

Federated Learning in ASR: Not as Easy as You Think

Wentao Yu¹, Jan Freiwald¹, Sören Tewes, Fabien Huennemeyer, Dorothea Kolossa

Institute of Communication Acoustics, Ruhr University Bochum, Germany

Email: {wentao.yu, jan.freiwald, soeren.tewes, fabien.huennemeyer, dorothea.kolossa}@rub.de

Abstract 现状 --> 隐私保护问题 --> on-device FL

With the growing availability of smart devices and cloud services, personal speech assistance systems are increasingly used on a daily basis. Most devices redirect the voice recordings to a central server, which uses them for upgrading the recognizer model. This leads to major privacy concerns, since private data could be misused by the server or third parties. Federated learning is a decentralized optimization strategy that has been proposed to address such concerns. Utilizing this approach, private data is used for on-device training. Afterwards, updated model parameters are sent to the server to improve the global model, which is redistributed to the clients. In this work, we implement federated learning for speech recognition in a hybrid and an end-to-end model. We discuss the outcomes of these systems, which both show great similarities and only small improvements, pointing to a need for a deeper understanding of federated learning for speech recognition.

Index Terms: speech recognition, privacy, federated learning, end-to-end models, hybrid ASR

1 Introduction

In order to obtain training data that is optimally matched to the typical use cases and application environments of automatic speech recognition (ASR) systems, data of the users is sent to a central server. In many commercial applications, like Amazon’s Alexa, user speech recordings are decoded on a central server, where the service provider can incorporate the speech samples into their training and testing datasets. This practice leads to massive datasets, spanning all kinds of different user characteristics and environments, which can be utilized to train highly optimized recognizer models.

In terms of privacy, this process also enables companies—or anyone, who may be able to access their data collections—to extract highly sensitive information about their users, such as identity, language, gender, personal conversation content and even emotional and health status. For example, in [1], the speaker’s age, height and weight information are extracted from the audio signal.

This leads directly into a dilemma: while big-data-driven AI shows exciting results in many different areas, centralized data processing leads to concerns about the security and the privacy of such systems and is not in line with the privacy-by-design principle stipulated in the European Union General Data Protection Regulation [2]. To counter this dilemma, the technique of federated learning (FL), as proposed by H. Brendan McMahan, et al. [3–6], can be used. The goal is to learn machine-learning models on data gathered on multiple private devices, while at the same time keeping all this re-training data private. For this purpose, federated learning proposes a decentralized optimization process. As shown in Figure 1, in this way, it ensures user

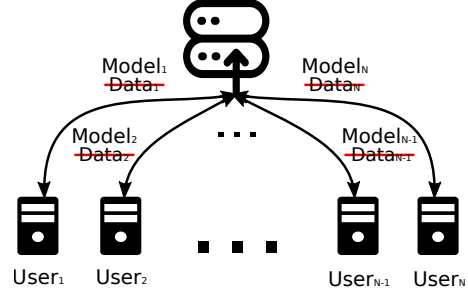


Figure 1: In a federated learning setup, no user data is transmitted to the central server. Instead, the server rolls out a model to the users, which is then updated locally and sent back to the server. The server calculates an updated model and repeats the cycle.

privacy by training on any private data locally, on-device, instead of sending private data to central servers. Each of the clients that are involved in re-training then sends a re-trained model to the central server, and the central server updates its global model based on these user models, while all the clients’ private data is kept on-device. In an optimal setup, the clients and the central server keep an appropriate communication frequency, maintaining a balance between communication cost and model performance [7, 8].

In recent years, federated learning was applied successfully in many different tasks. In [9, 10], federated learning is used in **wake-word detection** for a digital assistant. Federated learning was also applied in **speech emotion detection** [11]. In [12], federated learning is used for mobile keyboard prediction. Another example is in the **medical** area, where more and more data is available from hospitals, clinics, or patients, but strict privacy requirements limit its use. In this situation, federated learning promises a compromise that maintains privacy while allowing to improve the quality of healthcare through big-data-driven AI [13].

There are a few prior works on federated learning for large-vocabulary **ASR**. For example, in [14], federated learning is applied with an end-to-end (E2E) model on the French set of the Common Voice dataset [15]. In [16], an RNN-T architecture ASR [17] is used on the LibriSpeech corpus, and [18] applies FL to a hybrid ASR model. Additionally, in [19] Dimitriadis et.al. introduce a federated learning setup for ASR based on an end-to-end system.

