

Vishwa S. Parekh¹

VPAREKH@SOM.UMARYLAND.EDU

¹ *University of Maryland Medical Intelligent Imaging (UM2i) Center*

University of Maryland School of Medicine

Baltimore, MD 21201

² *Department of Biomedical Engineering, Johns Hopkins University, Baltimore, MD 21218*

Editors: Under Review for MIDL 2023

Abstract

Segmentation is one of the primary tasks in the application of deep learning in medical imaging, owing to its multiple downstream clinical applications. As a result, many large-scale segmentation datasets have been curated and released for the segmentation of different anatomical structures. However, these datasets focus on the segmentation of a subset of anatomical structures in the body, therefore, training a model for each dataset would potentially result in hundreds of models and thus limit their clinical translational utility. Furthermore, many of these datasets share the same field of view but have different subsets of annotations, thus making individual dataset annotations incomplete. To that end, we developed SegViz, a federated learning (FL) framework for aggregating knowledge from distributed medical image segmentation datasets with different and incomplete annotations into a ‘global’ meta-model. We evaluated the SegViz framework for the task of liver and spleen segmentation on CT scans. The experimental setup involved two nodes with CT datasets with similar fields of view but different and incomplete annotations. The first node consisted of the MSD liver segmentation dataset, while the second consisted of the MSD spleen segmentation dataset. The SegViz framework was trained to build a single model capable of segmenting both liver and spleen aggregating knowledge from both these nodes by aggregating the weights after every 10 epochs. The global SegViz model was tested on an external dataset, Beyond the Cranial Vault (BTCV), comprising both liver and spleen annotations using the dice similarity (DS) metric. The baseline individual segmentation models for spleen and liver trained on their respective datasets produced a DS score of 0.834 and 0.878 on the BTCV test set. In comparison, the SegViz model produced comparable mean DS scores of 0.829 and 0.899 for the segmentation of the spleen and liver, respectively on the BTCV dataset. Our results demonstrate SegViz as an essential first step towards training clinically translatable multi-task segmentation models from distributed datasets with disjoint incomplete annotations with excellent performance.

Keywords: Deep learning, federated learning, 3D segmentation, MONAI, UNet.

1. Introduction

Medical image segmentation is one of the primary tasks in automated medical image analysis as it forms the basis for many downstream applications, including diagnosis, prognosis, and treatment response assessment (Chen et al., 2021; Flores et al., 2021; Menze et al., 2014). As a result, many large-scale datasets have been curated and released for the segmentation of different tissue types and anatomical structures (Antonelli et al., 2022; Wasserthal et al., 2022; Sekuboyina et al., 2021). However, each of these datasets has been curated for a specific use case and therefore, focuses on segmenting only a subset of anatomical structures in the body. Consequently, developing and deploying algorithms for each use case would potentially result in hundreds of models, thereby limiting their clinical utility – imagine deploying a different algorithm for every type of cancer, injury, and other diseases. Furthermore, aggregating multiple datasets to a central repository to train a single multi-task segmentation model is challenging owing to patient privacy concerns and space/computational requirements.

