

Fast Federated Edge Learning with Overlapped Communication and Computation and Channel-Aware Fair Client Scheduling

Mehmet Emre Ozfatura[†], Junlin Zhao^{*}, and Deniz Gündüz[†]

[†]Department of Electrical and Electronic Engineering, Imperial College London, London, UK

^{*}School of Science and Engineering, the Chinese University of Hong Kong (Shenzhen), Shenzhen, China

Email: m.ozfatura@imperial.ac.uk, zhaojunlin@cuhk.edu.cn, d.gunduz@imperial.ac.uk

Abstract—We consider federated edge learning (FEEL) over wireless fading channels taking into account the downlink and uplink channel latencies, and the random computation delays at the clients. We speed up the training process by overlapping the communication with computation. With fountain coded transmission of the global model update, clients receive the global model asynchronously, and start performing local computations right away. Then, we propose a dynamic client scheduling policy, called MRTP, for uploading local model updates to the parameter server (PS), which, at any time, schedules the client with the minimum remaining upload time. However, MRTP can lead to biased participation of clients in the update process, resulting in performance degradation in non-iid data scenarios. To overcome this, we propose two alternative schemes with fairness considerations, termed as age-aware MRTP (A-MRTP), and opportunistically fair MRTP (OF-MRTP). In A-MRTP, the remaining clients are scheduled according to the ratio between their remaining transmission time and the update age, while in OF-MRTP, the selection mechanism utilizes the long term average channel rate of the clients to further reduce the latency while ensuring fair participation of the clients. It is shown through numerical simulations that OF-MRTP provides significant reduction in latency without sacrificing test accuracy.

Index Terms—Client selection, fair scheduling, federated edge learning

I. INTRODUCTION

Extensive research efforts have been devoted to overcome the communication bottleneck in federated learning (FL) [1]. Lossy compression techniques, including quantization [2]–[5] and sparsification [6]–[9] have been developed to reduce the communication cost. However, these approaches treat the communication channel connecting the clients to the parameter server (PS) as an error-free bit pipe, ignoring the wireless channel characteristics. However, efficient implementation of FL at the wireless edge, called *federated edge learning (FEEL)*, requires jointly optimizing the learning framework with the underlying communication framework taking into account channel characteristics and constraints [10], [11].

As common in wireless networks, different clients may have distinct computational capabilities, channel statistics, power budgets, etc., resulting in the *straggler effect* in FEEL. By carefully designing the resource allocation policy, the balance

between the training latency and energy consumption has been studied in [12]–[14]. By scheduling only a limited number of clients in each round in FEEL, the communication requirement can be reduced, and the straggler effect can also be alleviated [15]–[18]. In [15], the authors propose to maximize the number of scheduled clients in each round to speed up training. In [16], a function of the ages of client updates is minimized to schedule all the clients as often as possible. Convergence rates of three different scheduling policies, namely, random, round-robin, proportional fair scheduling, are analyzed in [17]. Different from the aforementioned works, [18] considers ‘update-aware’ client scheduling in FEEL, which takes the significance of model updates into account together with their channel states. Alternatively, in [19] clients are clustered around several access points to perform FL hierarchically in order to reduce the communication latency. Another approach to addressing the communication bottleneck is aggregating the local model updates using the superposition property of the wireless multiple access channel [20]–[22]. However, this approach requires accurate synchronization among the transmitting agents [23].

Unlike the previous literature on FEEL, we consider both the downlink and uplink channel variations, together with random computation delays. With the exception of [24], works on FEEL ignore the downlink communication channel and the associated latency. Our first contribution is to overlap the downlink and uplink transmissions with local computations at the clients. This is achieved by fountain-coded delivery of the global model update [25], such that clients with better downlink channel conditions can receive the global model quickly, and start computing right away. Then, the clients are scheduled for the uplink delivery of their model updates to the PS as soon as they complete their local computations. To further minimize the latency, we schedule the client with the minimal uplink latency at any point in time.