In this paper, we directly compare the federated learning performance of two widely used ASR model types, i.e. hybrid ASR based on the Pytorch-Kaldi system [20] and the E2E transformer/CTC model, trained via ESPnet [21]. We use the Librispeech corpus [22] for all experiments, which contains 960h of speech data from 2484 speakers. Specifically, **we consider a use case, where every client node only possesses data of one speaker, although scenarios with multiple speakers per client are conceivable (e.g. cross-silo federated learning in [6]).** FL的实验设定

¹The authors contributed equally.

We compare the performance and the communication cost of a number of federated learning strategies with different client-server communication frequencies and different model update strategies, finding that despite their large differences in approach and implementation, hybrid and E2E models show similar behavior, which raises the important question of how to achieve effective, decentralized learning in sequence-labeling tasks in general and ASR, specifically. To allow for easy reproducibility, we are also releasing the script that generates this dataset, as well as our training and test recipes¹.

The paper is organized as follows: The idea of federated learning and the different update strategies are detailed in Section 2, followed by a brief introduction to the PyTorch-Kaldi and ESPnet models. Section 3 describes our reorganization of the Librispeech dataset, which allows us to evaluate FL strategies in a meaningful manner. We present the results in Section 4 and discuss the outcomes and the questions that they raise in Section 5.

2 System overview

In this section, we take a look at federated averaging in general. Our experiments rely on two ASR toolkits, PyTorch-Kaldi and ESPnet. Hence we give a short introduction to both frameworks and describe our training configurations for both.

2.1 Federated averaging

Federated averaging (*FedAvg*) [3] is a central algorithmic step in federated learning. As shown in Algorithm 1, a pretrained model \mathbf{W}^0 initializes all clients at the beginning of the learning process. This is followed by an iterative optimization loop, which contains two steps per round: the first step is the local optimization, which is performed by each client. Each client’s model is initialized by the previous global model \mathbf{W}^{i-1} . To obtain the new client-side local model, each client updates their model with their own private dataset, which is kept on-device. The second step is the *FedAvg* step, which is applied over all selected client models. The averaged model is considered as the new global model in the next round. In contrast to other works, averaging the accumulated gradient updates (e.g. [3, 16]), we directly average the model parameters at the end of each round, which leads to the same outcome with an easier implementation:

$$\begin{aligned}\mathbf{W}^i &= \sum_n \alpha_n \cdot \mathbf{W}_n^i \\ &= \sum_n \alpha_n \cdot \mathbf{W}^{i-1} + \sum_n \alpha_n \cdot \Delta \mathbf{W}_n^i \\ &= \mathbf{W}^{i-1} + \sum_n \alpha_n \cdot \Delta \mathbf{W}_n^i,\end{aligned}\quad (1)$$

where α_n are the *FedAvg* weights described in Algorithm 1 and $\Delta \mathbf{W}_n^i$ is the model parameter update for client n in the i -th iteration.

In Algorithm 1, the federated learning round can be performed with different communication frequencies. We have identified the following levels for defining these:

- **Batch-level (B):** The communication is always performed after K mini-batches. In every communication

round, the selected client examines whether there are untrained utterances for this epoch. If the untrained utterances of the current speaker do not suffice for K mini-batches, the batch size is reduced. The client stops sending a new model, when there is no more training data in the current epoch. After an epoch is completed, a new round begins with the complete set of clients participating again.

- **Epoch-level (E):** the update is performed with the selected clients’ models after every epoch. In the ESPnet setup, we can select all client models for the *FedAvg* update, or alternatively, we can randomly select N client models for updating. In PyTorch-Kaldi, we can also perform updates after a fixed fraction of an epoch.
- **Convergence-level (C):** there is only one update, which is performed with the selected clients’ final models, after convergence.

The *FedAvg* weights α_n in Algorithm 1 can be defined as:

- **Mean (M):** $\alpha_n = \frac{1}{N}$, where N is the number of the selected clients.
- **Weighted (W):** $\alpha_n = \frac{|\mathbf{D}_n|}{\sum_i |\mathbf{D}_i|}$, where $|\mathbf{D}_n|$ is the number of utterances of the client-side dataset for client n .