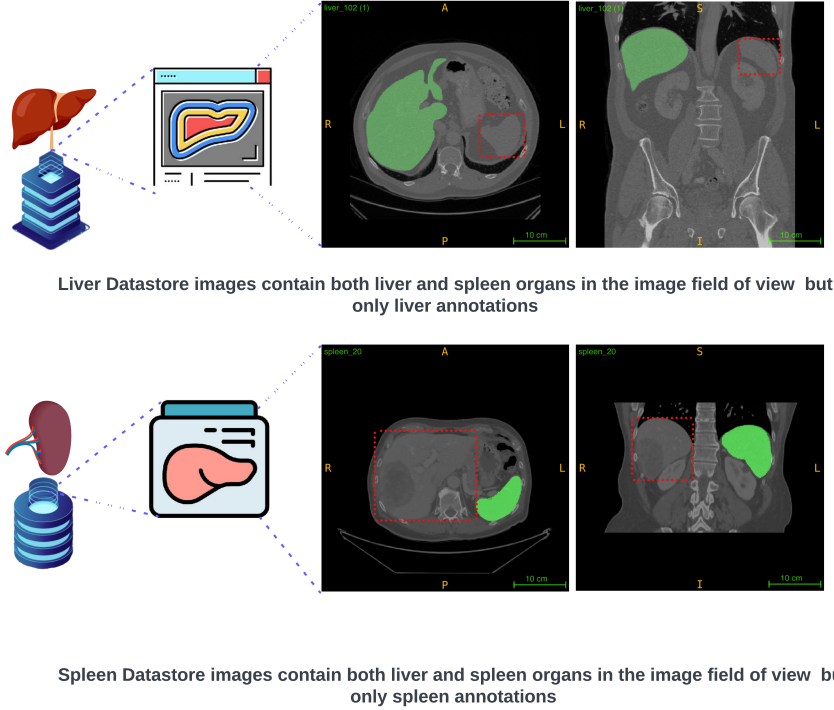


Figure 1: Illustration of an example federated learning setup with nodes containing datasets with a similar field of view but different and incomplete annotations.

These challenges can be addressed by training multi-task segmentation models from distributed datasets using collaborative learning. However, aggregating knowledge from

datasets curated at different imaging centers is challenging as different imaging centers may focus on related but different tasks; suppose one center is training a liver segmentation model while another center is training a spleen segmentation model from CT scans. These two datasets would contain images with a similar field of view but different, incomplete annotations, as illustrated in Figure 1. Such a situation, where one dataset has only a few organs annotated while another dataset contains no overlapping annotations with the first one is very common in medical imaging. Therefore, there is a need to develop methods that can be effective in sharing knowledge in a privacy-preserving, distributed setting, where each dataset can have inherent heterogeneity and contains missing annotations. To that end, we developed SegViz, a federated learning (FL) framework for aggregating knowledge from heterogeneous, distributed medical image segmentation datasets and incomplete annotations into a ‘global’ meta-model. We evaluated the SegViz framework for the task of liver and spleen segmentation on CT scans using distributed nodes containing MSD liver and MSD spleen segmentation datasets. The segmentations from the SegViz framework were compared to training separate baseline models for each individual task.

2. Related Work

Generating manual annotations for medical images is time-consuming, requires high skill, and is an expensive effort, especially for 3D images (Tajbakhsh et al., 2020). One potential solution is to curate datasets with partial annotations, wherein only a subset of structures is annotated for each image or volume. Furthermore, knowledge from similar partially annotated datasets from multiple groups can be aggregated to collaboratively train global models using Federated Learning (Chowdhury et al., 2022). Knowledge aggregation would not only save time but also allow different groups to benefit from each other’s annotations without explicitly sharing them. Consequently, different techniques have been proposed in the literature for aggregating knowledge from heterogeneous datasets with partial, incomplete labels (Parekh et al., 2021; Shen et al., 2021; Boutillon et al., 2022; Shen et al., 2022).

In Boutillon et al. (2022), the authors developed a multi-task multi-domain deep segmentation model for the segmentation of pediatric imaging datasets with excellent performance. However, the proposed technique was developed and evaluated for different anatomical regions in the body with no overlapping field of view or incomplete annotations. Similarly, the cross-domain medical image segmentation technique developed in Parekh et al. (2021) was focused on segmentation of the same anatomical structure and the proposed technique was not developed to tackle incomplete annotations. In contrast, (Shen et al., 2022) attempted to train global segmentation models from partially annotated distributed datasets. The global federated learning framework developed in their work, however, failed to accurately segment different anatomical structures on the external test set. For optimal performance, the authors used an ensemble of multiple local federated learning models, making it computationally expensive and practically challenging.

Therefore, we developed SegViz, a federated learning framework to address the shortcomings of current techniques in efficiently aggregating knowledge from heterogeneous datasets with incomplete annotations. The SegViz framework utilizes the intrinsic similarities between the different imaging datasets to learn a general representation across

multiple tasks. We evaluated the SegViz framework for the task of knowledge aggregation from MSD liver and MSD spleen datasets and tested its performance on an external BTCV dataset.

