

Adaptive Gradient Methods For Over-the-Air Federated Learning

Chenhao Wang^{§*}, Zihan Chen^{†*}, Howard H. Yang[§], and Nikolaos Pappas[‡]

[§] ZJU-UIUC Institute, Zhejiang University, Haining 314400, China

[†] Singapore University of Technology and Design, Singapore 487372, Singapore

[‡] Department of Computer and Information Science, Linköping University, Linköping 58183, Sweden

Abstract—Federated learning (FL) provides a privacy-preserving approach to realizing networked intelligence. But the performance of FL is often constrained by the limited communication resources, especially in the context of a wireless system. To tackle this communication bottleneck, recent studies propose an analog over-the-air (A-OTA) FL paradigm which employs A-OTA computations in the model aggregation step, significantly enhancing system scalability. However, these existing architectures mainly conduct model training via (stochastic) gradient descent, while adaptive optimization methods, which have achieved notable success in deep learning, remain unexplored. In this paper, we establish a distributed training paradigm that incorporates adaptive gradient methods into the A-OTA FL framework, aiming to enhance the system’s convergence performance. We derive an analytical expression for the convergence rate, capturing the effects of various system parameters on the convergence performances of the proposed method. We also perform several experiments to validate the efficacy of the proposed method.

I. INTRODUCTION

Federated learning (FL) is a promising technology that emerged from the intersection of artificial intelligence and edge computing, facilitating networked intelligence while preserving end-users’ data privacy [1]–[4]. However, a typical FL system consists of multiple clients, and the training process of FL requires frequent model parameter exchanges between the end-user devices (a.k.a. clients) and the server, which brings hefty communication overhead. For networks with limited communication resources, such communication bottleneck can throttle the scalability of the FL system and concurrently constrain the training efficiency of the system.

To address this issue, recent studies [5]–[8] suggest using analog over-the-air (A-OTA) computations in the FL system, in which clients modulate their local gradients onto a set of common waveforms and simultaneously send out the analog signals to the edge server for fast model aggregation via exploiting the superposition property of wireless waveforms. As a result, the A-OTA computation-based FL system achieves significantly improved spectral efficiency, reduced access latency, and enhanced privacy protections for the clients [6], [8], [9].

This work was supported in part by the National Natural Science Foundation of China under Grant 62271513. (*: Equal contribution; *Corresponding Author: Howard H. Yang.*)

However, A-OTA computations inevitably introduce signal distortion to the globally aggregated gradients [10]–[13]. Specifically, random channel fading and interference (i.e., channel noise) would be induced into the aggregated gradients at the edge server. The consequent noisy aggregated signal would result in performance degradation, such as a slower convergence rate and unstable training performance in the A-OTA FL system.

To harness such uncertainty while maintaining the benefits of the analog over-the-air computations, numerous work [12] have been proposed to improve the system performance, among which the majority of these approaches are built on FedAVG [1] or FedSGD [13] for gradient descent-based global model updates. But adaptive optimization methods, which are demonstrated with enhanced performance in non-FL setup compared to the standard gradient descent methods, are rarely explored in such FL scenarios. For instance, adaptive methods such as AdaGrad [14] and Adam [15] have been widely used in various fields with significant improvements in training effectiveness and robustness [16]–[18]. Adaptive algorithms are known to be interference-resistant in deep learning scenarios, which is a characteristic that also applies to A-OTA computing scenarios. Therefore, a natural question arises: *in the presence of noisy aggregated signals as well as the inherent heterogeneous data*, does the adaptive optimization methods still benefit the performance of A-OTA FL (and/or, does such a training method even converge)?

The present paper answers the above question by developing an A-OTA FL framework, namely, ADOTA-FL, which incorporates adaptive gradient descent techniques. To the best of our knowledge, this is the first work that explores the adaptive gradient method for the A-OTA FL setting. Our theoretical analysis provides the upper bound performance for the convergence rate of ADOTA-FL. The numerical results of extensive experiments also confirm the efficacy of our proposed framework.

II. SYSTEM MODEL

We consider a wireless edge network, as depicted in Fig. 1, that comprises one server and N clients. The clients communicate with the server through wireless transmissions over a shared spectrum. Every client $n \in \{1, \dots, N\}$ has its own local dataset $\mathcal{D}_n = \{(\mathbf{x}_i \in \mathbb{R}^u, y_i \in \mathbb{R})\}_{i=1}^{m_n}$ with size $|\mathcal{D}_n| = m_n$.

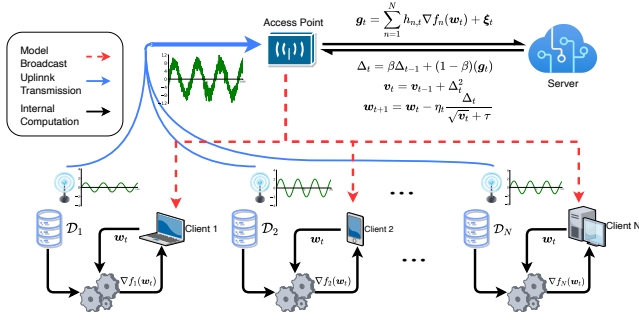


Fig. 1. An overview of the adaptive algorithm for analog over-the-air federated learning.