Algorithm 1 Federated averaging. There are N clients. The initial dataset is \mathbf{D}_0 . The federated learning training-set \mathbf{T}_n^i for client n in round i is derived from the complete client-side training set \mathbf{D}_n .

```

1:  $\mathbf{W}^0 = \text{train\_initial\_model}(\mathbf{D}_0)$ 
2: for round  $i \in [1, 2, 3, \dots]$  do
3:   for client  $n \in [1..N]$  do
4:      $\mathbf{W}_n^i = \mathbf{W}^{i-1}$ 
5:      $\mathbf{T}_n^i = \text{subset}(\mathbf{D}_n, i)$ 
6:      $\mathbf{W}_n^i = \text{train\_model}(\mathbf{T}_n^i, \mathbf{W}_n^i)$ 
7:   end for
8:    $\mathbf{W}^i = \sum_n \alpha_n \cdot \mathbf{W}_n^i$ , where  $\sum_n \alpha_n = 1$ 
9: end for
```

2.2 PyTorch-Kaldi

PyTorch-Kaldi[20] allows for high flexibility when using DNNs in training hybrid models for use with Kaldi. For this purpose, it provides a model configuration interface, as well as training recipes.

Our overall setup is divided into three parts, **Kaldi training**, **PyTorch-Kaldi training**, and **FedAvg**. First, the features of all relevant datasets are extracted, the corresponding GMM models are trained, and alignments are obtained in Kaldi [23], using MFCCs with first and second derivatives and cepstral mean and variance normalization on a per-speaker basis. In the second part, we use PyTorch-Kaldi to train a hybrid DNN-WFST model with the PyTorch library. We use the standard MLP structure that is provided with many of the examples in PyTorch-Kaldi, a feed-forward network with a context size of 5 frames [20]. This results in the initial model \mathbf{W}^0 , which is only learned on the `initial` dataset. After the server-side training is thus completed, federated learning is performed as described in Algorithm 1. Here, client-side trainings only re-train the DNN of the hybrid model, leaving the WFST unchanged.

The initial model is trained **for 24 epochs** and distributed to all clients, followed by 12 federated learning epochs with different communication frequencies. In order to use the

¹<https://github.com/rub-ksv/Federated-learning-ASR>

standard training script for federated learning, we changed it to pause and save a checkpoint model after a specifiable number of data chunks; **every clients' training data** is divided into $n = 8$ chunks. This allows us to train our network on each of our clients, then pause the training, average every updated model and redistribute the result to all clients. In this way, we can perform (potentially partial) epoch-level training, as described in Section 2.1.

The scripts are based on training recipes provided for Kaldi [23] and PyTorch-Kaldi [20]. We used, among others, the WSJ, Librispeech and CommonVoice recipes as a guide for structuring and parameterizing the training and decoding of our models.

2.3 ESPnet setup

The ESPnet end-to-end model [21] uses a sequence-to-sequence (S2S) transformer model with connectionist temporal classification (CTC) optimization. As described in [24], CTC learns to align features and transcription optimally, which leads to fast convergence. The model is updated using a linear combination of CTC and S2S loss terms as the **objective function**

$$L = \eta \cdot \log p_{\text{CTC}}(\mathbf{s}|\mathbf{o}) + (1 - \eta) \log p_{\text{S2S}}(\mathbf{s}|\mathbf{o}), \quad (2)$$

with \mathbf{s} as the states and the constant hyper-parameter η , which controls the strength of both terms. The joint S2S-CTC model performs token-by-token decoding, until the end-of-sentence symbol (EOS). In each step, the log-posteriors of the current token from the Transformer and CTC models are linearly combined with an RNN-language model $p_{\text{LM}}(\mathbf{s}_t)$

$$\log p^*(\mathbf{s}_t|\mathbf{o}) = \eta \log p_{\text{CTC}}(\mathbf{s}_t|\mathbf{o}) + (1 - \eta) \log p_{\text{S2S}}(\mathbf{s}_t|\mathbf{o}) + \theta \log p_{\text{LM}}(\mathbf{s}_t), \quad (3)$$

where **θ controls the contribution of the language model** and is set to 0.3.

All models use 80 log-mel filterbank features together with pitch, delta pitch, and the probability of voicing. These 83-dimensional features are extracted using 25 ms frame size and 10 ms frame shift.

For the language model, we use a 4-layer recurrent neural network, which predicts a character at a time and receives the previous character as the input. Each layer consists of 2048 units. The language model is trained on the complete LibriSpeech corpus for 20 epochs using SGD.

The dataset is divided into a server-side dataset and user-side datasets (see Section 3). The **initial** training set is used to train the initial model \mathbf{W}^0 . All ESPnet experiments are trained using NVIDIA's Volta-based DGX-1 multi GPU system. The initial model is trained on 3 Tesla V100 GPUs, each with 32 GB memory, for 100 epochs at a batch size of 32. The upper bound reference model, trained on the **complete** set, uses the same model configuration and hyperparameter settings. For the initial and the reference model, training takes two days, and four days, respectively. In the federated learning stage, the batch size is set to 8. The hyper-parameter η is set to 0.3 during training and $\eta = 0.5$ during decoding.

