

Scaling Federated Learning for Fine-tuning of Large Language Models

Agrin Hilmkil*, Sebastian Callh*, Matteo Barbieri*, Leon René Sütfield+, Edvin Listo Zec+, and Olof Mogren+

*Peltarion, peltarion.com

+RISE Research Institutes of Sweden, ri.se

Abstract

Federated learning (FL) is a promising approach to distributed compute, as well as distributed data, and provides a level of privacy and compliance to legal frameworks. This makes FL attractive for both consumer and healthcare applications. While the area is actively being explored, few studies have examined FL in the context of larger language models and there is a lack of comprehensive reviews of robustness across tasks, architectures, numbers of clients, and other relevant factors. In this paper, we explore the fine-tuning of Transformer-based language models in a federated learning setting. We evaluate three popular BERT-variants of different sizes (BERT, ALBERT, and DistilBERT) on a number of text classification tasks such as sentiment analysis and author identification. We perform an extensive sweep over the number of clients, ranging up to 32, to evaluate the impact of distributed compute on task performance in the federated averaging setting. While our findings suggest that the large sizes of the evaluated models are not generally prohibitive to federated training, we found that the different models handle federated averaging to a varying degree. Most notably, DistilBERT converges significantly slower with larger numbers of clients, and under some circumstances, even collapses to chance level performance. Investigating this issue presents an interesting perspective for future research.

1 Introduction

Transformer-based architectures such as BERT have recently lead to breakthroughs in a variety of language-related tasks, such as document classification, sentiment analysis, question answering, and various forms of text-mining (Vaswani et al., 2017; Devlin et al., 2019; Adhikari et al., 2019; Sun et al., 2019a; Yang et al., 2019; Lee et al., 2020). These models create semantic representations of

text, which can subsequently be used in many downstream tasks (Devlin et al., 2019). The training process for Transformers typically includes two phases: During pre-training, the model learns to extract semantic representations from large, task-independent corpora. The pre-training is followed by task-specific fine-tuning on a separate dataset to optimize model performance further.

In this paper, we study the effects of fine-tuning Transformer-based architectures in a federated learning (FL) setting. In FL, models are trained in a decentralized fashion on a number of local compute instances, called clients, and intermittently aggregated and synchronized via a central server. As such, FL is a solution for distributed compute, as well as distributed data, and provides a level of privacy with regards to the sharing of personal or otherwise sensitive data. Model aggregation is commonly performed via averaging of the weights of the individual client models, called Federated Averaging (FEDAVG) (McMahan et al., 2017a).

Depending on the application, the number of clients in an FL setting can differ wildly. In instances where smartphones are used as clients, their number can reach into the millions (Hard et al., 2018), whereas settings with higher compute requirements and more data per client will often range between a handful and a few dozens of clients. Here, we focus on the latter, as training large language models requires a lot of compute. A potential application of this is the medical field, in which automated analyses of electronic health records yield enormous potential for diagnostics and treatment-related insights (Zeng et al., 2018).

Our contribution: We provide a comprehensive overview of the applicability of the federated learning setting to large language models. To this end, we work with a fixed computation budget for each task, and use a fixed total amount of data while

varying the number of clients between which the data is split up. This way, we isolate the effects of distributing data over several clients for distributed compute. We leave comparisons with a fixed amount of data per client and varying non-i.i.d. data distributions between clients for future work. The main contributions of this paper are the following: (1) We provide a comparison of three popular Transformer-based language models in the federated learning setting, using the IMDB, Yelp F, and AG News datasets. (2) We analyze how the number of clients impacts task performance across tasks and model architectures.

2 Related work

Federated optimization was first introduced by (Konečný et al., 2015). The key challenges in this paradigm are communication efficiency when learning from many clients, privacy concerns with respect to leakage of client data, and variability in data distributions between clients (non-i.i.d. setting). FEDAVG (McMahan et al., 2017a) solves the federated optimization problem by building a global model based on local stochastic gradient descent updates and has been shown to work on non-i.i.d. data. Since then, many adaptations have arisen (Li et al., 2019; Mohri et al., 2019; Karimireddy et al., 2019). Guha et al. (2019) proposes a one-shot FL algorithm, learning a global model efficiently in just one communication round. Zhao et al. (2018), Hsu et al. (2019) and Listo Zec et al. (2020) study effects of FEDAVG and non-i.i.d. client data. McMahan et al. (2017b) and Hard et al. (2018) train large recurrent language models with user-level differential privacy guarantees and for mobile keyboard prediction, respectively. Ge et al. (2020) use federated learning for named entity recognition with heterogeneous medical data.

