

# AGE-BASED SCHEDULING POLICY FOR FEDERATED LEARNING IN MOBILE EDGE NETWORKS

Howard H. Yang<sup>†</sup>, Ahmed Arafa<sup>\*</sup>, Tony Q. S. Quek<sup>†</sup>, and H. Vincent Poor<sup>‡</sup>

<sup>†</sup> Singapore University of Technology and Design, Singapore 487372

<sup>\*</sup> University of North Carolina at Charlotte, NC 28223, USA

<sup>‡</sup> Princeton University, NJ 08544, USA

## ABSTRACT

Federated learning (FL) is a machine learning model that preserves data privacy in the training process. Specifically, FL brings the model directly to the user equipments (UEs) for local training, where an edge server periodically collects the trained parameters to produce an improved model and sends it back to the UEs. However, since communication usually occurs through a limited spectrum, only a portion of the UEs can update their parameters upon each global aggregation. As such, new scheduling algorithms have to be engineered to facilitate the full implementation of FL. In this paper, based on a metric termed the age of update (AoU), we propose a scheduling policy by jointly accounting for the staleness of the received parameters and the instantaneous channel qualities to improve the running efficiency of FL. The proposed algorithm has low complexity and its effectiveness is demonstrated by Monte Carlo simulations.

**Index Terms**— Federated learning, mobile edge computing, scheduling policy, age-of-update.

## 1. INTRODUCTION

Due to the sheer volume of data generated by the end-user devices, as well as the increasing concerns about sharing private information, a new machine learning model, namely the federated learning (FL) [1–4], has emerged from the intersection of artificial intelligence and edge computing [5–7]. In stark contrast to the conventional machine learning methods that run in a data center, FL usually operates at the network edge and brings the models directly to the devices for training, where only the resultant parameters shall be sent to the edge servers that reside in an access point (AP). This salient feature of *on-device training* brings along great advantages of eliminating the large communication overheads as well as preserving data privacy, and hence making FL particularly relevant for mobile applications [6–10]. However, as the AP needs to link a large number of user equipments (UEs) over a limited

spectrum, only a portion of the UEs can be selected to access the radio channel and send their trained updates in each global aggregation [5, 11]. To this end, the stragglers issue is more crucial for FL than for conventional trainings run in data centers [3, 4, 12].

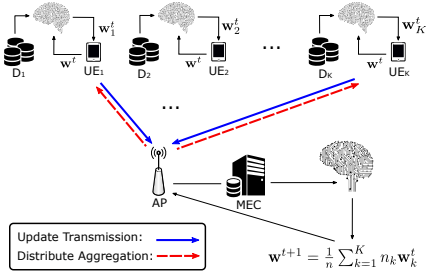
Recognizing such criticality, a host of studies have been carried out and resulted in various scheduling protocols for FL, ranging from minimizing the transmission latency [13], maximizing the spectral utility [14], to opportunistically alternating between selecting the UEs with advantageous and disadvantageous channel conditions [15]. Even though positive gains have been demonstrated, these works put the main focus on exploiting the spectral resources so as to maximize the number of updates collectible by the AP in each round of global communication but ignore the staleness of these updates. As pointed out by [16], the staleness of updates has a significant impact on the convergence rate of distributed machine learning models such as FL. Therefore, it is reasonable to expect a scheduling algorithm that accounts for both the staleness and the communication quality can accelerate the convergence of FL in a mobile edge network. On this purpose, we introduce a metric, referred to as the age-of-update (AoU)<sup>1</sup>, to measure the staleness associated with each update and based on that, devise a scheduling protocol that allows the AP to collect timely updates from the UEs for FL training. The algorithm is shown to possess low complexity, and its effectiveness is amply illustrated via Monte Carlo simulations.