For C-level averaging, where each client trains their model to convergence, early stopping is used, and triggered if the evaluation accuracy does not improve over 10 epochs. Both M-averaging and W-averaging are applied in the C-level experiments. As described in Section 2.1, we either

use all clients to update the model or randomly select N of them in each round. To assess the impact of this choice, we repeat the E-level training with $N = 100$, $N = 350$ and $N = \text{all}$.

3 Dataset

For this work, we reorganized the Librispeech[22] dataset, creating dedicated training, test and development sets to **simulate a realistic federated learning environment**.

Therefore, we created a server-side training set for our initial and alignment models (**initial**). Additionally, we created user datasets, **which contain training, development, and test sets for every client**. All 1372 clients are disjoint from the initial training set and from each other. Furthermore, the union of the federated and initial training sets is termed **complete** and used to obtain an upper performance bound, corresponding to what would be seen if all data were located on a central server, i.e. with standard, non-privacy-preserving learning.

For testing, we created three sets: **unseen** is composed of previously unseen speakers, **initial** contains speakers from the training set of the initial model, and **federated** contains speakers from the federated training set, so it simulates clients that can profit from their data contribution.

Statistics of this dataset partitioning are shown in Tab. 1.

Table 1: Dataset composition of all training and test sets. The proportion of male and female speakers is close to 50%.

	Dataset	Utterances	Duration [h]
Training	initial	37768	134.5
	federated	110570	367.8
	complete	148338	502.3
Test	initial	487	1.7
	federated	682	2.3
	unseen	476	1.1

4 Results

4.1 Hybrid ASR

The results of the PyTorch-Kaldi experiments are summarized in Table 2. As expected, the performance improves along the training phases, from the GMM-HMM alignment model **via the initial model to the FL-models**. In all cases, the upper bound reference model performs best. FL-training was performed at the E-level, where the number in parentheses shows the number of epochs after which communication takes place. We find that more frequent update rounds of *FedAvg* slightly improve the overall results, with best performance observed for updates after every 1/2 or 1/4 of an epoch. The highest improvement through FL-training is seen on the **federated** test set, which contains utterances of the FL speakers. The recognition quality of the other test sets does not degrade after FL-training.

4.2 End-to-end learning

Table 3 shows the experimental results using the ESPnet toolkit. Analogously to the PyTorch-Kaldi experiments, the upper-bound reference model provides best performance. For C-level FL, both averaging methods are applied, but as there is no large difference between W- and M-averaging, we applied only W-averaging in the later experiments.

Table 2: WER (%) of hybrid models trained conventionally and using E-level federated learning, with associated communication cost [GB]. The final two rows show the WERs of upper bound reference models trained on the `complete` set. Ref. 1 learns its decision tree on the `complete` set; Ref. 2 uses the decision tree of the initial model.

	unseen	federated	initial	cost[GB]
Aligner	32.76	33.22	26.96	-
Initial	19.06	18.18	14.05	-
E(2)-M	19.08	17.86	14.05	326
E(1)-M	19.01	17.79	14.05	651
E(1/2)-M	19.02	17.67	13.90	1302
E(1/4)-M	18.94	17.71	13.94	2604
Ref. 1	17.77	15.78	13.24	-
Ref. 2	17.70	16.13	13.29	-

The E-level experiments (E-100-W, E-350-W, and E-all-W) are trained on randomly selected $N = [100, 350, \text{all}]$ speakers per round. As expected, we see that a larger number of clients N tends towards a better performance.

We also tested the B-level strategy (**B**), updating the model after every minibatch, which we hoped would give optimal performance, as it best approximates the situation in non-FL training. However, while the communication cost became higher as expected, there are slight performance improvements at best, which do not justify this extra effort.

Table 3: WER (%) of end-to-end models trained on the `initial` set and subsequently applying FL. Last row: WER of a reference model trained on the `complete` set. C, E, B represent C-, E-, and B-level updates; M and W indicate M- and W-averaging.

	unseen	federated	initial	cost [GB]
Initial	8.98	6.94	4.23	-
C-M	9.00	6.88	4.18	317
C-W	9.01	6.88	4.15	317
E-100-W	8.91	6.89	4.17	925
E-350-W	8.87	6.91	4.14	1618
E-all-W	8.81	6.86	4.17	4758
B-W	8.93	6.81	4.15	13319
Ref	6.08	3.37	3.07	-

We also see that the C-level updates, most effective in terms of communication cost, can still improve the recognition accuracy compared to the pretrained model, but the improvement does not reach that of the E-level updates. E-level updates with N randomly selected speakers may also be closest to a realistic scenario, as we would not wish to ask every client to participate in every one of the model updates.