Regarding model size, most architectures used in FL to date are relatively small (e.g., CIFG for mobile keyboard prediction: 1.4M parameters (Hard et al., 2018)), compared to BERT-based language models with hundreds of millions of parameters. How these very large models behave under FEDAVG remains underexplored. To the best of our knowledge, Lin et al. (2020) and Liu and Miller (2020) are the first ones to train large Transformer models in a federated setting. Liu and Miller (2020) trained BERT on a medical corpus and showed that both pre-training and fine-tuning could be done in a federated manner with only minor declines in task

performance. Nonetheless, the study is mainly a proof-of-concept and does not explore many of the factors that can be expected in real-world scenarios. For instance, the authors only used five clients, and evaluated them only on i.i.d. data. Lin et al. (2020) introduces FedDF, an ensemble distillation algorithm for model fusion. The authors train a central model through unlabeled data on the client models outputs, and perform fine-tuning on a pre-trained DistilBERT (Sanh et al., 2019) in a federated setting as a baseline. To the best of our knowledge, no systematic variation of the number of clients and other relevant factors has previously been explored in this context.

3 Method

3.1 Federated learning

Federated learning aims to solve the optimization problem

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{K} \sum_{k=1}^K F_k(\theta), \quad (1)$$

where $F_k(\theta) = \mathbb{E}_{x \sim \mathcal{D}_k} [\ell_k(\theta; x)]$ is the expected loss on client k and \mathcal{D}_k is the data distribution of client k . In FEDAVG, a global model f_θ is initialized on a central server and distributed to all K clients, each of which then trains its individual copy of the network using SGD for E local epochs with local batch size B . The clients' updated parameters are then averaged on the central server, weighted by the local data size at each client. The averaged model is distributed to the clients again, and the process is repeated for a defined number of communication rounds.

We implement FEDAVG using distributed PyTorch (Paszke et al., 2019). For each experiment we start from a pre-trained model, and fine-tune it with federated averaging on the current task.

3.2 Models

We include BERT with 110M parameters, 12 layers (Devlin et al., 2019), ALBERT with 11M parameters, 12 layers (Lan et al., 2020) and DistilBERT with 65M parameters, 6 layers (Sanh et al., 2019). This allows us to study the effect that both the parameter count and the number of layers have on FEDAVG. All models are the corresponding base models pre-trained on (cased) English. In particular, it should be noted that while the models have similar architectures, they have some key differences. ALBERT introduces factorized embed-

ding parameterization and cross-layer parameter sharing, while the DistilBERT model is a student network trained with knowledge distillation from BERT. We use the weights and implementations of the models available in the Huggingface Transformers library (Wolf et al., 2019).

3.3 Datasets

We performed experiments on three standard datasets to assess the performance of the proposed approach on different tasks. All of them pose classification problems with a different number of target categories and dataset sizes. For each dataset, we use the test set specified by the source.

IMDB. The Large Movie Review Dataset (Maas et al., 2011) contains of a collection of 50,000 movie reviews and their associated binary sentiment polarity labels (either “positive” or “negative”), which is used to train a sentiment classifier.

Yelp F. This dataset (Zhang et al., 2015) contains reviews of local businesses and their associated rating (1-5). The task is posed as a text classification task, from the review text to its associated rating.

AG News. The AG’s corpus of news articles¹ consists of over one million news articles gathered from more than 2,000 news sources, divided into a number of categories depending on their content. We used the common subset (Zhang et al., 2015) of the whole dataset, consisting of a total of 120,000 samples equally divided in four categories.