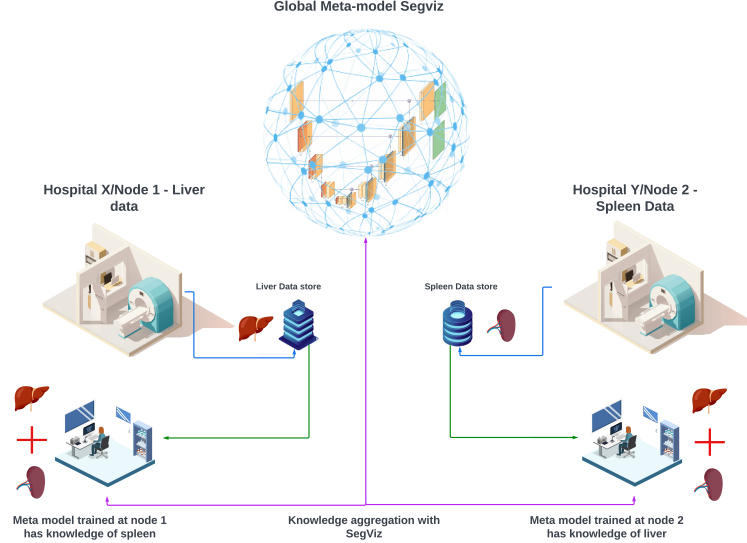


Figure 2: Illustration of the proposed SegViz framework: Client nodes update the global meta-model where knowledge aggregation occurs after every 10 iterations of the local model. The weights of the global model are then shared with the client models allowing both nodes to share knowledge without sharing data.

3. Methods

3.1. SegViz

We developed SegViz as a multi-task federated learning framework to learn a diverse set of tasks from distributed nodes with incomplete annotations, as illustrated in Figure 2. The global SegViz model is initialized at the server with two distinct blocks - a representation block and a task block. The goal of the representation block is to learn a generalized representation of the underlying dataset while the goal of the task block is to learn individual tasks distributed across different nodes. Every client is initialized with a subset of the SegViz model, comprising the representation block and a subset of the task block representing the client’s tasks. During training, the weights of the representation block are always aggregated by the server and redistributed back to the client nodes. On the other hand, the weights of the task block are directly copied from the corresponding client nodes containing the corresponding task, thereby preserving the task-related information for each node in their task block.

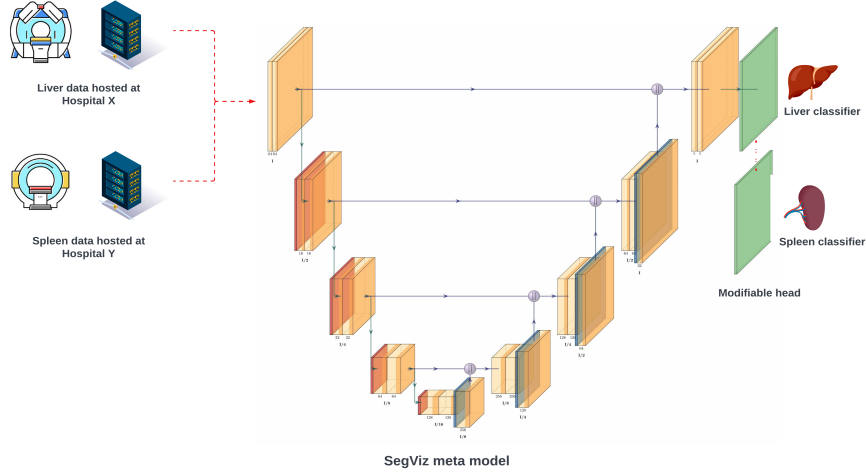


Figure 3: Illustration of the modified 3D-UNet configuration: The representation block refers to all the layers except the last two convolutional layers, while the task block refers to the final two layers including the classifier. The classifier (green) can be fine-tuned for either the spleen or liver task. Figure generated using (Iqbal, 2018)

3.1.1. SEG VIZ MODEL ARCHITECTURE

The backbone of the SegViz model architecture was constructed using a modified version of the multi-head 3D-UNet (Çiçek et al., 2016) configuration for all our experiments. Each U-Net has 5 layers with down/up-sampling at each layer by a factor of 2. Unlike how U-Net implementations typically operate, these down or up-sampling operations happen at the beginning of each block instead of at the end. The U-Net also contains 2 convolutional residual units at the layers and uses Batch Normalization at each layer. The task block comprised a multi-head architecture with each head consisting of two layers, including the final classification layer. The SegViz model was implemented using the MONAI (Cardoso et al., 2022) framework and the pre-processing and training were done using Pytorch. The SegViz model architecture has been illustrated in Figure 3.

4. Experiments and Data

4.1. Clinical Data

The SegViz framework was evaluated using the spleen and liver training datasets from the Medical Segmentation Decathlon (MSD) challenge (Antonelli et al., 2022). The Spleen MSD dataset consisted of 61 3D Computed tomography (CT) volumes with spleen annotations. The Liver MSD dataset consisted of 201 3D CT volumes with liver and liver tumor annotations. For this study, the liver tumor annotations were discarded. We considered all 30 training image volumes from the Beyond the Cranial Vault (BTCV) dataset (Landman

et al., 2015) as an external test set for all our experiments. All the image volumes were resized to $256 \times 256 \times 128$, and the voxel values normalized between 0 and 1.