## 2. SYSTEM MODEL

Let us consider a mobile edge network that consists of one AP and  $K$  user equipments (UEs), as depicted in Fig. 1, all are capable of performing local computing. For a generic UE  $k$ , we consider that it is equipped with single antenna and has local data set  $\mathcal{D}_k = \{\mathbf{x}_i \in \mathbb{R}^d, y_i \in \mathbb{R}\}_{i=1}^{n_k}$  with size

<sup>1</sup>Such age metric has been previously used, mainly, in the context of networking and has been shown to improve the notion of data freshness in various applications, see, e.g., the original treatment in [17], and the survey in [18].

This work was supported in part by the U.S. National Science Foundation under Grants CCF-0939370 and CCF-1513915.



**Fig. 1.** An illustration of the employed network model.

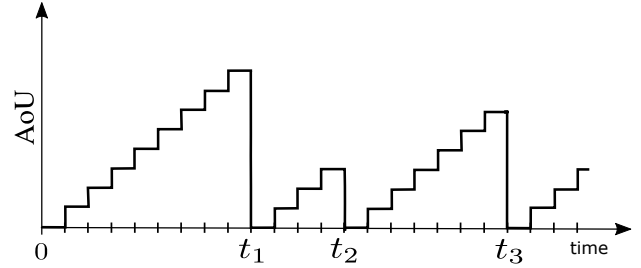
$n_k = |\mathcal{D}_k|$ , where  $|\cdot|$  denotes the cardinality. In this network, the UEs communicate with the AP through a common spectrum, which is divided into  $N$  equal-length subchannels. Here, we assume  $N < K$ . The AP can possibly assign different UEs a different number of subchannels. We adopt a block-fading propagation model, where the channels between any pair of antennas are assumed independently and identically distributed (i.i.d.) and quasi-static, i.e., the channel is constant during one transmission block, and varies independently from block to block.

The goal of the AP is to learn a statistical model over data that reside on the  $K$  associated UEs. Specifically, the AP needs to fit a vector  $\mathbf{w} \in \mathbb{R}^d$  so as to minimize a particular loss function. Formally, such a task can be expressed as follows:

$$\min_{\mathbf{w} \in \mathbb{R}^d} \left\{ P(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}; \mathbf{x}_i, y_i) + \xi r(\mathbf{w}) \right\} \quad (1)$$

where  $n = \sum_{i=1}^K n_i$  is the size of the whole data set,  $\xi$  is the regularizing parameter and  $r(\mathbf{w})$  a deterministic penalty function. The function  $\ell(\cdot)$  represents the loss function. Several examples of loss functions used in popular machine learning models are summarized in [19].

Due to concerns of data privacy, the raw data set is generally not available at the AP and hence (1) shall be solved using FL. In particular, the training procedure (cf. Algorithm 3) comprises three major steps: *i*) the UEs conduct local computing aimed to solve (1) by referring to the data resided on device and a downloaded global model, and send the resultant parameters to the AP, *ii*) the AP aggregates the received updates to produce an improved global model, and *iii*) the new model is sent back to UEs, and the process is repeated. After a sufficient amount of training and update exchanges, usually termed *communication rounds*, between the AP and UEs, the objective function (1) can converge to the global optimal. Nevertheless, because the wireless medium is resource-constrained, the AP can only select a subgroup of UEs for parameter updates in each communication round, and the choice of selection can largely affect the convergence rate of FL [19]. In that respect, it languishes a dire need for appropriate scheduling algorithms.



**Fig. 2.** An example of the AoU evolution in time, where the transmissions take place at time stamps  $t_1$ ,  $t_2$ , and  $t_3$ .

### 3. SCHEDULING POLICY DESIGN

In this section, we elaborate the concept of AoU, and detail the development of a scheduling policy based on this quantity.

#### 3.1. Design Metric

In order to quantify the staleness of each update, we leverage the concept of information freshness [17] and define a metric termed Age-of-Update (AoU). For a generic UE  $k$ , its AoU evolves as follows:

$$T_k[t+1] = (T_k[t] + 1)(1 - S_k[t]), \quad S_k[t] \in \{0, 1\} \quad (2)$$

where  $T_k[0] = 0$ , and  $S_k[t]$  takes value 1 if UE  $k$  is selected by the AP for update during communication round  $t$ , and takes value 0 otherwise. It is important to note that the AoU of UE  $k$  measures the time elapsed since its latest update is received by the AP, and hence a larger AoU value indicates a higher degree of staleness associated with that UE.