3.4 Experiments and hyperparameters

We construct several experiments to evaluate how well Federated Learning scales to an exponentially increasing number of clients. In all experiments, the respective dataset is evenly partitioned into a number of subsets equal to the number of clients. Data points are uniformly sampled on each client (i.i.d.) like (McMahan et al., 2017a). We do not perform any hyperparameter tuning, and instead, keep all other hyperparameters constant for an unbiased comparison. As baselines we run for each task and BERT-variant a non-federated scenario (clients = 1) with the same configuration.

We run the baselines for a fixed number of rounds based on our compute budget. The test set performance for the baselines are then compared against varying number of participating clients at

¹http://groups.di.unipi.it/~gulli/AG_corpus_of_news_articles.html

the same number of rounds. Finally, since runs with a larger number of clients converge more slowly, we allow those runs to continue to a second threshold and report the number of rounds required to reach 90% of the baseline performance, similar to McMahan et al. (2017a). Runs not reaching 90% of the baseline performance within the second threshold are reported as failures.

We run the baseline for 100 rounds for both IMDB and AG News while setting the second threshold to 200 rounds. However, we only run Yelp F baselines for 50 rounds due to its large size and set the second threshold at 100 rounds. Like Lin et al. (2020), we avoid momentum, weight decay, and dynamic learning rates for simplicity. Instead, all experiments are performed with SGD. Based on Sun et al. (2019b) we choose the constant learning rate $2 \cdot 10^{-5}$, maximum sequence length 128 and batch size (B) of 32. Furthermore, the number of local epochs (E) is set to 2 per round.

4 Results

4.1 Fixed compute budget

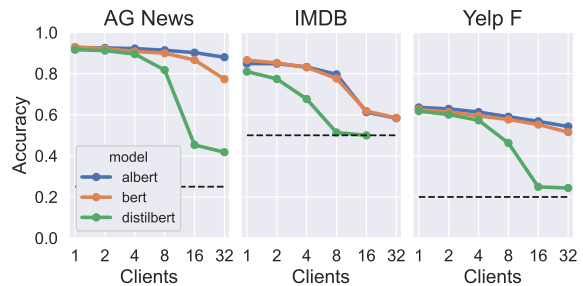


Figure 1: Accuracy at a fixed compute budget of 100 rounds for AG, IMDB, and 50 rounds for Yelp F. The expected accuracy of a random classifier for each task has been highlighted in the dashed line. Higher is better.

In Figure 1, we study the effect of increasing the number of clients. It shows the final accuracy after 100 rounds IMDB and AG News, and 50 rounds of the much larger Yelp F, with an exponentially increasing number of clients. Both ALBERT and BERT are well behaved and exhibit a gradual decrease with an increasing number of clients. However, DistilBERT shows a much steeper decline when moving past 4 clients for all datasets, down to the random classifier baseline (IMDB, Yelp F).

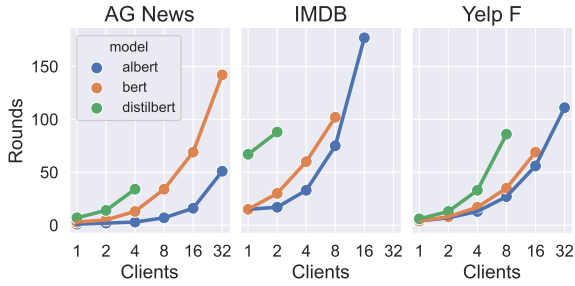


Figure 2: Number of training rounds required to reach 90% of the non-federated baseline accuracy. Omissions occur when the target is not reached in 100 (Yelp F) or 200 rounds (AG News, IMDB). Lower is better.

4.2 Rounds until target performance

Examining the number of rounds necessary to achieve 90% of the non-federated baseline accuracy (Figure 2) yields a similar observation. While all models perform worse with more clients, ALBERT and BERT mostly reach the target accuracy within the allocated number of rounds until 32 clients are used. DistilBERT on the other is unable to reach the target accuracy at 16 clients for Yelp F, and as low as 4 clients for IMDB.

