

# ADAPTIVE LOCAL TRAINING IN FEDERATED LEARNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Federated Learning is a machine learning paradigm where multiple clients collaboratively train a global model by exchanging their locally trained model weights instead of raw data. In the standard setting, every client trains the local model for the same number of epochs. We introduce ALT (Adaptive Local Training), a simple yet effective feedback mechanism that could be introduced at the client side to limit unnecessary and degrading computations. ALT dynamically adjusts the number of training epochs for each client based on the similarity between their local representations and the global one, ensuring that well-aligned clients can train longer without experiencing client drift. We evaluated ALT on federated partitions of the CIFAR-10 and TinyImageNet datasets, demonstrating its effectiveness in improving model convergence and stability.

## 1 INTRODUCTION

Federated learning (FL) (McMahan et al., 2017) has emerged as a machine learning approach prioritizing privacy while fostering collaborative training, avoiding centralized data storage concerns.

Typically, in FL every client performs training for the same number of local epochs (Karimireddy et al., 2020; Li et al., 2021; Shenaj et al., 2023). However, different clients could have different computational capabilities, and communication speeds, and the server might request updates when the training is not yet concluded. In addition to that, when the clients' data distribution is very heterogeneous, training each client for a fixed pre-defined number of steps leads to client drift and complicates the aggregation.

For this reason, we train for a variable number of local epochs across clients and training rounds, similarly to (Li et al., 2020; Michieli et al., 2022). In particular, we study the effect of dynamic local epochs on the client side and propose a client-side control strategy to mitigate representation drift, reduce communication costs, and improve performances.

**Related Work.** FedProx (Li et al., 2020) introduces a flexible framework that allows clients to perform variable amounts of local training, assuming that some clients may not complete their updates within a fixed time window. Additionally, it employs a dynamic regularizer in the local objective to mitigate the impact of inconsistent local updates.

Similarly, SCAFFOLD (Karimireddy et al., 2020) tackles client drift by estimating and correcting update directions for both the server and clients, ensuring more stable local updates. However, while these methods offer improvements, they fail to achieve significant performance gains over FedAvg in

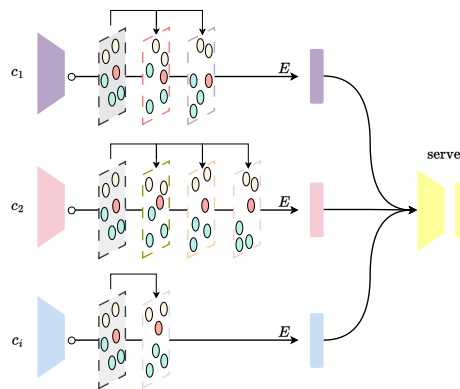
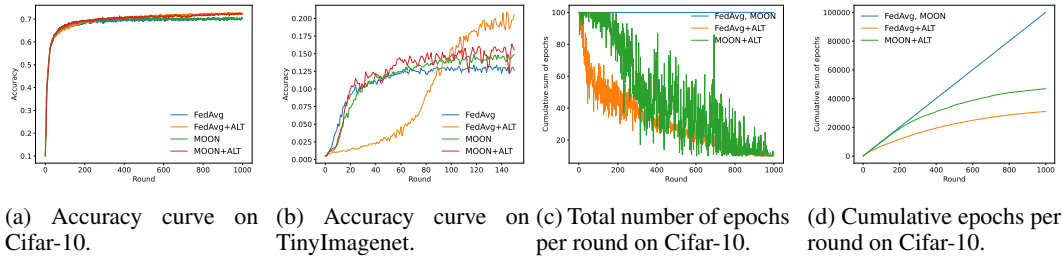


Figure 1: Overview of the proposed Federated learning strategy with dynamic local training epochs.



deep neural networks training, particularly when applied in realistic computer vision applications. To address this limitation, MOON (Li et al., 2021) proposes a model-based contrastive learning approach to enhance local training in non-IID settings. By enforcing similarity between the current local model and the incoming global model, while discouraging similarity with the previous local model, MOON effectively stabilizes training in deep networks.