As illustrated in Fig. 2, the evolution of AoU is affected by the scheduling algorithm. Moreover, as the UEs are usually situated at various geographic locations, their AoU will progress differently due to their distinct communication conditions. In this regard, the different AoUs can be utilized as a reference to prioritize the channel access for the UEs.

#### 3.2. Problem Formulation

Using the notion of AoU, we devise our scheduling protocol such that the AP is able to keep the collection of all the updates “as fresh as possible” by minimizing the sum of certain functions of the AoUs in each communication round. To be concrete, upon the global aggregation at round  $t$ , the AP selects the set of UEs for parameter update, denoted by  $\mathbf{S}[t] = \{S_1[t], S_2[t], \dots, S_K[t]\}$ , according to the following criteria:

$$\min_{\mathbf{S}[t], \mathbf{P}} \sum_{k=1}^K f_\alpha(T_k[t])(1 - S_k[t]) \quad (3a)$$

$$\text{s.t.} \left( \sum_{n \in \mathcal{C}_k} \frac{1}{2} \log(1 + G_{k,n} P_{k,n}) - R_s \right) S_k[t] \geq 0 \quad (3b)$$

$$\sum_{n \in \mathcal{C}_k} P_{k,n} \leq P_{\text{TX}} \quad (3c)$$

$$S_k[t] \in \{0, 1\}, \quad \forall k \in \{1, \dots, K\} \quad (3d)$$

$$\mathcal{C}_i \cap \mathcal{C}_j = \emptyset, \text{ if } i \neq j \text{ and } \cup_{k=1}^K \mathcal{C}_k \subseteq \{1, \dots, N\} \quad (3e)$$

where  $G_{k,n}$  is the gain attainable by UE  $k$  on the  $n$ -th subchannel,  $P_{k,n}$  the correspondingly injected power, and  $\mathcal{C}_k$  is the set of subchannels assigned to UE  $k$ . Because the transmission of updates need to be finished within a certain time, the selected UEs should have their rates individually exceeding a threshold  $R_s$ , representing a required quality-of-service. Moreover, the UEs are subject to a maximum transmit power through constraint (3c), while the constraint (3e) stipulates that the UEs shall follow an orthogonal channel access methodology to avoid interfering with each other. The function  $f_\alpha(\cdot)$  in (3) represents the *sensitivity* of the edge server to the staleness of local updates. In this paper, we set the function as follows [20]:

$$f_\alpha(x) = \begin{cases} \frac{x^{1-\alpha}}{1-\alpha}, & \text{if } \alpha \neq 1, \\ \log(1+x), & \text{otherwise.} \end{cases} \quad (4)$$

The essence of the above function is to ensure a fairness treatment among the UEs based on their AoUs.

Two special cases are noteworthy: *i*) if  $f_\alpha(x)$  is set as a constant, then solving problem (3) is equivalent to packing the maximum number of UEs into the spectrum during each communication round [14], and *ii*) when the communication is conducted under an ideal environment, i.e., if  $R_s = 0$ , solutions of the optimization problem (3) lead to an  $N$ -Round-Robin policy [19]. Therefore, in the general scenario, the design problem given by (3) balances the tradeoff between fairness and quality on the radio channel access. It is also worthwhile to note that the design per problem (3) is not only suitable for scheduling updates in FL, but can also be leveraged in the scheduling of wireless traffic [21].

### 3.3. Algorithm

The problem given by (3) is a combinatorial problem with mixed integers that does not possess low complexity solutions in general. Here, we tackle the problem via a greedy algorithm. The approach is constituted of two major steps.

