# FedCostWAvg: A new averaging for better Federated Learning

Leon Mächler[6], Ivan Ezhov[1,3], Florian Kofler[1,2,3], Suprosanna Shit[1,3], Johannes C. Paetzold[1,3,5], Timo Loehr[1,3], Claus Zimmer[2], Benedikt Wiestler[2], and Bjoern H. Menze[1,3,4]

[1] Department of Informatics, Technical University Munich, Germany
[2] Department of Diagnostic and Interventional Neuroradiology, School of Medicine, Klinikum rechts der Isar, Technical University of Munich, Germany
[3] TranslaTUM - Central Institute for Translational Cancer Research, Technical University of Munich, Germany
[4] Department of Quantitative Biomedicine, University of Zurich, Switzerland
[5] ITERM Institute Helmholtz Zentrum Muenchen, Neuherberg, Germany
[6] École Normale Supérieure, Paris, France
`leon-philipp.machler@ens.fr`

**Abstract.** We propose a simple new aggregation strategy for federated learning that won the MICCAI Federated Tumor Segmentation Challenge 2021 (FETS), the first ever challenge on Federated Learning in the Machine Learning community. Our method addresses the problem of how to aggregate multiple models that were trained on different data sets. Conceptually, we propose a new way to choose the weights when averaging the different models, thereby extending the current state of the art (FedAvg). Empirical validation demonstrates that our approach reaches a notable improvement in segmentation performance compared to FedAvg.

**Keywords:** Federated Learning · Brain Tumor Segmentation · Multi-Modal Medical Imaging · MRI · MICCAI Challenges · Machine Learning

## 1 Introduction

### 1.1 Motivation

Preserving data privacy is of paramount importance for confidentiality-critical fields such as the medical domain. Today it is not uncommon that large volumes of private medical records are illegally released to the *dark web*[1]. To prevent such incidents, often large amounts of resources are allocated but cannot guarantee full security. Among many precautions, reducing human (including IT specialists) exposure to the data is highly desirable to reduce the chance of compromising data protection by human failure.
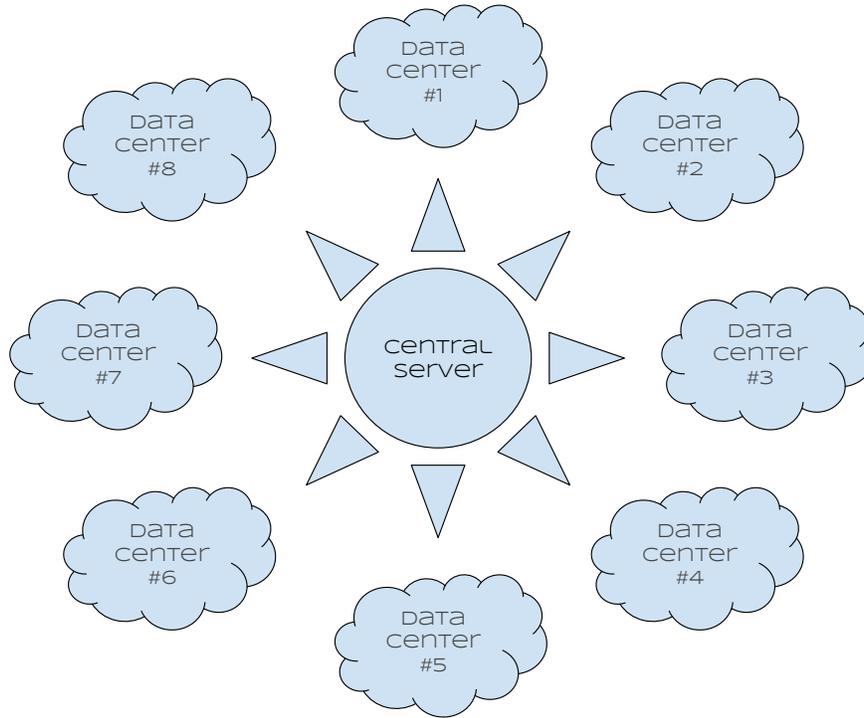
**Fig. 1.** Schematic illustration of the federated learning concept. Within multiple data centers, a model is trained for our task. Next, parameters are sent to the central server, where aggregation of the parameters takes place. An aggregated global configuration of the parameters is broadcasted back to the centers. The procedure repeats until convergence or some other limit is reached.