Building on these insights, we introduce a simple yet effective mechanism to control the length of local training, opening new research opportunities in adaptive and dynamic federated learning.

## 2 METHOD

Let us assume a set of clients  $\mathcal{K}$ , where each client  $k \in \mathcal{K}$  has access to a local set of samples  $(\mathcal{X}_k, \mathcal{Y}_k)$  from a dataset  $\mathcal{D}_k$  with  $|\mathcal{D}_k| = n_k$ . At each communication round  $r \in \{1, \dots, R\}$ , the server selects a subset of clients  $S \subset \mathcal{K}$  and sets an adaptive threshold  $T_h(r) = a + \frac{b \cdot r}{R}$ , that linearly increases during the training (we set  $a = 0.1$  and  $b = 0.8$ , see Appendix B for more details). Each client  $s \in S$  then initializes its local model  $\theta_s^r$  with the global model  $\theta^r$  and trains for  $E_s^r$  local epochs, where  $E_s^r$  varies across clients and rounds.

At each training step, as usual in neural networks training, the local model is updated using gradient descent on mini-batches  $\mathcal{B} \subset \mathcal{D}_s$ , i.e.,  $\theta_s^r \leftarrow \theta_s^r - \eta \nabla \mathcal{L}_s(\theta_s^r; \mathcal{B})$ .

Let us denote with  $p_s = f(w_s, \mathcal{B})$  and  $p_g = f(w_g, \mathcal{B})$  the feature embeddings of the local and global models for the considered batch: the local training halts as soon as the similarity condition  $\cos(p_s, p_g) < T_h(r)$  is met, i.e., when the difference between the embeddings is smaller than the threshold. If the condition is never met, it will stop as usual after the maximum number of local epochs  $E$  is reached. Note that the threshold increases with time, i.e., the criteria becomes more and more strict while the difference w.r.t. the starting model typically increases.

After local training, the client models are sent to the server, which aggregates them using standard federated averaging, i.e.,:  $\theta^{r+1} = \sum_{s \in S} \frac{n_s}{n} \theta_s^r$ . The process is repeated for  $R$  rounds. The algorithm is detailed in Appendix C.

## 3 RESULTS

We evaluate the performance of our algorithm on the CIFAR-10 (Krizhevsky et al., 2009), and TinyImagenet (Le & Yang, 2015) datasets. As a strong baseline for local training, we consider MOON. We consider  $|\mathcal{K}| = 100$  clients (with 10% participation), and the data is partitioned according to the Dirichlet distribution with the concentration parameter  $\alpha = 100$ . By looking at the Figure, we can notice that our method allows to reduce substantially carbon footprint by reducing the cumulative epochs per round (sum of all clients epochs), and leads also to improved performance. In particular, we notice that FedAvg + ALT (FedALT) could work similarly or better then MOON, while being much more efficient.

## 4 CONCLUSION

In this work, we introduced a representation learning feedback mechanism to control the number of local epochs reducing energy consumption and communication costs which can be seamlessly integrated into FL algorithms.

108 REFERENCES

- 109
- 110 Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and
- 111 Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In
- 112 *International conference on machine learning*, pp. 5132–5143. PMLR, 2020.
- 113
- 114 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
- 115 2009.
- 116
- 117 Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3, 2015.
- 118
- 119 Qinbin Li, Bingsheng He, and Dawn Song. Model-contrastive federated learning. In *Proceedings of*
- 120 *the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10713–10722, 2021.
- 121
- 122 Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith.
- 123 Federated optimization in heterogeneous networks. *Proceedings of Machine learning and sys-*
- 124 *tems*, 2:429–450, 2020.
- 125
- 126 Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas.
- 127 Communication-efficient learning of deep networks from decentralized data. In *Artificial Intelli-*
- 128 *gence and Statistics (AISTATS)*, pp. 1273–1282. PMLR, 2017.
- 129
- 130 Umberto Michieli, Marco Toldo, and Mete Ozay. Federated learning via attentive margin of semantic
- 131 feature representations. *IEEE Internet of Things Journal*, 10(2):1517–1535, 2022.
- 132
- 133 Donald Shenaj, Giulia Rizzoli, and Pietro Zanuttigh. Federated learning in computer vision. *IEEE*
- 134 *Access*, 2023.
- 135