While this approach, called MRTP, minimizes the latency, it may result in a loss in test accuracy as clients close to the PS would dominate the training process thanks to their statistically better channel conditions. We mitigate this bias by introducing both short-term and long-term fairness condi-

tions. In particular, we utilize the *age* and *frequency* metrics, respectively, for short- and long-term fairness, where age, in a broad sense, measures the staleness of the client model update, and frequency measures the long term participation statistics of the clients. Our numerical results show that the proposed overlapped and fair scheduling policy significantly speeds up FEEL without sacrificing the final test accuracy.

II. SYSTEM MODEL

A. FL model

Consider K clients collaboratively training a model parameter vector of dimension d , $\theta \in \mathbb{R}^d$, with periodic communication with a PS to minimize the empirical loss function, $F(\theta) = 1/K \sum_{k=1}^K F_k(\theta)$, where $F_k(\theta)$ is the loss function of client k . Let \mathcal{D}_k denote the dataset at client k with $D_k \triangleq |\mathcal{D}_k|$ samples. The empirical loss function at client k is $F_k(\theta) = 1/D_k \sum_{\mathbf{u} \in \mathcal{D}_k} f(\theta, \mathbf{u})$, where $f(\theta, \mathbf{u})$ is the task-dependent loss function measured with the model parameter vector θ and data \mathbf{u} . At each communication round $n \in \{1, 2, \dots\}$, each participant client performs τ -step stochastic gradient descent (SGD) on its local dataset to minimize the loss function based on the received global model parameter $\theta(n)$. At the m -th step of the local SGD of participant client k in communication round n , the local model is updated as

$$\theta_k^{m+1}(n) = \theta_k^m(n) - \eta_k^m(n) \nabla F_k(\theta_k^m(n), \xi_k^m(n)), \quad (1)$$

where $\eta_k^m(n)$ is the learning rate, and $\nabla F_k(\theta_k^m(n), \xi_k^m(n))$ is the gradient computed with the local model parameter $\theta_k^m(n)$ and mini-batch data $\xi_k^m(n)$, which is uniformly randomly selected from \mathcal{D}_k and follows $\mathbb{E}_{\xi} \{\nabla F_k(\theta_k^m(n), \xi_k^m(n))\} = \nabla F_k(\theta_k^m(n))$.

Each participant client then forwards its updated local model to the PS. Denoting the model vector at client k updated after τ -steps by $\theta_k(n)$, $k = 1, \dots, K$, the global model is updated by the PS as $\theta(n+1) = \frac{1}{K} \sum_{k=1}^K \theta_k(n)$.

B. Communication Model

A block fading channel model is assumed, where the channels between the PS and the clients remain unchanged in each communication round of the FEEL process.

1) *Asynchronous global model transmission*: Note that we have a multicast channel when transmitting the global model from the PS to the clients. To speed up the training process, we use fountain coded multicasting of the global model to the clients [25]. The downlink rate at client k in communication round n is given by $R_k^{\text{dl}}(n) = \log_2 \left(1 + \frac{P|h_k^{\text{dl}}(n)|^2}{\sigma_k^2} \right)$, $\forall k$, where $h_k^{\text{dl}}(n)$ is the complex-valued downlink channel coefficient between the PS and client k in round n , P is the transmit power at the PS, and σ_k^2 is the noise variance.

With asynchronous global model transmission in the downlink, a client recovers the global model with a latency dependent on its channel gain, and immediately starts local computations. Assuming that the global model is compressed into Q bits, it will take $Q/R_k^{\text{dl}}(n)$ seconds for client k to receive the model update. This will allow us to parallelize global model

transmission and computations, instead of targeting the worst client to guarantee all the clients receive the global model simultaneously.

In the rest of the paper, to simplify notation we will drop the round index n when it is clear from the context.

2) *Local model update*: In the uplink, where clients upload their model parameters to the PS, we assume a time-division framework; that is, only a single client is scheduled at any point in time. The instantaneous rate of each client is given by $R_k^{\text{ul}}(n) = \log_2 \left(1 + P_k |h_k^{\text{ul}}(n)|^2 / \sigma_0^2 \right)$, where P_k is the transmit power at client k , and $h_k^{\text{ul}}(n)$ is the uplink channel coefficient from client k to the PS.