### 1.2   The typical training scenario

In machine learning, a common scenario today looks like this: One or more institutions (companies, research institutes, governments, etc.) gather data, share it with data scientists who, in turn, train some sort of a model using the data. For example, a group of hospitals share MRI scans of tumors with the medical community to help with the development of an automatic tumor segmentation model. One problem with this approach is that the data, once it is shared, might get leaked, misused, or stolen from the developers. Other hurdles include legal reasons that might make it impossible for the hospitals to share and pool the data in the first place.

### 1.3   Federated Learning

Conventional machine learning requires exposing training data to a learning algorithm and its developers. When several data sources are involved, the pooling

together of the data to create a single data set is also required. New approaches like *Federated learning* (FL) [2] allow to separate model training from developer access while also not requiring any pooling of data. FL was introduced in a series of seminal works starting from 2015 [3,4,5]. FL is a protocol consisting of two alternating steps: a) independent training of models on local entities with their respective unique corpus of data, and b) broadcasting back of only the weights of the trained models to a central entity where the weights are aggregated and a new model is redistributed. The choice of which type of model or network to perform step (a) is dictated by the task (e.g., classification, segmentation, etc.) and can be made based on the state-of-the-art in the respective task. The new FL scenario looks like this: A developer sends his or her model to all the institutions that own training data, the institutions locally train the model for the developer and send the newly trained models back. In this way, the developer can train their model while never getting any access to the data. In this setting however a new problem arises.

### 1.4   The aggregation problem

How to aggregate the different models that come back? A naive approach to solve the problem would be:

1. Send an initial model to the first data center
2. Get back a newly trained model and send it to the second data center
3. Repeat until all data centers have trained the model once

Approaches like this are called sequential learning and fail due to a phenomenon called "catastrophic forgetting" [6]. Effectively what would happen is that the final model would only be trained on the data of the last center and would not have generalized to the entire corpus of data. It would simply forget what was learned in the previous center as soon as it gets trained by the next. The state-of-the-art approach tries to avoid this phenomenon by including feedback from every center in each update.

### 1.5   State of the art

The seminal work of learning deep networks from decentralized data [5] proposed as a solution a plain coordinate-wise mean averaging (FedAvg) of the model weights coming separately from multiple centers. Recently [7] proposed a valuable extension to FedAvg, which takes invariance of network weights to permutations into account. In [8] (FedProx), the authors adjust the training loss of a local model to enforce closeness of local and global model updates. Despite methodological advances, there is neither theoretical nor practical evidence for the right recipe when choosing an aggregation strategy. In this paper, we propose a new idea on how to do aggregation. Similar to other initiatives [9,10,11,12], the FETS challenge[1] [13] is organized to benchmark different weight

---

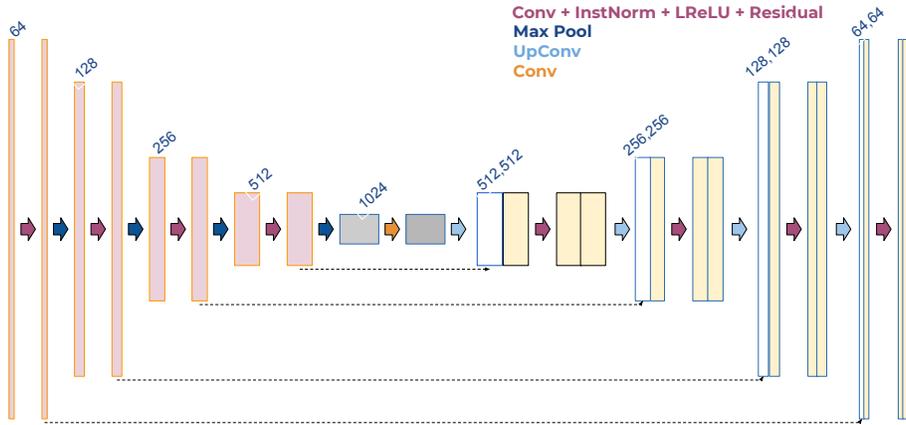[1] https://fets-ai.github.io/Challenge/

**Fig. 2.** 3D U-net architecture as provided by the FETS challenge.

aggregation strategies on the clinically important glioma segmentation problem [14,13,15,16,17]. We contribute to the initiative by proposing an effective extension to the FedAvg strategy. When compared with the other submissions, our model significantly outperformed all of them and won the challenge. On top of that we tested the model locally on a smaller corpus of data to compare it to FedAvg. It notably improves performance compared with FedAvg at no additional compute time.