5 Conclusions

We have shown that federated learning can slightly improve the recognition rates in automatic speech recognition, both using hybrid and end-to-end models. We suspect that one reason for the disappointing performance is to be found at the heart of the federated averaging approach. During FL-training, federated averaging leads to smaller gradients and, hence, to slow improvements. While it is possible

that better performance may be attained by engaging more clients and allowing for more model updates, this comes at an excruciating communication cost. Therefore, we believe that it is necessary to focus on more effective collaborative learning strategies, e.g., based on suitable client selection criteria or on re-training only parts of the networks.

Furthermore, McMahan et al. [6] point to an inherent problem of the federated learning approach, when it is applied on non-Independent, Identically Distributed (non-IID) client data. Since on-device training data is naturally grouped by speaker, there are different statistical properties in every client, which may degrade the model performance. To counter this problem, a cross-silo training approach can be beneficial and will be subject of our future investigations. Federated Variational Noise (FVN) will also be investigated, which has shown its value in recovering FL-ASR performance despite non-IID client data [16].

In addition to this global view, it is also interesting to look at the model architectures in detail:

For hybrid models, federated learning imposes several difficulties and impracticalities. Due to the comparatively smaller size of the initial training data, decision tree learning results in a smaller output layer size of the FL-trained models in contrast to the full model. This also hints at a similar issue in the calculation of the language model, which we did not consider here, but which would also likely have a smaller size and larger perplexity if it is not updated—while simultaneously, any updates of the language model by nature lead to a leakage of private speech content.

We also found that fewer iteration steps between *FedAvg* improve the quality of the model, but impose a heavy computational burden. Finally, the alignment model is not updated at the moment, and no realignments of client training data were performed after model averaging, which may also be a factor in the observed performance issues. This additionally raises the interesting question, in how far—explicit or implicit—alignment is an issue in federated learning for time-series modeling in general.

In contrast, the end-to-end-model suffers neither from an impacted decision-tree nor from a missing explicit re-alignment. We conducted many experiments to assess the influence of different parameters and hyper-parameters of federated training, finding that allowing more participants in each round of model averaging yields better performance, and that in trying to balance the performance and cost, epoch-level averaging appears most promising.

All in all, however, we see the outcomes of our experiment as a pointer to underlying issues of the idea of federated averaging, when it is applied to complex, highly variable time series data, as seen in large-vocabulary speech recognition. We are releasing our source code in the hope of setting a starting point for further investigations into the important, yet also challenging question of how to achieve a reasonable balance of performance and cost in privacy-preserving machine learning.

Acknowledgements

This work was funded by the PhD School “SecHuman - Security for Humans in Cyberspace” by the federal state of NRW, and by the German Federal Ministry of Education and Research (BMBF) within the “Innovations for Tomorrow’s Production, Services, and Work” Program (02L19C200), a project that is implemented by the Project Management Agency Karlsruhe (PTKA). The authors are responsible for the content of this publication.