C. Local computations

We assume that the computation speeds at the clients are also random following a distribution that is independent across clients and rounds. In each round, we denote by $\pi(j)$ the client with the j -th smallest accumulated latency in downlink transmission and local computation, $j = 1, \dots, K$, and by T_j the corresponding latency since the beginning of the communication round, with $T_0 \triangleq 0$.

D. Client scheduling

With asynchronous downlink transmission and heterogeneous local computation speeds, clients complete their local model updates in a sequential manner. The clients that have completed local computations and are thus available to upload their local models to the PS are referred to as *idle* clients. We denote the set of idle clients at time t within each round by $\mathcal{C}^{\text{idle}}(t)$. Hence, client $\pi(j)$ is added to the idle set at time T_j , and remains there until it uploads its model to the PS.

Let $s(t) \in \{0, 1, \dots, K\}$ denote the index of the client scheduled for transmission at time t , where $s(t) = 0$ means no client is scheduled. We have $s(t) \neq \pi(j)$ if $t < T_j$; that is, a client cannot be scheduled for upload before it completes its local computations. Let $Q_k(t)$ denote the remaining size of model parameter vector (measured in bits) at time t that has not yet been uploaded to the PS. We have $Q_{\pi(k)}(T_k) = Q$, and

$$Q_k(t) = Q - R_k^{\text{ul}} \int_0^t \mathbb{1}_{\{s(t)=k\}} dt, \quad (2)$$

where $\mathbb{1}_{\{x\}}$ is the indicator function, which is 1 when x holds, and 0 otherwise.

Let t_k denote the time client k completes uploading its model to the PS, i.e., $t_k \triangleq \min_t \{Q_k(t) = 0\}$. Client k is removed from the idle set at time t_k . Each round continues until N out of K clients upload their models to the PS, and therefore, some clients may never be added to the idle set, or may not leave the set at the end of the round. We have $t_k = \infty$ for those clients. Let $\mathcal{K}(t)$ denote the set of clients that have completed their upload by time t within round n , i.e., $\mathcal{K}(t) = \{k : t_k \leq t\}$. For a specific scheduling policy, the completion time of round n is given by $T(n) \triangleq \min_t \{|\mathcal{K}(t)| = N\}$, while the set of clients scheduled in round n are given, with slight abuse of notation, by $\mathcal{K}_n \triangleq \mathcal{K}(T(n))$.

Algorithm 1 MRTP in round n

Input: $K, N, Q, \{R_k^{\text{ul}}(n)\}_{k=1}^K$
Output: $\mathcal{K}_n, T(n)$

- 1: Initialization: $t = 0, \mathcal{C}^{\text{idle}}(t) = \emptyset, \mathcal{K}_n = \emptyset, Q_k(t) = Q$
- 2: **while** $|\mathcal{K}_n| < N$ **do**
- 3: **if** $\mathcal{C}^{\text{idle}}(t) \neq \emptyset$ **then**
- 4: Schedule $k^* = \arg \min \frac{Q_k(t)}{R_k^{\text{ul}}}$ for $k \in \mathcal{C}^{\text{idle}}(t)$
- 5: **while** No new arrival and $Q_{k^*}(t) > 0$ **do**
- 6: Schedule client k^* , update $t, Q_{k^*}(t)$ as in (2)
- 7: **if** new arrival k **then**
- 8: Update $\mathcal{C}^{\text{idle}}(t) \leftarrow \mathcal{C}^{\text{idle}}(t) \cup \{k\}$
- 9: **else**
- 10: Update $\mathcal{C}^{\text{idle}}(t) \leftarrow \mathcal{C}^{\text{idle}}(t) \setminus \{k^*\}$
- 11: Update $\mathcal{K}(t) \leftarrow \mathcal{K}(t) \cup \{k^*\}$
- 12: **else**
- 13: Wait until arrival of another client
- 14: **if** new arrival k **then**
- 15: $t \leftarrow T_k$
- 16: Update $\mathcal{C}^{\text{idle}}(t) \leftarrow \mathcal{C}^{\text{idle}}(t) \cup \{k\}$

III. SCHEDULING POLICIES

Our goal is to come up with scheduling policies that would not only minimize the completion time of each round, but also lead to the fastest convergence in terms of the wall-clock time.