## 2   Methodology

### 2.1   Segmentation network

The segmentation network is a 3D-Unet. It was provided by the challenge organizers and remained unchanged during all experiments. The architecture is composed of an encoder with residual branches followed by a decoder. We use the *LeakyReLu* activation function [18] along with instance normalization [19] - for mitigating the covariate shift. Dice serves as loss function. Fig. 2 illustrates the schematic of the network.

### 2.2   Federated Cost Weighted Averaging (FedCostWAvg)

The gold standard federated averaging (FedAvg) approach updates the global model as an average of all local models weighted by the respective sizes of the training data set. The new model $M_{i+1}$ is calculated as follows:

$$M_{i+1} = \frac{1}{S} \sum_{j=1}^{n} s_j M_i^j \tag{1}$$

where $s_j$ is the number of samples that model $M^j$ was trained on in round $i$ and $S = \sum_j s_j$. We propose a new weighting strategy that includes the amount by which the cost function decreased during the last step. Using FedCostWAvg, the new model $M_{i+1}$ is calculated as following:

$$M_{i+1} = \sum_{j=1}^{n} (\alpha \frac{s_j}{S} + (1 - \alpha) \frac{k_j}{K}) M_i^j \tag{2}$$

with:

$$k_j = \frac{c(M_{i-1}^j)}{c(M_i^j)}, K = \sum_j k_j \tag{3}$$

where $c(M_i^j)$ is the cost of the model $j$ at timestep $i$ that is simply calculated from the cost function that is being used to train the models locally. $\alpha$ is a parameter ranging between 0 and 1 that can be chosen to determine the balance between data size and cost improvement. In our experiments, a value of $\alpha = 0.5$ performed best. Intuitively, this weighting strategy adjusts not only for the training data set size but also for the size of the local improvements that were made during the last training round. Local updates which only marginally improved the local cost will influence the global update to a lesser extent than those which had a bigger impact.

## 3   Results

The method won the challenge and significantly outperformed all other submitted methods; tables 1 and 2 summarize the performance upon convergence.

In addition we used the provided data (which is a smaller subset of the challenge data) to test the performance of FedCostWAvg against FedAvg in order to visualize the convergence behaviour. We trained and validated the model on 369 samples which were unevenly distributed over 17 data centers. The training-validation split was 80/20, the learning rate was $1e - 4$ and we did 10 epochs per federated round. Please note that computational resources were limited so no exhaustive grid search to find optimal hyperparameters was feasible, also training could not run long enough to achieve maximal performance. Figures 3, 4 and 5 depict the performances over communication rounds. Also note that of course the most informative comparison between methods was done in the challenge itself with more data and many different initialisations. This comparison serves only as a visualization of how different convergence behaviours look like for one initialisation. We observe an improvement for almost all classes and metrics, when using our proposed method. The exemption is the DICE Enhanced Tumor Metric. Note though that the difference is not significant and the methods have not yet converged.

## 3.1    Discussion

While these results already show a clear improvement over FedAvg, it is unclear whether other hyperparameters would have achieved an even better result. Due to limitations in training resources a proper grid search was not feasible.