First, amongst all the UEs, the AP needs to form a candidate list  $\mathcal{J}_s$  for further selection, by choosing those UEs that are able to achieve the targeted rate  $R_s$  under the currently available spectrum budget. Moreover, each UE in the candidate list shall occupy the least amount of subchannels. As

---

#### Algorithm 1 Candidate UE List & Subchannel Allocation

---

**Input:** The set of currently available subchannels  $\bar{N} \subseteq \mathcal{N}$   
**Initialize:**  $\mathcal{J}_s = \emptyset$  and  $\mathcal{C}_k = \emptyset, \forall k \in \{1, 2, \dots, K\}$   
**for** each UE  $k \in \{1, 2, \dots, K\}$  **do**  
    set  $n_k^* = 1$  and order the channel gains as  $G_{k,\tau(1)} \geq G_{k,\tau(2)} \geq \dots \geq G_{k,\tau(\bar{N})}$  where  $\tau(n) \in \bar{N}$   
    **while**  $n_k^* \leq \bar{N}$  **do**  
        Compute the following power updates:

$$P_{k,\tau(n)} = \begin{cases} \frac{P(n_k^*)}{n_k^*} + \frac{1}{G_{k,n_k^*}} - \frac{1}{G_{k,\tau(n)}}, & 1 \leq n \leq n_k^*, \\ 0, & n_k^* < n \leq \bar{N} \end{cases} \quad (5)$$

in which  $P(n_k^*)$  is given by

$$P(n_k^*) = \left[ P_{\text{TX}} - \sum_{n=1}^{n_k^*-1} \left( \frac{1}{G_{k,n_k^*}} - \frac{1}{G_{k,\tau(n)}} \right) \right]^+ \quad (6)$$

where  $[\cdot]^+ = \max(\cdot, 0)$

**if**  $\sum_{n=1}^{n_k^*} \log(1 + G_{k,\tau(n)} P_{k,\tau(n)}) \geq R_s$  **then**  
    Update  $\mathcal{J}_s = \mathcal{J}_s \cup \{k\}$  and  $\mathcal{C}_k = \cup_{n=1}^{n_k^*} \{\tau(n)\}$   
    Set  $n_k^* = \bar{N} + 1$   
**else**  
    Update  $n_k^* = n_k^* + 1$

**Output:**  $\mathcal{J}_s, \{\mathcal{C}_k, k \in \mathcal{J}_s\}$

---



---

#### Algorithm 2 Age-Based Scheduling (ABS)

---

**Input:**  $\mathcal{J}_s, \{\mathcal{C}_k, k \in \mathcal{J}_s\}$ , and  $\mathbf{T}[t] = \{T_k[t], k \in \mathcal{J}_s\}$   
**Initialize:**  $S_k[t] = 0, k \in \{1, 2, \dots, K\}$   
**if**  $\mathcal{J}_s \neq \emptyset$  **then**  
    Assign  $S_{k^*}[t] = 1$ , in which the index  $k^*$  is chosen as  $k^* = \arg \max_{k \in \mathcal{J}_s} f_\alpha(T_k[t]) / \|\mathcal{C}_k\|_0$   
    Update the available subchannels as  $\bar{N} = \bar{N} \setminus \mathcal{C}_{k^*}$   
    Call Algorithm 1 with updated input  $\bar{N}$  and obtain the new outputs  $\mathcal{J}_s, \{\mathcal{C}_k, k \in \mathcal{J}_s\}$   
    Call Algorithm 2 with updated inputs  $\mathcal{J}_s, \{\mathcal{C}_k, k \in \mathcal{J}_s\}$   
**Output:**  $\mathbf{S}[t] = \{S_1[t], S_2[t], \dots, S_K[t]\}$