A. Minimum remaining time-based policy (MRTP)

Note that, at any time instant within a round, we can schedule any client from the idle set. However, it is easy to see that there is no loss of optimality making scheduling decisions only when the idle set is updated, i.e., when a new client becomes idle, or one of the idle clients completes uploading its model.

In MRTP, each time the idle set is updated, we schedule the client with the minimum remaining time to upload its update. Specifically, at time t , client k^* is scheduled if

$$k^* = \arg \min_{k \in \mathcal{C}^{\text{idle}}(t)} \frac{Q_k(t)}{R_k^{\text{ul}}}. \quad (3)$$

Details of the MRTP can be found in Algorithm 1.

B. Age-aware MRTP (A-MRTP)

While MRTP minimizes the upload time, clients with statistically better channel qualities are more likely to be scheduled, which results in non-uniform sampling of clients and overfitting due to excessive use of limited amount of data. This may drastically degrade the performance of FEEL, especially when the data is not i.i.d., which is common in practice. Therefore, to strike a balance between the latency and model accuracy, we propose two alternative schemes by taking the ‘age of update’ into consideration.

For the short term fairness, we utilize the *age* metric, where the age of a client at round n , denoted by $a_{k,n}$, represents the number of rounds since the last time it was scheduled. The age parameter evolves as follows:

$$a_{k,n+1} = \begin{cases} a_{k,n} + 1, & \text{if } k \notin \mathcal{K}_n \\ 1, & \text{if } k \in \mathcal{K}_n \end{cases}. \quad (4)$$

In A-MRTP, αN clients are firstly selected to upload their models as in MRTP to minimize the latency, where $\alpha \in [0, 1]$ is a tuning parameter. Then, to promote selecting clients that are less frequently scheduled, we select the client with the minimum ratio between the remaining time and the age of its update, which can be considered as the average latency for each timely update in the short term. Therefore, a balance between the efficiency in training and fairness is achieved by tuning α .

C. Opportunistic Fair MRTP (OF-MRTP)

The main drawback of A-MRTP is that it utilizes only instantaneous rate $R_k^{\text{ul}}(n)$ for client scheduling. However, when the clients are located at different distances from the PS, their participation frequency will still depend on their locations. We propose an opportunistic policy that utilizes the relative channel condition, denoted by $\gamma_{k,n}$, which measures the ratio of instantaneous channel rate of client k to its long term average value, \bar{R}_k^{ul} ; that is, we have $\gamma_{k,n} = R_k^{\text{ul}}(n)/\bar{R}_k^{\text{ul}}$. Hence, instead of scheduling clients based on their instantaneous channel states, we use $\gamma_{k,n}$ for scheduling. Further, we consider two metrics to ensure both short-term and long-term fairness among the clients. We use the *age metric* for short-term fairness, that is, to promote uniform participation of the clients, and define a *frequency metric* for long-term fairness, that is, to promote equal participation of clients. The frequency metric, $f_{k,n}$, denotes the participation frequency of the k th client at round n . We define $f_{k,n} \triangleq l_k(n)/(n-1)$, where $l_k(n)$ denotes the total number of rounds that client k has participated until round n .

In order to introduce long-term fairness, we consider a subset of the idle clients as follows

$$\tilde{\mathcal{C}}^{\text{idle}}(t) = \{k : k \in \mathcal{C}^{\text{idle}}(t), f_{k,n} < f_{\max}\}, \quad (5)$$

where f_{\max} is a maximum frequency constraint. For the opportunistic policy, we further consider the following subset of clients

$$\hat{\mathcal{C}}^{\text{idle}}(t) \triangleq \{k : k \in \tilde{\mathcal{C}}^{\text{idle}}(t), a_{k,n} > a_{th}, \gamma_{k,n} > \gamma_{\min}\}, \quad (6)$$

where a_{th} is the minimum age constraint introduced to ensure short-term fairness and γ_{\min} is the rate constraint for opportunistic scheduling. If there are multiple clients in $\hat{\mathcal{C}}^{\text{idle}}(t)$, then client k with the maximum $\gamma_{k,n}$ value, that is the one with the best relative channel condition is scheduled.