The simple and straightforward interpretation of the mechanism of FedCost-WAvg is amplification of more informative updates against less informing ones. It could be seen as a diminishing returns acknowledging method. A deeper insight might be the interpretation as resembling a PID controller[2] [20]. When one reframes the federated learning problem as a control problem, then the central server that does the averaging is equivalent to a control unit that is included in a feedback loop. When one would then extend this logic to the averaging approach, it might be intelligent to view FedCostWAvg as an approximation of a PID controller, where the newly added term corresponding to the drop in cost is effectively functioning as the derivative part and the data size term as the proportional one. Future research could try to include the integral term as well.

| Label | DICE WT | DICE ET | DICE TC | Sens. WT | Sens. ET | Sens. TC |
|---|---|---|---|---|---|---|
| Mean | 0,8248 | 0,7476 | 0,7932 | 0,8957 | 0,8246 | 0,8269 |
| StdDev | 0,1849 | 0,2444 | 0,2643 | 0,1738 | 0,2598 | 0,2721 |
| Median | 0,8936 | 0,8259 | 0,9014 | 0,948 | 0,9258 | 0,9422 |
| 25th quantile | 0,8116 | 0,7086 | 0,8046 | 0,9027 | 0,7975 | 0,8258 |
| 75th quantile | 0,9222 | 0,8909 | 0,942 | 0,9787 | 0,9772 | 0,9785 |

**Table 1.** Final performance of FedCostWAvg in the FETS Challenge, DICE and Sensitivity

| Label | Spec WT | Spec ET | Spec TC | H95 WT | H95 ET | H95 TC | Comm. Cost |
|---|---|---|---|---|---|---|---|
| Mean | 0,9981 | 0,9994 | 0,9994 | 11,618 | 27,2745 | 28,4825 | 0,723 |
| StdDev | 0,0024 | 0,0011 | 0,0014 | 31,758 | 88,566 | 88,2921 | 0,723 |
| Median | 0,9986 | 0,9996 | 0,9998 | 5 | 2,2361 | 3,0811 | 0,723 |
| 25th quantile | 0,9977 | 0,9993 | 0,9995 | 2,8284 | 1,4142 | 1,7856 | 0,723 |
| 75th quantile | 0,9994 | 0,9999 | 0,9999 | 8,6023 | 3,5628 | 7,0533 | 0,723 |

**Table 2.** Final performance of FedCostWAvg in the FETS Challenge, Specificity, Hausdorff95 Distance and Communication Cost

---

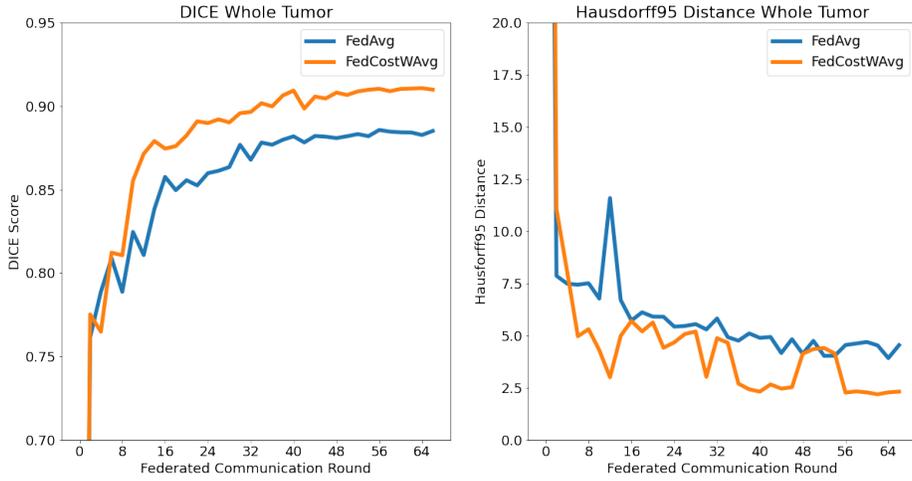[2] The credit for this observation goes to David Naccache.