4.3 Dynamics of fine-tuning

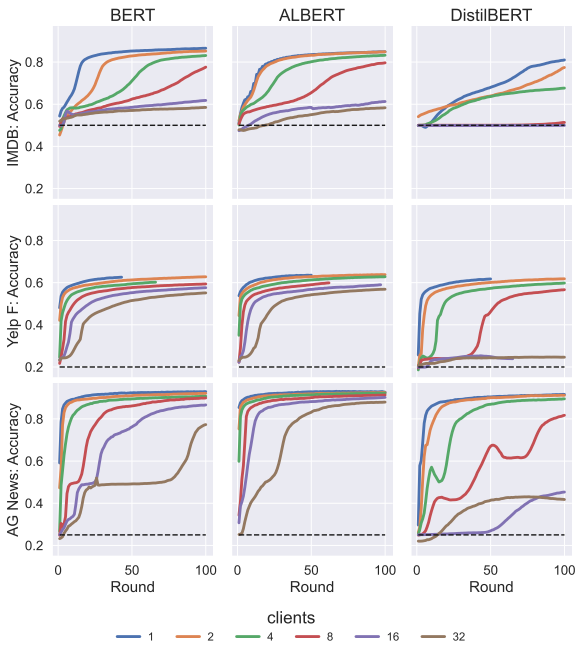


Figure 3: Test accuracy (higher is better) over communication rounds for our scenarios. The random classifier baseline is shown as a dashed line.

The test accuracy during fine-tuning (Figure 3) allows a more complete understanding of how well FEDAVG scales for language model fine-tuning.

While some scenarios (e.g. Yelp F. with BERT) show a gradual degradation with the number of clients, others configurations are more adversely affected by the increasing number of clients. In some instances the accuracy stays constant over a large period, sometimes even at the random classifier baseline for the whole (DistilBERT on IMDB) or part (DistilBERT on AG News) of the experiment when the number of clients is high.

5 Discussion

In this paper, we have evaluated the performance of Transformer-based language models fine-tuned in a federated setting. While BERT and ALBERT seem to learn each task quickly (Figure 3), DistilBERT has a much slower learning progression in the federated setup. A possible explanation is the process of distillation during pre-training, but further research is needed to fully understand why this happens. We demonstrated that BERT and ALBERT scale well up to 32 clients with no sharp decline in performance (Figure 1), but found DistilBERT to struggle at 16 clients in the Yelp F and AG News tasks, and with as low as 4 clients in the IMDB task, with a substantial drop in performance compared to the baseline. Furthermore, DistilBERT takes more rounds to achieve the same performance. These results demonstrate that we have obtained a higher communication efficiency for BERT and ALBERT as compared to DistilBERT. Further work is required to get a good picture of exactly what affects the communication efficiency for federated learning of Transformer-based language models.

Conversely, the sudden drop in performance in some scenarios indicates that FL can be sensitive to the number of clients. The cause for the instability has not been fully determined. It may be related to both smaller partitions and contradicting models. This highlights the importance of evaluating FL with a varying number of clients at these scales.

For all three models, the performance decline from the baseline is steeper for IMDB compared to the other tasks. This may be related to the variability in the movie review data, adding to a larger inter-client difference in data distribution when data is put into smaller partitions, resulting in a larger difference between the client models taking part in the federated averaging.

In conclusion, we have demonstrated the applicability of the federated learning paradigm and evaluated it on a number of Transformer-based models

up to 32 clients. Our findings show that the relatively large sizes of these models are not prohibitive for federated learning.

6 Acknowledgements

This work was funded by VINNOVA (grant 2019-05156). We would also like to thank AI Sweden and CGit for providing us with compute resources.