Similarly to A-MRTP, the proposed opportunistic policy consists of two steps. In the initial step, α portion of the clients are scheduled from $\tilde{\mathcal{C}}^{\text{idle}}(t)$ according to MRTP, while the remaining clients are scheduled from $\hat{\mathcal{C}}^{\text{idle}}(t)$ based on $\gamma_{k,n}$. However, if $\hat{\mathcal{C}}^{\text{idle}}(t) = \emptyset$, then client is scheduled according to MRTP from $\tilde{\mathcal{C}}^{\text{idle}}(t)$. As shown in (5), clients with excessive participation frequency are excluded from scheduling, which increases the participation of less frequently selected clients, and thus, improves the fairness in scheduling. Overall, the proposed opportunistic scheduling strategy OF-MRTP is defined by four system parameters $\alpha, a_{th}, \gamma_{\min}$, and f_{\max} .

Age threshold	MRTP fraction	frequency constraint	Accuracy (mean \pm std)	Average per round latency
$a_{th} = 10$	$\alpha = 0.25$	$f_{max} = 0.3$	83.857 ± 0.43	84.81 ± 3.41 seconds
$a_{th} = 10$	$\alpha = 0.25$	$f_{max} = 0.4$	83.405 ± 0.29	86.6 ± 3.13 seconds
$a_{th} = 10$	$\alpha = 0.5$	$f_{max} = 0.3$	83.65 ± 0.49	84.49 ± 3.78 seconds
$a_{th} = 10$	$\alpha = 0.5$	$f_{max} = 0.4$	83.157 ± 0.61	86.33 ± 4.77 seconds
$a_{th} = 5$	$\alpha = 0.25$	$f_{max} = 0.3$	84.10 ± 0.43	141.35 ± 9.27 seconds
$a_{th} = 5$	$\alpha = 0.25$	$f_{max} = 0.4$	83.734 ± 0.4	138.05 ± 7.77 seconds
$a_{th} = 5$	$\alpha = 0.5$	$f_{max} = 0.3$	83.84 ± 0.39	87.25 ± 5.65 seconds
$a_{th} = 5$	$\alpha = 0.5$	$f_{max} = 0.4$	84.02 ± 0.3	70.04 ± 4.17 seconds

TABLE I: Test accuracy and average per round latency of OF-MRTP for 5000 rounds, averaged over 10 trials.

IV. NUMERICAL RESULTS

A. Simulation Setup

1) *Objective and network setup*: We consider image classification on the CIFAR-10 dataset, which contains 50,000 training and 10,000 test images from 10 classes. We employ a convolutional neural network (CNN) architecture consisting of 4 convolutional layers followed by 4 fully connected layers. We set $\tau = 4$ local iterations using batchsize of 32.

We consider 100 client devices randomly distributed around the PS, and assume that the training dataset is distributed disjointly among the client devices in a non-iid manner, such that each client has 500 distinct training images from at most 4 different classes. Finally, we set the participation ratio to 20%, which means, at each round, $N = 20$ devices, out of 100, are scheduled to upload their model to the PS.

2) *Computation latency*: To model the computation latency at the clients, we consider the commonly employed shifted exponential distribution [26], where the probability of completing τ local updates by time t is given by

$$I(t) = \begin{cases} 1 - e^{-\mu(\frac{t}{\tau} - T_{min})}, & \text{if } t \geq \tau T_{min}, \\ 0, & \text{otherwise.} \end{cases} \quad (7)$$

where T_{min} is the minimum computation latency to perform one local update and μ is the average delay for one local update that is $\mu = \bar{T} - T_{min}$, where \bar{T} is the mean computation time for one local update. In order to obtain practically relevant values for μ and T_{min} , we measured the time for computation on a CPU using `time.time()` command, according to given batchsize and the network model, for over 100000 trials.