**Fig. 3.** Comparison of the DICE Whole Tumor metric per federated round for Fed-CostWAvg vs. FedAvg. Note of course that the bigger the DICE score, the better and the smaller the Hausdorff95 distance, the better.
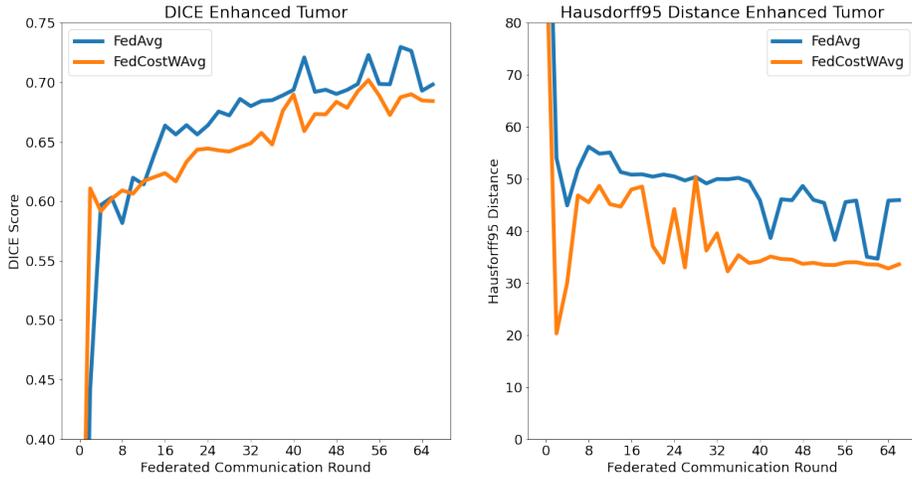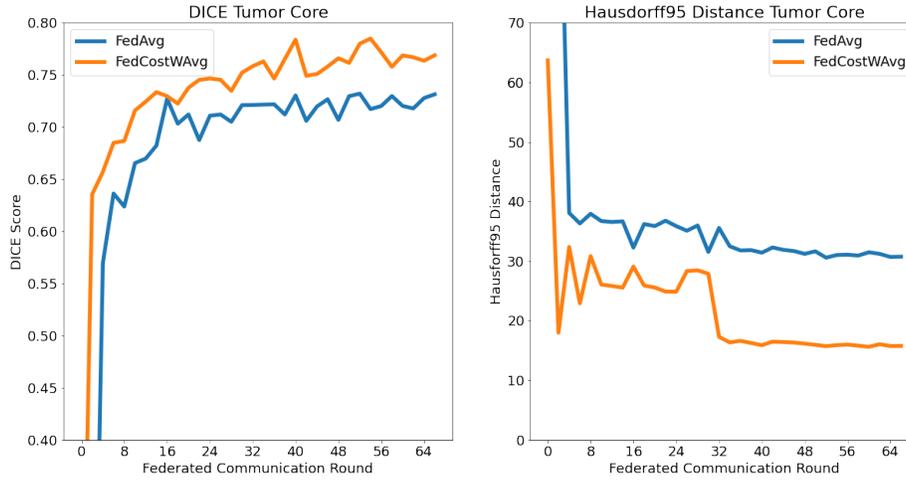


**Fig. 4.** Comparison of the DICE Enhanced Tumor metric per federated round for FedCostWAvg vs. FedAvg. Note of course that the bigger the DICE score, the better and the smaller the Hausdorff95 distance, the better.

**Fig. 5.** Comparison of the DICE Tumor Core metric per federated round for FedCost-WAvg vs. FedAvg. Note of course that the bigger the DICE score, the better and the smaller the Hausdorff95 distance, the better.

## 4   Conclusion

In this paper, we describe a method for model aggregation developed for the MICCAI Federated Tumor Segmentation Challenge (FETS). The novelty of the method lays in including local cost improvements when calculating the weights for averaging models which are trained at different centers. The approach is validated on a brain tumor segmentation task and achieves the best performance among all participating teams.

## Acknowledgements

# References

1. Healthcareitnews.com: Tens of thousands of patient records posted to dark web `https://www.healthcareitnews.com/news/tens-thousands-patient-records-posted-dark-web`, accessed: 2021-07-16.

2. Rieke, N., Hancox, J., Li, W., Milletari, F., Roth, H.R., Albarqouni, S., Bakas, S., Galtier, M.N., Landman, B.A., Maier-Hein, K., et al.: The future of digital health with federated learning. NPJ digital medicine **3**(1) (2020) 1–7

3. Konečný, J., McMahan, H.B., Yu, F.X., Richtárik, P., Suresh, A.T., Bacon, D.: Federated learning: Strategies for improving communication efficiency. CoRR **abs/1610.05492** (2016)

4. Konečný, J., McMahan, B., Ramage, D.: Federated optimization: Distributed optimization beyond the datacenter. CoRR **abs/1511.03575** (2015)