---

such, a candidate UE  $k$  needs to satisfy the following:

$$\min_{\mathcal{C}_k, \mathbf{P}_k} \|\mathcal{C}_k\|_0 \quad (7a)$$

$$\text{s.t.} \sum_{n \in \mathcal{C}_k} \frac{1}{2} \log(1 + G_{k,n} P_{k,n}) \geq R_s, \quad (7b)$$

$$\sum_{n \in \mathcal{C}_k} P_{k,n} \leq P_{\text{TX}} \quad (7c)$$

where  $\|\cdot\|_0$  is the L-0 norm [22] and  $\mathbf{P}_k = \{P_{k,n}, n \in \mathcal{C}_k\}$  is the transmit power on the allotted subchannels. For each selection of  $\mathcal{C}_k$ , subproblem (7) can be solved by a water-filling algorithm, and the details of forming the candidate list is summarized in Algorithm 1.

---

**Algorithm 3** Wireless Federated Learning

**Parameters:**  $\tau$  = number of local steps per communication round,  $\eta$  = step size for local learning

**Initialize:** The AP assigns  $T_k[0] = 0, \forall k \in \{1, \dots, K\}$  and randomly initializes  $\mathbf{w}^0 \in \mathbb{R}^d$

**for**  $t = 0, 1, 2, \dots, T - 1$  **do**

The AP runs ABS (Algorithm 2) to get  $\mathbf{S}[t]$ , and sends the parameter  $\mathbf{w}^t$  to all the UEs in the updating list

**for** each UE with  $S_k[t] = 1$  **do**

Initialize  $\mathbf{w}_k^t = \mathbf{w}^t$

**for**  $s = 1$  to  $\tau$  **do**

Sample  $i \in \mathcal{D}_k$  uniformly at random, and update the local parameter  $\mathbf{w}_k^t$  as follows:

$$\mathbf{w}_k^t = \mathbf{w}_k^t - \eta(\nabla \ell(\mathbf{w}_k^t; \mathbf{x}_i, y_i) + \xi \nabla r(\mathbf{w}_k^t)) \quad (9)$$

The UE updates the parameter  $\mathbf{w}_k^t$  to the AP

The AP collects the trained weights  $\{\mathbf{w}_k^t\}_{k \in \mathbf{S}[t]}$ , and updates the AoUs  $\mathbf{T}[t]$  and the global parameter  $\mathbf{w}^{t+1}$  respectively as follows:

$$\mathbf{T}[t+1] = (\mathbf{T}[t] + \mathbf{1}) \odot (\mathbf{1} - \mathbf{S}[t]), \quad (10)$$

$$\mathbf{w}^{t+1} = \frac{\sum_{k \in \mathbf{S}[t]} |\mathcal{D}_k| \mathbf{w}_k^t}{\sum_{k \in \mathbf{S}[t]} |\mathcal{D}_k|} \quad (11)$$

where  $\mathbf{1} \in \mathbb{R}^K$  is the all-one vector and  $\odot$  represents the Hadamard product

**Output:**  $\mathbf{w}^T$

---

Next, from the set of candidate UEs, the AP shall select the UEs that have relatively large AoU while requiring a relatively small amount of subchannels for the transmissions. Formally, the AP will select UE  $i$  instead of UE  $j$  if the following holds:

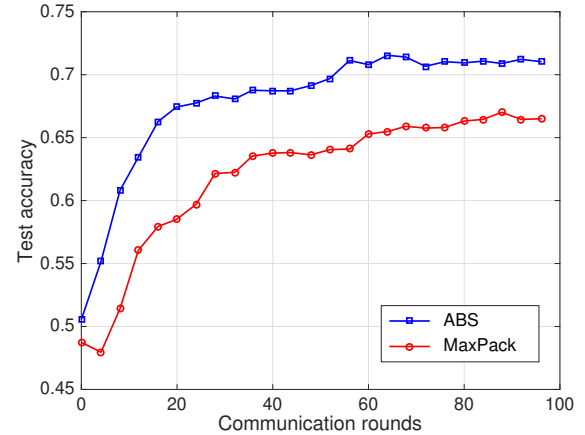
$$f_\alpha(T_i[t]) / \|\mathcal{C}_i\|_0 > f_\alpha(T_j[t]) / \|\mathcal{C}_j\|_0. \quad (8)$$