3) *Path loss and noise*: K clients are uniformly randomly distributed in the region within 500 meters around the PS. The path loss of client k is given by $PL_k = 148.1 + 37.6 \log_{10} d_k$, where d_k is the distance of k to the PS measured in kilometers. We set $P = 15$ dBm, and $P_k = 10$ dBm, $\forall k$. The channels of clients across different time slots are modeled as i.i.d. fading. The standard deviation of the channel gain of client k is $\phi_k = (10^{-PL_k/10})^{1/2}$, $\forall k$, and the channel coefficient is

MRTP fraction	Accuracy (mean \pm std)	Average latency
$\alpha = 0.9$	80.67 ± 0.73	315.14 ± 151.4 seconds
$\alpha = 0.7$	82.01 ± 0.375	414.5 ± 164.6 seconds

TABLE II: Test accuracy and average per round latency of A-MRTP for 5000 rounds, averaged over 10 trials.

modeled as $h_k(n) = \phi_k \tilde{h}_k(n)$, where $\tilde{h}_k(n)$ denotes an i.i.d. random variable accounting for Rayleigh fading of unit power in communication round n . The variances of noise at the PS and the clients are set as $\sigma_0^2 = \sigma_k^2 = 7.96 \times 10^{-14}$ Watts.

B. Simulation Results

We first consider the A-MRTP scheme. For the experiments, we consider $\alpha = \{0.7, 0.9\}$. The average per round latency and test accuracy results are presented in Table II. As predicted, we observe that by decreasing α ; that is, scheduling more clients according to the age-metric, both the test accuracy and the average latency increase.

We then consider OF-MRTP by setting $\alpha = \{0.25, 0.5\}$, $a_{th} = \{5, 10\}$, $f_{max} = \{0.3, 0.4\}$, and $\gamma_{min} = 1$. The test accuracy and average per round latency results are presented in Table I. We note that, if a round-robin scheduler is employed with participation ratio 20%, then the maximum age will be $a_{t,k} = 5$, and similarly, the maximum participation frequency will be $f_{k,t} = 0.2$. Therefore, to setup the parameter values a_{th} and f_{max} , we consider these values as our reference points. One can easily observe from Table I that with OF-MRTP the average latency as well as the variation on the latency significantly drops, since client k is scheduled only if $R_k^u(n) > \bar{R}_k^u(n)$. Besides, thanks to the control on the participation frequency of the clients, we also observe a significant improvement in the test accuracy. We also want to remark that when $a_{th} = 10$, we observe similar test accuracy and latency results for other parameters. The reason is that since MRTP is used for client scheduling when $\hat{C}^{idle}(t) = \emptyset$, at the same time larger a_{th} increases the probability of $\hat{C}^{idle}(t)$ being empty, thus more clients are scheduled according to MRTP. This observation is also backed by the simulation results with $a_{th} = 5$, where the impact of the parameter α is more visible. Not surprisingly the minimum latency is achieved when more clients are scheduled according to MRTP and we allow larger participation frequencies. Interestingly, in this case, where $\alpha = 0.5$ and $f_{max} = 0.4$, we do not observe a compromise on the test accuracy.

For comparison, we consider MRTP and random scheduling as benchmark. Note that these are the optimal strategies, respectively, from the latency and fairness perspectives. In Fig. 1, we compare the convergence behaviour of the OF-MRTP, A-MRTP and MRTP schemes. As one would expect, the test accuracy increases faster with MRTP; however, due to the non-iid distribution of the data it converges to a sub-optimal model