4.2. Baseline segmentation

We trained two baseline U-Net models, one each on the liver and spleen segmentation datasets. The training and internal validation splits were considered from the overall training data in an 80:20 split. All 30 training samples from the BTCV dataset were used as the external test set. Random foreground patches of size $128 \times 128 \times 32$ were extracted from each volume such that the center voxel of each patch belonged to either the foreground or background class. The batch size was set to 2 and the learning rate was initially set to $1e-4$ with the Adam optimizer and CosineAnnealingLR (Loshchilov and Hutter, 2016) as the scheduler. The Dice Loss was used as the loss function. The average Dice Score was chosen as the final evaluation metric. Each model was trained for 500 epochs.

4.3. SegViz segmentation

The training and validation data split, as well as the test set used for evaluating SegViz segmentations, were the same as the baseline segmentation experiment. For the Segviz setup, the batch size at the nodes was set to 2 and the learning rate was initially set to $1e-4$ with the Adam optimizer with CosineAnnealingLR (Loshchilov and Hutter, 2016) as the scheduler. Similar to the baseline experiment, the Dice Loss was used as the loss function and the average Dice Score was chosen as the final evaluation metric. The SegViz framework was trained to aggregate knowledge from the liver and spleen datasets using the following procedure:

1. Initialize the individual local models on the local nodes and the global SegViz model at the central server
2. Each node trains its respective model for 10 epochs locally.
3. The weights of the models at each node are transmitted to the global model
4. The weights of the two models are aggregated using federated averaging (McMahan and more, 2017) for all but the last 2 layers (task block) by the global model.
5. The global model then shares the updated weights with the local nodes.
6. Steps 2 to 5 are repeated for 1000 iterations.

5. Results

Figure 4 illustrates the performance of the SegViz model compared to the baseline segmentations and ground truth for an example set of three patients. The overlap between the ground truth and predictions have color-coded as yellow. The red and green colocoding represent the over-segmented and under-segmented regions by the segmentation algorithms, respectively. As shown in Figure 4, the SegViz model segmented both the liver and spleen with an excellent performance compared to individual baseline models for each anatomical

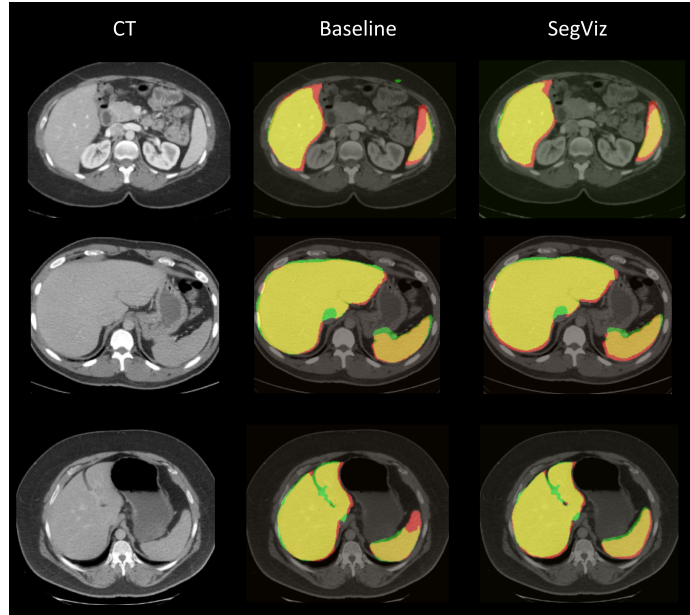


Figure 4: A comparison of the ground truth segmentation masks with the masks generated by the baseline and SegViz models. Each image contains both liver and spleen annotations but either one was present during training. The Red regions indicate the ground truth annotations, the Green regions indicate the output of the trained model, and the Yellow regions indicate the overlapping regions between the ground truth and trained models.

Table 1: Results $N=30$ on the BTCV test set: SegViz refers to the proposed federated-learning-based training of the model and baseline refers to centralized training of the model.

Method	Average Dice Score
SegViz spleen	0.829
Baseline spleen	0.834
SegViz liver	0.899
Baseline liver	0.878

region. The baseline models trained only to segment the liver and spleen individually had an average dice score of 0.834 and 0.878 for the spleen and liver respectively on the test BTCV image set. In comparison, the SegViz FL model demonstrated a similar performance with an average dice score of 0.829 and 0.899 for the segmentation of the spleen and liver, respectively on the same BTCV test set. Figure 5 shows the boxplots of the average dice scores for each experiment. The dice scores for all the models have been tabulated in Table 1.