References

- [1] S. Galgali, S. S. Priyanka, B. R. Shashank, and A. P. Patil, "Speaker profiling by extracting paralinguistic parameters using mel frequency cepstral coefficients," in *Proc. International Conference on Applied and Theoretical Computing and Communication Technology*, 2015, pp. 486–489.
- [2] European Parliament and Council, "Regulation (EU) 2016/679 of the European Parliament and of the Council of 27 April 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46/EC (General Data Protection Regulation)," 2016.
- [3] H. B. McMahan, E. Moore, D. Ramage, and B. Agüera y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Artificial Intelligence and Statistics*. PMLR, 2017, pp. 1273–1282.
- [4] J. Konečný, H. B. McMahan, D. Ramage, and P. Richtárik, "Federated optimization: Distributed machine learning for on-device intelligence," *arXiv preprint 1610.02527*, 2016.
- [5] J. Konečný, H. B. McMahan, F. Yu, P. Richtárik, A. T. Suresh, and D. Bacon, "Federated learning: Strategies for improving communication efficiency," *arXiv preprint 1610.05492*, 2016.
- [6] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode *et al.*, "Advances and open problems in federated learning," *Foundations and Trends® in Machine Learning*, vol. 14, no. 1, 2021.
- [7] J. Konečný, H. B. McMahan, F. Yu, P. Richtárik, A. T. Suresh, and D. Bacon, "Federated learning: Strategies for improving communication efficiency," *arXiv preprint arXiv:1610.05492*, 2016.
- [8] J. Hamer, M. Mohri, and A. T. Suresh, "Fedboost: A communication-efficient algorithm for federated learning," in *International Conference on Machine Learning*. PMLR, 2020, pp. 3973–3983.
- [9] D. Leroy, A. Coucke, T. Lavril, T. Gisselbrecht, and J. Dureau, "Federated learning for keyword spotting," *arXiv 1810.05512*, 2019.
- [10] A. Hard, K. Partridge, C. Nguyen, N. Subrahmanya, A. Shah, P. Zhu, I. L. Moreno, and R. Mathews, "Training keyword spotting models on non-IID data with federated learning," *arXiv 2005.10406*, 2020.
- [11] S. Latif, S. Khalifa, R. Rana, and R. Jurdak, "Poster abstract: Federated learning for speech emotion recognition applications," in *Proc. IPSN*, 2020, pp. 341–342.
- [12] A. Hard, K. Rao, R. Mathews, S. Ramaswamy, F. Beaufays, S. Augenstein, H. Eichner, C. Kiddon, and D. Ramage, "Federated learning for mobile keyboard prediction," *arXiv preprint arXiv:1811.03604*, 2018.
- [13] J. Xu, B. S. Glicksberg, C. Su, P. Walker, J. Bian, and F. Wang, "Federated learning for healthcare informatics," *Journal of Healthcare Informatics Research*, vol. 5, no. 1, pp. 1–19, 2021.
- [14] Y. Gao, T. Parcollet, J. Fernandez-Marques, P. P. de Gusmao, D. J. Beutel, and N. D. Lane, "End-to-end speech recognition from federated acoustic models," *arXiv preprint arXiv:2104.14297*, 2021.
- [15] R. Ardila, M. Branson, K. Davis, M. Henretty, M. Kohler, J. Meyer, R. Morais, L. Saunders, F. M. Tyers, and G. Weber, "Common voice: A massively-multilingual speech corpus," *arXiv preprint arXiv:1912.06670*, 2019.
- [16] D. Guliani, F. Beaufays, and G. Motta, "Training speech recognition models with federated learning: A quality/cost framework," 2020.
- [17] Y. He, T. N. Sainath, R. Prabhavalkar, I. McGraw, R. Alvarez, D. Zhao, D. Rybach, A. Kannan, Y. Wu, R. Pang, Q. Qiao Liang, D. Bhatia, Y. Shangguan, B. Li, G. Pundak, K. C. Sim, T. Bagby, S.-Y. Chang, K. Rao, and A. Gruenstein, "Streaming end-to-end speech recognition for mobile devices," in *Proc. ICASSP*, 2019, pp. 6381–6385.
- [18] X. Cui, S. Lu, and B. Kingsbury, "Federated acoustic modeling for automatic speech recognition," *arXiv 2102.04429*, 2021.
- [19] D. Dimitriadis, K. Kumatani, R. Gmyr, Y. Gaur, and S. Eskimez, "A federated approach in training acoustic models," in *Proc. Interspeech*, 2020.
- [20] M. Ravanelli, T. Parcollet, and Y. Bengio, "The PyTorch-Kaldi Speech Recognition Toolkit," in *Proc. ICASSP*, 2019.
- [21] S. Watanabe, T. Hori, S. Karita, T. Hayashi, J. Nishitoba, Y. Unno, N. Enrique, J. Heymann, M. Wiesner, N. Chen, A. Renduchintala, and T. Ochiai, "ESPnet: End-to-End Speech Processing Toolkit," in *Proc. Interspeech*, 2018, pp. 2207–2211.
- [22] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: An ASR corpus based on public domain audio books," in *Proc. ICASSP*, 2015, pp. 5206–5210.
- [23] D. Povey, A. Ghoshal, G. Boulianne, L. Burget, O. Glembek, N. Goel, M. Hannemann, P. Motlicek, Y. Qian, P. Schwarz, J. Silovsky, G. Stemmer, and K. Vesely, "The Kaldi speech recognition toolkit," in *Proc. ASRU*, 2011.
- [24] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, "Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks," in *Proceedings of the 23rd International Conference on Machine learning*, 2006, pp. 369–376.