References

- Ashutosh Adhikari, Achyudh Ram, Raphael Tang, and Jimmy Lin. 2019. Docbert: Bert for document classification. *arXiv preprint arXiv:1904.08398*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT (1)*, pages 4171–4186.
- Suyu Ge, Fangzhao Wu, Chuhan Wu, Tao Qi, Yongfeng Huang, and Xing Xie. 2020. Fedner: Medical named entity recognition with federated learning. *arXiv preprint arXiv:2003.09288*.
- Neel Guha, Ameet Talwalkar, and Virginia Smith. 2019. One-shot federated learning. *arXiv preprint arXiv:1902.11175*.
- Andrew Hard, Kanishka Rao, Rajiv Mathews, Françoise Beaufays, Sean Augenstein, Hubert Eichner, Chloé Kiddon, and Daniel Ramage. 2018. Federated learning for mobile keyboard prediction. *CoRR*, abs/1811.03604.
- Tzu-Ming Harry Hsu, Hang Qi, and Matthew Brown. 2019. Measuring the effects of non-identical data distribution for federated visual classification. *arXiv preprint arXiv:1909.06335*.
- Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank J Reddi, Sebastian U Stich, and Ananda Theertha Suresh. 2019. Scaffold: Stochastic controlled averaging for on-device federated learning. *arXiv preprint arXiv:1910.06378*.
- Jakub Konečný, Brendan McMahan, and Daniel Ramage. 2015. Federated optimization: Distributed optimization beyond the datacenter. *arXiv preprint arXiv:1511.03575*.
- Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2020. Albert: A lite bert for self-supervised learning of language representations. In *International Conference on Learning Representations*.
- Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2020. Biobert: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36(4):1234–1240.
- Tian Li, Maziar Sanjabi, Ahmad Beirami, and Virginia Smith. 2019. Fair resource allocation in federated learning. *arXiv preprint arXiv:1905.10497*.
- Tao Lin, Lingjing Kong, Sebastian U Stich, and Martin Jaggi. 2020. Ensemble distillation for robust model fusion in federated learning. *arXiv preprint arXiv:2006.07242*.
- Edvin Listo Zec, Olof Mogren, John Martinsson, Leon René Sütfield, and Daniel Gillblad. 2020. Federated learning using a mixture of experts.
- Dianbo Liu and Tim Miller. 2020. Federated pretraining and fine tuning of bert using clinical notes from multiple silos.
- Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. 2011. Learning word vectors for sentiment analysis. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 142–150, Portland, Oregon, USA. Association for Computational Linguistics.
- Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. 2017a. Communication-Efficient Learning of Deep Networks from Decentralized Data. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, volume 54, pages 1273–1282, Fort Lauderdale, FL, USA. PMLR.
- H Brendan McMahan, Daniel Ramage, Kunal Talwar, and Li Zhang. 2017b. Learning differentially private recurrent language models. *arXiv preprint arXiv:1710.06963*.
- Mehryar Mohri, Gary Sivek, and Ananda Theertha Suresh. 2019. Agnostic federated learning. In *Proceedings of Machine Learning Research*, volume 97, pages 4615–4625, Long Beach, California, USA. PMLR.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. In *EMC² Workshop at NeurIPS 2019*.

- Chi Sun, Luyao Huang, and Xipeng Qiu. 2019a. Utilizing bert for aspect-based sentiment analysis via constructing auxiliary sentence. *arXiv preprint arXiv:1903.09588*.
- Chi Sun, Xipeng Qiu, Yige Xu, and Xuanjing Huang. 2019b. How to fine-tune bert for text classification? In *Chinese Computational Linguistics*, pages 194–206, Cham. Springer International Publishing.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2019. Huggingface’s transformers: State-of-the-art natural language processing. *ArXiv*, abs/1910.03771.
- Wei Yang, Yuqing Xie, Aileen Lin, Xingyu Li, Luchen Tan, Kun Xiong, Ming Li, and Jimmy Lin. 2019. End-to-end open-domain question answering with bertserini. *arXiv preprint arXiv:1902.01718*.
- Zexian Zeng, Yu Deng, Xiaoyu Li, Tristan Naumann, and Yuan Luo. 2018. Natural language processing for ehr-based computational phenotyping. *IEEE/ACM transactions on computational biology and bioinformatics*, 16(1):139–153.
- Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015. [Character-level convolutional networks for text classification](#). In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems 28*, pages 649–657. Curran Associates, Inc.
- Yue Zhao, Meng Li, Liangzhen Lai, Naveen Suda, Damon Civin, and Vikas Chandra. 2018. Federated learning with non-iid data. *arXiv preprint arXiv:1806.00582*.