The above criteria brings us to the age-based scheduling (ABS) policy, summarized in Algorithm 2. Note that the ABS adopts a recursive search for UEs, since the occupation of a particular subchannel can incur other UEs in the candidate list to violate their rate constraint and hence get removed. Therefore, the candidate list  $\mathcal{J}_s$  needs to be recalculated via Algorithm 1 whenever a UE is selected. Nevertheless, as the UE selection and subchannel assignments are conducted in a greedy manner, the algorithm has relatively low complexity.

Finally, we present the FL approach to solving (1) in the setting of a mobile edge network, as per Algorithm 3.

#### 4. NUMERICAL RESULTS

In this section, we carry out simulations to verify the effectiveness of the proposed scheduling policy. Particularly, we



**Fig. 3.** Impact of scheduling policies on the convergence of FL in terms accuracy value.

evaluate the efficiency of training FL on a support vector machine (SVM) over the MNIST data set, which consists of handwritten numerals, under the proposed ABS policy, and a state-of-the-art approach [14], coined as the MaxPack policy. The whole data set is participated in 100 non-overlapped portions and assigned to the  $K = 100$  UEs that are randomly scattered within a disc of a radius 100 m, where the AP locates at the center. There are in total  $N = 20$  available subchannels, and the transmission on any subchannel is subject to Rayleigh fading with unit mean and the path loss that follows a power law, with the exponent being  $\beta = 3.5$ .

Fig. 3 plots the test accuracy as a function of the communication rounds. From this figure, we can observe that FL running under the ABS policy is able to attain a marked improvement in the accuracy value compared to the MaxPack policy. Moreover, the gain achieved by the ABS policy is especially pronounced in the initial stage of the training, e.g., when the communication rounds are less than 40, which demonstrates that the proposed scheme is able to accelerate the convergence of FL and hence improve the learning efficiency. This gain can be essentially attributed to that while the MaxPack policy exploits the multi-user diversity [23] by allocating the best subchannels and power to maximize the total number of receivable updates, the ABS strikes a balance between the aggressive channel use and short term fairness and hence achieves the gradient diversity [24].

#### 5. CONCLUSION

In this paper, we have proposed a scheduling policy for FL in the context of mobile edge networks. By adopting a metric termed AoU, our scheme is able to account for both the staleness of updates and the instantaneous channel qualities so as to accelerate the convergence rate of FL.