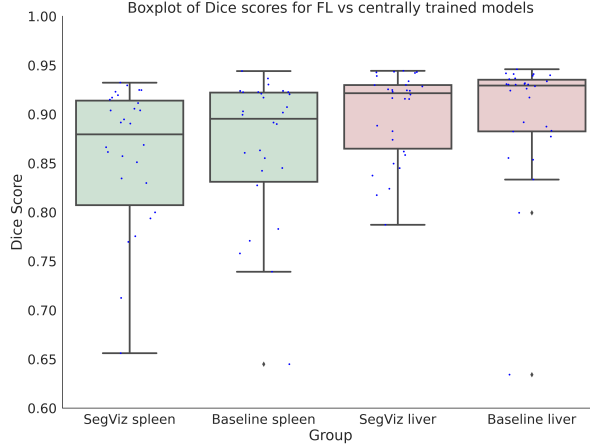


Figure 5: Boxplot of dice scores between the baseline and SegViz models for the liver and spleen datasets

6. Discussion

The SegViz framework proposed in this work demonstrated excellent performance in aggregating knowledge from heterogeneous datasets with different, incomplete labels. We evaluated our pipeline in a scenario where liver and spleen data was distributed across two nodes with incomplete annotations. Our approach successfully aggregated knowledge from both nodes with little to no drop in the performance of the global meta-model in terms of the average dice score. The comparable performance between the SegViz segmentations and multiple baseline model segmentations illustrates a preliminary example of constructing multi-task segmentation models with clinical applicability from dispersed datasets with disjoint partial annotations.

SegViz can be extended to multiple nodes, each with distinct heterogeneous data. Image segmentation from heterogeneous datasets with incomplete annotations has many potential benefits. For example, SegViz can potentially reduce labeling time by $1/\eta$ where η is the number of distinct labels in the distributed data sets by allowing the transfer of knowledge between each client. This would not only save time but also allow different research groups potentially benefit from each others' annotations without explicitly sharing them.

In recent years, different techniques have emerged in the literature for collaboratively training image segmentation models from heterogeneous datasets (Parekh et al., 2021; Shen et al., 2021; Boutillon et al., 2022; Shen et al., 2022). However, most of these techniques have either focused on aggregating knowledge from different domains with non-overlapping field of view or focused on similar tasks (Parekh et al., 2021; Shen et al., 2021; Boutillon et al., 2022). More recently, (Shen et al., 2022) attempted to aggregate knowledge from incompletely labeled datasets. However, the proposed global federated learning model failed to accurately segment different structures which were overcome by a more computationally expensive ensemble of multiple local federated learned models. In contrast, SegViz was able

to train global federated learning to accurately segment different structures with comparable performance to multiple independent segmentation models for each structure.

Our study has certain limitations. This is a preliminary study with experiments performed on small-scale publicly available datasets to establish a proof of concept. In the future we plan to extend this work to a framework with multiple nodes, each containing datasets with different, overlapping, or incomplete annotations.