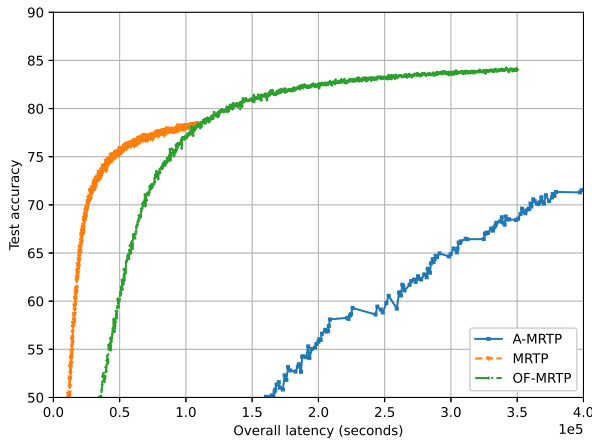


Fig. 1: Comparison of the MRTP, A-MRTP and OF-MRTP

and eventually diverges¹. While A-MRTP eventually reaches an accuracy level above MRTP (see Table II), it introduces significant latency due to scheduling the clients with poor channel conditions to reduce their age. The convergence behaviour of random scheduling is not included in the figure since the average per-round latency is very high, instead, we compare OF-MRTP and random scheduling based on the final test accuracy results which is averaged over 10 trails. We observe that the average final test accuracy with random scheduling is 83.92 ± 0.37 . In comparison with the test accuracy results in Table I, we can conclude that OF-MRTP does not compromise the accuracy while significantly reducing the latency.

V. CONCLUSION

We proposed novel global model transmission and client scheduling techniques to speed up wall-clock training time for FEEL without sacrificing final test accuracy. In particular, we ensured fair participation of the clients to achieve high test accuracy, and reduced the overall latency, which includes the computation time and model uplink/downlink latencies. To this end, we first introduced a fountain coded asynchronous model downlink strategy to allow clients to start local computations without waiting for others to download the global model. We then introduced MRTP, which adaptively schedules the client that can upload its local model to the PS in the fastest manner. MRTP and asynchronous downlink strategy, together, help to overlap computation and communication time, thus reduce the overall latency. However, as we experimentally show, client selection that solely focuses on the latency may lead to divergence when certain clients participate in the model update more frequently than others. Hence, we further employed the ‘update age’ and ‘update frequency’ as fairness metrics, which are opportunistically used to speed up training without sacrificing accuracy. Finally, through extensive simulations we show that it is possible to significantly reduce the overall latency without compromising the test accuracy.

¹The test accuracy is plotted until the divergence point.