We assume the local datasets are statistically independent across the clients. The goal of all the entities in this system is to collaboratively learn a statistical model using data from all the clients without violating their data privacy. Thus, the clients need to solve an optimization problem of the following form [1]:

$$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N f_n(\mathbf{w}; \mathcal{D}_n), \quad (1)$$

where $f_n(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}$ is the local empirical loss function of client n , constructed from its own dataset \mathcal{D}_n and $\mathbf{w} \in \mathbb{R}^d$ is the global model parameter with dimension d . We denote the optimal solution to (1) by \mathbf{w}^* , i.e.,

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} f(\mathbf{w}). \quad (2)$$

To obtain \mathbf{w}^* while concurrently preserving the data privacy of each client, the entities need to carry out the model training in an FL fashion. Specifically, the clients shall train their models locally and upload the intermediate gradients to the server rather than local private data. The server aggregates the clients' gradients and further improves the global model. Then, the server broadcasts the model to all the clients for another round of local training. Such interactions will repeat until the global model converges.

Due to the limited spectral resources, the efficiency of such federated training procedure is often throttled by the communication bottleneck, i.e., the server needs to select a portion of clients for parameter uploading in each communication round, which becomes cumbersome when the number of participants becomes large. In the following section, we introduce a model training framework that addresses this bottleneck using the A-OTA computation method. Additionally, it is devised on the basis of the adaptive gradient descent method, which has the potential to accelerate the model training process. Owing to these two attributes, we coin our method as *adaptive over-the-air federated learning (ADOTA-FL)*.

III. ADAPTIVE OVER-THE-AIR FEDERATED LEARNING

This section elaborates on the training process of ADOTA-FL. In particular, this method employs analog over-the-air

computations in the global gradient aggregation stage, substantially reducing access latency and facilitating (theoretically unlimited) algorithm scalability. In the global model improvement stage, it leverages ideas from AdaGrad methods to potentially speed up convergence. The general steps are summarized in Algorithm 1. And more details are provided below.

1) *Local Model Training*: Without loss of generality, we assume the system has progressed to the k -th round of global training, upon which the clients just received the global model parameters \mathbf{w}_k from the edge server.¹ Then, every client n updates its local gradient $\nabla f_n(\mathbf{w}_k)$ by taking the global model as an input.

2) *Analog Gradient Aggregation*: In this system, the clients use analog transmissions to upload their locally trained gradients. Specifically, once $\nabla f_n(\mathbf{w}_k)$ is computed, client n modulates this gradient vector entry-by-entry onto the magnitudes of a common set of orthogonal baseband waveforms, arriving at the following analog signal

$$x_n(t) = \langle \mathbf{s}(t), \nabla f_n(\mathbf{w}_k) \rangle, \quad (3)$$

where $\langle \cdot, \cdot \rangle$ denotes the inner product between two vectors and $\mathbf{s}(t) = (s_1(t), \dots, s_d(t))$, $t \in [0, \tau_0]$ has its entries satisfying

$$\int_0^{\tau_0} s_i^2(t) dt = 1, \quad i = 1, 2, \dots, d, \quad (4)$$

$$\int_0^{\tau_0} s_i(t) s_j(t) dt = 0, \quad i \neq j. \quad (5)$$

τ_0 represents the total time of signal duration.

Once the analog waveforms $\{x_n(t)\}_{n=1}^N$ have been assembled, the clients transmit them simultaneously to the edge server. Thanks to the superposition property of electromagnetic waves, the signal received at the edge server obeys the following form:

$$y(t) = \sum_{n=1}^N h_{t,n} P_n x_n(t) + \xi(t), \quad (6)$$

where $h_{t,n}$ is the channel fading experienced by client n , P_n is the corresponding transmit power, and $\xi(t)$ denotes the additive noise. In this work, we assume the channel fading is i.i.d. across clients, with mean and variance being μ_c and σ_c^2 , respectively. Besides, the transmit power of each client is set to compensate for the large-scale path loss. Additionally, we assume the noise follows a Gaussian distribution with variance σ_n^2 . This received signal will be passed through a bank of match filters, with each branch tuning to $s_i(t)$, $i = 1, 2, \dots, d$. On the output side, the server obtains the following vector:

$$\mathbf{g}_t = \frac{1}{N} \sum_{n=1}^N h_{t,n} \nabla f_n(\mathbf{w}_t) + \boldsymbol{\xi}_t, \quad (7)$$

in which $\boldsymbol{\xi}_t$ is a d -dimensional random vector with each entry being i.i.d. and follows the Gaussian distribution. It

¹Since the server can broadcast its signal at a high transmit power, we assume the global model can be successfully received by all the clients.