References

- Michela Antonelli, Annika Reinke, Spyridon Bakas, Keyvan Farahani, Annette Kopp-Schneider, Bennett A Landman, Geert Litjens, Bjoern Menze, Olaf Ronneberger, Ronald M Summers, et al. The medical segmentation decathlon. *Nature communications*, 13(1):1–13, 2022.
- Arnaud Boutillon, Pierre-Henri Conze, Christelle Pons, Valérie Burdin, and Bhushan Borotikar. Generalizable multi-task, multi-domain deep segmentation of sparse pediatric imaging datasets via multi-scale contrastive regularization and multi-joint anatomical priors. *Medical Image Analysis*, 81:102556, 2022.
- M Jorge Cardoso, Wenqi Li, Richard Brown, Nic Ma, Eric Kerfoot, Yiheng Wang, Benjamin Murrey, Andriy Myronenko, Can Zhao, Dong Yang, et al. Monai: An open-source framework for deep learning in healthcare. *arXiv preprint arXiv:2211.02701*, 2022.
- Xuming Chen, Shanlin Sun, Narisu Bai, Kun Han, Qianqian Liu, Shengyu Yao, Hao Tang, Chupeng Zhang, Zhipeng Lu, Qian Huang, et al. A deep learning-based auto-segmentation system for organs-at-risk on whole-body computed tomography images for radiation therapy. *Radiotherapy and Oncology*, 160:175–184, 2021.
- Alexander Chowdhury, Hasan Kassem, Nicolas Padoy, Renato Umeton, and Alexandros Karargyris. A review of medical federated learning: Applications in oncology and cancer research. In *International MICCAI Brainlesion Workshop*, pages 3–24. Springer, 2022.
- Özgün Çiçek, Ahmed Abdulkadir, Soeren S Lienkamp, Thomas Brox, and Olaf Ronneberger. 3d u-net: learning dense volumetric segmentation from sparse annotation. In *International conference on medical image computing and computer-assisted intervention*, pages 424–432. Springer, 2016.
- Mona Flores, Ittai Dayan, Holger Roth, Aoxiao Zhong, Ahmed Harouni, Amilcare Gentili, Anas Abidin, Andrew Liu, Anthony Costa, Bradford Wood, et al. Federated learning used for predicting outcomes in sars-cov-2 patients. *Research Square*, 2021.
- Haris Iqbal. Harisiqbal88/plotneuralnet v1.0.0, December 2018. URL <https://doi.org/10.5281/zenodo.2526396>.
- Bennett Landman, Zhoubing Xu, J Igelsias, Martin Styner, T Langerak, and Arno Klein. Miccai multi-atlas labeling beyond the cranial vault—workshop and challenge. In *Proc. MICCAI Multi-Atlas Labeling Beyond Cranial Vault—Workshop Challenge*, volume 5, page 12, 2015.

- Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016.
- Brendan McMahan and more. Communication-Efficient Learning of Deep Networks from Decentralized Data. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, volume 54 of *Proceedings of Machine Learning Research*, pages 1273–1282. PMLR, 20–22 Apr 2017.
- Bjoern H Menze, Andras Jakab, Stefan Bauer, Jayashree Kalpathy-Cramer, Keyvan Farahani, Justin Kirby, Yuliya Burren, Nicole Porz, Johannes Slotboom, Roland Wiest, et al. The multimodal brain tumor image segmentation benchmark (brats). *IEEE transactions on medical imaging*, 34(10):1993–2024, 2014.
- Vishwa S Parekh, Shuhao Lai, Vladimir Braverman, Jeff Leal, Steven Rowe, Jay J Pillai, and Michael A Jacobs. Cross-domain federated learning in medical imaging. *arXiv preprint arXiv:2112.10001*, 2021.
- Anjany Sekuboyina, Malek E Hussein, Amirhossein Bayat, Maximilian Löffler, Hans Liebl, Hongwei Li, Giles Tetteh, Jan Kukačka, Christian Payer, Darko Štern, et al. Verse: A vertebrae labelling and segmentation benchmark for multi-detector ct images. *Medical image analysis*, 73:102166, 2021.
- Chen Shen, Pochuan Wang, Holger R Roth, Dong Yang, Daguang Xu, Masahiro Oda, Weichung Wang, Chiou-Shann Fuh, Po-Ting Chen, Kao-Lang Liu, et al. Multi-task federated learning for heterogeneous pancreas segmentation. In *Clinical Image-Based Procedures, Distributed and Collaborative Learning, Artificial Intelligence for Combating COVID-19 and Secure and Privacy-Preserving Machine Learning*, pages 101–110. Springer, 2021.
- Chen Shen, Pochuan Wang, Dong Yang, Daguang Xu, Masahiro Oda, Po-Ting Chen, Kao-Lang Liu, Wei-Chih Liao, Chiou-Shann Fuh, Kensaku Mori, et al. Joint multi organ and tumor segmentation from partial labels using federated learning. In *International Workshop on Distributed, Collaborative, and Federated Learning, Workshop on Affordable Healthcare and AI for Resource Diverse Global Health*, pages 58–67. Springer, 2022.
- Nima Tajbakhsh, Laura Jeyaseelan, Qian Li, Jeffrey N Chiang, Zhihao Wu, and Xiaowei Ding. Embracing imperfect datasets: A review of deep learning solutions for medical image segmentation. *Medical Image Analysis*, 63:101693, 2020.
- Jakob Wasserthal, Manfred Meyer, Hanns-Christian Breit, Joshy Cyriac, Shan Yang, and Martin Segeroth. Totalsegmentator: robust segmentation of 104 anatomical structures in ct images. *arXiv preprint arXiv:2208.05868*, 2022.