5. McMahan, B., Moore, E., Ramage, D., Hampson, S., y Arcas, B.A.: Communication-efficient learning of deep networks from decentralized data. In: Artificial intelligence and statistics, PMLR (2017) 1273–1282

6. McCloskey, M., Cohen, N.J.: Catastrophic interference in connectionist networks: The sequential learning problem. In Bower, G.H., ed.: Psychology of Learning and Motivation. Volume 24. Academic Press (1989) 109–165

7. Yurochkin, M., Agarwal, M., Ghosh, S., Greenewald, K., Hoang, N., Khazaeni, Y.: Bayesian nonparametric federated learning of neural networks. In: International Conference on Machine Learning, PMLR (2019) 7252–7261

8. Sahu, A.K., Li, T., Sanjabi, M., Zaheer, M., Talwalkar, A., Smith, V.: On the convergence of federated optimization in heterogeneous networks. arXiv preprint arXiv:1812.06127 **3** (2018) 3

9. Sekuboyina, A., et al.: Verse: a vertebrae labelling and segmentation benchmark. arXiv preprint arXiv:2001.09193 (2020)

10. Payette, K., de Dumast, P., Kebiri, H., Ezhov, I., Paetzold, J.C., Shit, S., Iqbal, A., Khan, R., Kottke, R., Grehten, P., et al.: A comparison of automatic multi-tissue segmentation methods of the human fetal brain using the feta dataset. arXiv e-prints (2020) arXiv–2010

11. Paetzold, J.C., McGinnis, J., Shit, S., Ezhov, I., Büschl, P., Prabhakar, C., Todorov, M.I., Sekuboyina, A., Kaissis, G., Ertürk, A., et al.: Whole brain vessel graphs: A dataset and benchmark for graph learning and neuroscience (vesselgraph). arXiv preprint arXiv:2108.13233 (2021)

12. Bilic, P., Christ, P., Vorontsov, E., et al.: The liver tumor segmentation benchmark (lits). arXiv preprint arXiv:1901.04056 (2019)

13. Pati, S., Baid, U., Zenk, M., Edwards, B., Sheller, M., Reina, G.A., Foley, P., Gruzdev, A., Martin, J., Albarqouni, S., et al.: The federated tumor segmentation (fets) challenge. arXiv preprint arXiv:2105.05874 (2021)

14. Bakas, S., Akbari, H., Sotiras, A., Bilello, M., Rozycki, M., Kirby, J.S., Freymann, J.B., Farahani, K., Davatzikos, C.: Advancing the cancer genome atlas glioma mri collections with expert segmentation labels and radiomic features. Scientific data **4**(1) (2017) 1–13

15. Reina, G.A., Gruzdev, A., Foley, P., Perepelkina, O., Sharma, M., Davidyuk, I., Trushkin, I., Radionov, M., Mokrov, A., Agapov, D., et al.: Openfl: An open-source framework for federated learning. arXiv preprint arXiv:2105.06413 (2021)

16. Sheller, M.J., Edwards, B., Reina, G.A., Martin, J., Pati, S., Kotrotsou, A., Milchenko, M., Xu, W., Marcus, D., Colen, R.R., et al.: Federated learning in medicine: facilitating multi-institutional collaborations without sharing patient data. Scientific reports **10**(1) (2020) 1–12

17. Kofler, F., Berger, C., Waldmannstetter, D., Lipkova, J., Ezhov, I., Tetteh, G., Kirschke, J., Zimmer, C., Wiestler, B., Menze, B.H.: Brats toolkit: translating brats brain tumor segmentation algorithms into clinical and scientific practice. Frontiers in neuroscience **14** (2020) 125
18. Maas, A.L., Hannun, A.Y., Ng, A.Y., et al.: Rectifier nonlinearities improve neural network acoustic models. In: Proc. icml. Volume 30., Citeseer (2013) 3
19. Ulyanov, D., Vedaldi, A., Lempitsky, V.: Instance normalization: The missing ingredient for fast stylization. arXiv preprint arXiv:1607.08022 (2016)
20. Bellman, R.E.: Adaptive control processes. Princeton university press (2015)