136

137 APPENDIX

138 A IMPLEMENTATION DETAILS

139

140 We adopt as feature extractor  $f$  a simple CNN with 2 convolutional layers for CIFAR-10, and

141 ResNet-50 for Tiny-Imagenet. We use the SGD optimizer with a learning rate 0.01 for all ap-

142 proaches. The SGD weight decay is set to 0.00001 and SGD momentum is set to 0.9. The batch size

143 is set to 64. The maximum number of local epochs  $E$  is set to 10, and the number of communication

144 rounds  $R$  is 1000 for Cifar-10 and 150 for TinyImagenet.

145

146

147

148

149 B SELECTION OF THE THRESHOLD FUNCTION

150

151 We experimented different threshold functions  $T_h(r)$ , including 3 simple approaches:

152

- 153 • Linear Increasing:  $T_h(r) = 0.1 + 0.8 \cdot \frac{r}{R}$
  - 154 • Linear Decreasing:  $T_h(r) = 0.9 - 0.8 \cdot \frac{r}{R}$
  - 155 • Fixed:  $T_h(r) = c$ , where  $c$  is a constant.
- 156
- 157

158 Among these strategies, the most effective was the linear increasing threshold. This approach en-

159 ables a ‘slow start,’ where early training rounds allow for more flexibility in local updates before

160 gradually enforcing stricter similarity constraints. This progressive tightening helps balance explo-

161 ration and convergence, leading to improved overall model performance. However, in practice, any

thresholding (including a fixed rule) is useful to reduce carbon footprint.

## C ALGORITHM

```

162
163
164
165 1: Input: Initialize model parameters  $\theta$ 
166 2: Initialize maximum rounds  $R$ 
167 3: Initialize threshold parameters  $a, b$ 
168 4: Server executes:
169 5: for  $r = 1$  to  $R$  do
170 6:  $S \leftarrow$  Random subset of clients
171 7:  $T_h(r) \leftarrow a + \frac{b * r}{R}$ 
172 8: for each client  $i \in S$  do
173 9:  $\theta_i \leftarrow$  ClientUpdate( $i, \theta, r, T_h(r)$ )
174 10: end for
175 11:  $\theta \leftarrow \sum_{i \in S} \frac{n_i}{n} \theta_i$ 
176 12: end for
177 13: ClientUpdate( $i, \theta, r, T_h(r)$ ):
178 14: if  $r = 1$ :
179 15:  $\theta_i \leftarrow \theta$ 
180 16:  $\theta_i := \{w_i, v_i\}$ 
181 17:  $\theta_g \leftarrow \theta$ 
182 18:  $\theta_g := \{w_g, v_g\}$ 
183 19: stop  $\leftarrow False$ 
184 20: for  $j = 1, 2, \dots, E$  and stop = False do
185 21: for each batch  $\mathcal{B}$  in  $\mathcal{D}_i$  do
186 22:  $p_i \leftarrow f(w_i, \mathcal{B})$ 
187 23:  $p_g \leftarrow f(w_g, \mathcal{B})$ 
188 24: if  $\cos(p_i, p_g) < T_h(r)$ : stop  $\leftarrow True$ 
189 25:  $\theta_g \leftarrow \theta_g - \eta \nabla \mathcal{L}(\theta_g; \mathcal{B})$ 
190 26:  $\theta_i \leftarrow \theta_g$ 
191 27: return  $\theta_i$ 
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

```

**Algorithm 1:** Implementation of FedAvg with the ALT stopping criteria (FedALT)