## 6. REFERENCES

- [1] J. Konečný, H. B. McMahan, F. X. Yu, P. Richtárik, A. T. Suresh, and D. Bacon, “Federated learning: Strategies for improving communication efficiency,” *arXiv preprint arXiv:1610.05492*, 2016.
- [2] H. B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. Arcas, “Communication-efficient learning of deep networks from decentralized data,” *Available as ArXiv:1602.05629*, 2016.
- [3] Y. Lin, S. Han, H. Mao, Y. Wang, and W. J. Dally, “Deep gradient compression: Reducing the communication bandwidth for distributed training,” in *Int. Conf. Learn. Represent. (ICLR)*, Vancouver, Canada, May 2018.
- [4] V. Smith, S. Forte, C. Ma, M. Takáč, M. I Jordan, and M. Jaggi, “CoCoA: A general framework for communication-efficient distributed optimization,” *J. Machine Learning Research*, vol. 18, no. 230, pp. 1–49, 2018.
- [5] Z. Zhao, C. Feng, H. H. Yang, and X. Luo, “Federated learning-enabled intelligent fog-radio access networks: Fundamental theory, key techniques, and future trends,” *IEEE Wireless Commun. Mag.*, 2020.
- [6] S. Wang, T. Tuor, T. Salonidis, K. K Leung, C. Makaya, T. He, and K. Chan, “Adaptive federated learning in resource constrained edge computing systems,” *IEEE J. Sel. Areas Commun.*, vol. 37, no. 6, pp. 1205–1221, Jun. 2019.
- [7] K. B. Letaief, W. Chen, Y. Shi, J. Zhang, and Y. J. A. Zhang, “The roadmap to 6G–AI empowered wireless networks,” *IEEE Commun. Mag.*, vol. 57, no. 8, pp. 84–90, Aug. 2019.
- [8] M. Chen, Z. Yang, W. Saad, C. Yin, H. V. Poor, and S. Cui, “A joint learning and communications framework for federated learning over wireless networks,” *Available as ArXiv:1909.07972*, 2019.
- [9] Y. Du, S. Yang, and K. Huang, “High-dimensional stochastic gradient quantization for communication-efficient edge learning,” *Available as ArXiv:1910.03865*, 2019.
- [10] M. M. Amiri and D. Gunduz, “Federated learning over wireless fading channels,” *Available as ArXiv:1907.09769*, 2019.
- [11] Y. Mao, C. You, J. Zhang, K. Huang, and K. B. Letaief, “A survey on mobile edge computing: The communication perspective,” *IEEE Commun. Surveys & Tutorials*, vol. 19, no. 4, pp. 2322–2358, Aug. 2017.
- [12] S. Ha, J. Zhang, O. Simeone, and J. Kang, “Coded federated computing in wireless networks with straggling devices and imperfect CSI,” *Available as ArXiv:1901.05239*, 2019.
- [13] T. Nishio and R. Yonetani, “Client selection for federated learning with heterogeneous resources in mobile edge,” in *Proc. IEEE Int. Conf. Commun.*, Shanghai, China, May 2019, pp. 1–6.
- [14] K. Yang, T. Jiang, Y. Shi, and Z. Ding, “Federated learning via over-the-air computation,” *IEEE Trans. Wireless Commun.*, 2019.
- [15] G. Zhu, Y. Wang, and K. Huang, “Broadband analog aggregation for low-latency federated edge learning,” *IEEE Trans. Wireless Commun.*, 2019.
- [16] W. Dai, Y. Zhou, N. Dong, H. Zhang, and E. P. Xing, “Toward understanding the impact of staleness in distributed machine learning,” in *Int. Conf. Learn. Represent. (ICLR)*, New Orleans, Louisiana, May 2019, pp. 1–6.
- [17] S. Kaul, R. Yates, and M. Gruteser, “Real-time status: How often should one update?,” in *Proc. IEEE INFOCOM*, Orlando, FL, Mar. 2012, pp. 2731–2735.
- [18] A. Kosta, N. Pappas, and V. Angelakis, “Age of information: A new concept, metric, and tool,” *Foundations and Trends in Networking*, vol. 12, no. 3, pp. 162–259, 2017.
- [19] H. H. Yang, Z. Liu, T. Q. S. Quek, and H. V. Poor, “Scheduling policies for federated learning in wireless networks,” *IEEE Trans. Commun.*, 2019.
- [20] A. B. Atkinson, “On the measurement of inequality,” *J. Econ. Theory*, vol. 2, no. 3, pp. 244–263, Sept. 1970.
- [21] Y. Cao and V. O. K. Li, “Scheduling algorithms in broadband wireless networks,” *Proc. IEEE*, vol. 89, no. 1, pp. 76–87, Jan. 2001.
- [22] E. J. Candès, J. Romberg, and T. Tao, “Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information,” *IEEE Trans. Inf. Theory*, vol. 52, no. 2, pp. 489–509, Feb. 2006.
- [23] D. Tse and P. Viswanath, *Fundamentals of wireless communication*, Cambridge university press, 2005.
- [24] D. Yin, A. Pananjady, M. Lam, D. Papailiopoulos, K. Ramchandran, and P. L. Bartlett, “Gradient diversity: A key ingredient for scalable distributed learning,” in *Proc. Machine Learning Research*, Playa Blanca, Lanzarote, Canary Islands, Apr. 2018, pp. 1998–2007.