REFERENCES

- [1] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, “Communication-efficient learning of deep networks from decentralized data,” in *Artificial Intel, and Stat.* PMLR, 2017, pp. 1273–1282.
- [2] F. Seide, H. Fu, J. Droppo, G. Li, and D. Yu, “1-bit stochastic gradient descent and its application to data-parallel distributed training of speech DNNs,” in *Conf. of Int’l Speech Comm. Assoc.*, 2014.
- [3] D. Alistarh, D. Grubic, J. Li, R. Tomioka, and M. Vojnovic, “QSGD: Communication-efficient SGD via gradient quantization and encoding,” *arXiv preprint arXiv:1610.02132*, 2016.
- [4] W. Wen, C. Xu, F. Yan, C. Wu, Y. Wang, Y. Chen, and H. Li, “Terngrad: Ternary gradients to reduce communication in distributed deep learning,” *arXiv preprint arXiv:1705.07878*, 2017.
- [5] S. Zhou, Y. Wu, Z. Ni, X. Zhou, H. Wen, and Y. Zou, “Dorefa-net: Training low bitwidth convolutional neural networks with low bitwidth gradients,” *arXiv preprint arXiv:1606.06160*, 2016.
- [6] N. Strom, “Scalable distributed DNN training using commodity GPU cloud computing,” in *Annual Conf. of Int’l Speech Comm. Assoc.*, 2015.
- [7] N. Dryden, T. Moon, S. Jacobs, and B. Van Essen, “Communication quantization for data-parallel training of deep neural networks,” in *Wrkshp. on Mach. Learn. in HPC Environ. (MLHPC)*, 2016, pp. 1–8.
- [8] A. F. Aji and K. Heafield, “Sparse communication for distributed gradient descent,” *arXiv preprint arXiv:1704.05021*, 2017.
- [9] H. Wang, S. Sievert, Z. Charles, S. Liu, S. Wright, and D. Papailiopoulos, “Atom: Communication-efficient learning via atomic sparsification,” *arXiv preprint arXiv:1806.04090*, 2018.
- [10] D. Gunduz, D. B. Kurka, M. Jankowski, M. M. Amiri, E. Ozfatura, and S. Sreekumar, “Communicate to learn at the edge,” *IEEE Commun. Mag.*, vol. 58, no. 12, pp. 14–19, 2020.
- [11] M. Chen, D. Gunduz, K. Huang, W. Saad, M. Bennis, A. V. Feljan, and H. V. Poor, “Distributed learning in wireless networks: Recent progress and future challenges,” *arXiv cs.LG:2104.02151*, 2021.
- [12] Q. Zeng, Y. Du, K. Huang, and K. K. Leung, “Energy-efficient radio resource allocation for federated edge learning,” in *IEEE Int’l Conf. on Comm. Workshops*, 2020, pp. 1–6.
- [13] Z. Yang, M. Chen, W. Saad, C. S. Hong, and M. Shikh-Bahaei, “Energy Efficient Federated Learning Over Wireless Communication Networks,” *arXiv*, pp. 1–30, 2019.
- [14] M. Chen, H. V. Poor, W. Saad, and S. Cui, “Convergence Time Optimization for Federated Learning over Wireless Networks,” *IEEE Trans. Wirel. Commun.*, pp. 1–30, 2020.
- [15] T. Nishio and R. Yonetani, “Client selection for federated learning with heterogeneous resources in mobile edge,” in *IEEE Int’l Conf. on Communications (ICC)*, 2019, pp. 1–7.
- [16] H. Yang, A. Arafa, T. Quek, and H. V. Poor, “Age-Based Scheduling Policy for Federated Learning in Mobile Edge Networks,” *IEEE Int. Conf. Acoust. Speech Signal Proc. (ICASSP)*, pp. 8743–8747, 2020.
- [17] H. H. Yang, Z. Liu, T. Q. Quek, and H. V. Poor, “Scheduling Policies for Federated Learning in Wireless Networks,” *IEEE Trans. Commun.*, vol. 68, no. 1, pp. 317–333, 2020.
- [18] M. M. Amiri, D. Gündüz, S. R. Kulkarni, and H. Vincent Poor, “Convergence of update aware device scheduling for federated learning at the wireless edge,” *IEEE Trans. Wireless Comm.*, pp. 1–1, 2021.
- [19] M. S. H. Abad, E. Ozfatura, D. Gunduz, and O. Ercetin, “Hierarchical federated learning across heterogeneous cellular networks,” in *IEEE Int’l Conf. Acous., Speech and Sig. Proc. (ICASSP)*, 2020, pp. 8866–8870.
- [20] M. M. Amiri and D. Gunduz, “Machine Learning at the Wireless Edge: Distributed Stochastic Gradient Descent Over-the-Air,” *IEEE Trans. Signal Proc.*, vol. 68, no. 4, pp. 2155–2169, 2020.
- [21] G. Zhu, Y. Wang, and K. Huang, “Broadband Analog Aggregation for Low-Latency Federated Edge Learning,” *IEEE Trans. Wirel. Commun.*, vol. 19, no. 1, pp. 491–506, 2020.
- [22] M. M. Amiri and D. Gündüz, “Federated learning over wireless fading channels,” *IEEE Trans. Wireless Comm.*, vol. 19, no. 5, pp. 3546–3557, 2020.
- [23] Y. Shao, D. Gündüz, and S. C. Liew, “Federated edge learning with misaligned over-the-air computation,” *arXiv cs.IT:2102.13604*, 2021.
- [24] M. M. Amiri, D. Gündüz, S. R. Kulkarni, and H. V. Poor, “Convergence of federated learning over a noisy downlink,” *arXiv 2008.11141*, 2020.
- [25] J. Castura and Y. Mao, “Rateless coding over fading channels,” *IEEE Communications Letters*, vol. 10, no. 1, pp. 46–48, 2006.
- [26] N. Ferdinand and S. C. Draper, “Hierarchical coded computation,” in *2018 IEEE Int. Symp. Inf. Theory (ISIT)*, June 2018, pp. 1620–1624.