is worthwhile to stress that the vector given in Eq. (7) is a distorted version of the globally aggregated gradient.²

3) *Global Model Update*: Using \mathbf{g}_t , the server can update the global model via adaptive optimization methods at the end of communication round t . However, due to channel fading and interference noise perturbations, the globally aggregated gradient may experience a significant distortion, and the general adaptive method may not perform well. To abbreviate the impact of such distortions, we store and update an intermediate global model as follows:

$$\Delta_t = \beta \Delta_{t-1} + (1 - \beta) \mathbf{g}_t, \quad (8)$$

where β is a controlling factor that adjusts the portion between the historical information and the newly acquired information. Notably, the operation in Eq. (8) introduces a momentum-like approach to smooth out the fluctuation in the aggregated gradient. As such, a smaller value of β leads to a faster convergence but also incurs stronger volatility.

Aided by Δ_t , we construct a vector \mathbf{v}_t as follows:

$$\mathbf{v}_t = \mathbf{v}_{t-1} + \Delta_t^2, \quad (9)$$

where in Δ_t^2 , the square is performed entry-wise. The vector \mathbf{v}_t is used as the updated global model's denominator, where it adjusts the learning rate corresponding to each entry in the global model. The purpose of updating \mathbf{v}_t using Δ_t is to make each dimension of the gradients has its learning rate separately, which is related to the historical accumulation of the gradients. The specific impact on the algorithm will be described in Eq. (10).

Finally, using \mathbf{v}_t and Δ_t , we update the global model as:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \frac{\Delta_t}{\sqrt{\mathbf{v}_t} + \tau}, \quad (10)$$

where η_t is the learning rate factor, and τ is a constant that prevents ill-conditioning. It is important to note that the operation $\sqrt{\mathbf{v}_t}$ is taken with respect to each entry of \mathbf{v}_t . As mentioned earlier, the squared gradient information is accumulated in the corresponding dimension in each iteration by \mathbf{v}_t . This cumulative value, as a denominator, results in smaller dimensions having smaller learning decays. That means, the more gentle the gradient direction, the larger step it will get. In this way, we can optimize the system better.

The new global model \mathbf{w}_{t+1} will be broadcasted to all the clients for the next round of local computing. And the clients and server will repeat this process for multiple rounds until the global model converges.

IV. ANALYSIS

In this section, we analyze the convergence rate of the proposed ADOTA-FL. We assume each client's local objective function $f_n(\cdot)$ is convex. As such, the global objective function $f(\cdot)$ is also convex. We denote by $\mathbf{g}_{t,n} = h_{t,n} \nabla f_n(\mathbf{w}_t) + \boldsymbol{\xi}_t$ as the noisy gradient of client n in the t -th global iteration.

²Note that such a system is practically achievable. A recent prototype of the A-OTA FL system is presented in [8].

Algorithm 1 Adaptive Over-the-Air FL (ADOTA-FL)

Input: Initial delay vector \mathbf{v}_{-1} , initial global model \mathbf{w}_0 , T, η_0

- 1: **for** $t = 0, 1, 2, \dots, T - 1$ **do**
- 2: **for** each client $n \in N$ **in parallel do**
 # Train model locally and upload gradients
- 3: $\nabla f_n(\mathbf{w}_t) \leftarrow \text{CLIENTUPDATE}(t, \mathbf{w}_t)$
- 4: $\mathbf{g}_t = \frac{1}{N} \sum_{n=1}^N h_{t,n} \nabla f_n(\mathbf{w}_t) + \boldsymbol{\xi}_t$
- 5: $\Delta_t = \beta \Delta_{t-1} + (1 - \beta) \mathbf{g}_t$
- 6: $\mathbf{v}_t = \mathbf{v}_{t-1} + \Delta_t^2$
- 7: $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \frac{\Delta_t}{\sqrt{\mathbf{v}_t} + \tau}$

Output: \mathbf{w}_T

The following assumptions are adopted in our proof.

Assumption 1. There exists a constant C such that $\|\mathbf{w}_t - \mathbf{w}^*\|_2 \leq C$ for any $t \in \{1, \dots, T\}$.

Assumption 2. There exists a constant D such that $\|\frac{\mathbf{g}_{t,n}}{\sqrt{\mathbf{v}_t} + \tau}\|_2 \leq D$ for any $t \in \{1, \dots, T\}$ and $n \in \{1, \dots, N\}$.

For any vector \mathbf{a} , a^j denotes the j -th coordinate of \mathbf{a} . For example, \mathbf{g}_t^j means the j -th coordinate of gradient \mathbf{g}_t , and \mathbf{v}_t^j means the j -th coordinate of gradient \mathbf{v}_t . For the ease of expository, we assume $\beta = 0$, though our analysis can be directly extended to $\beta > 0$.

Lemma 1. Let Assumptions 1, 2 hold, then we can get the auxiliary lemma by induction:

$$\sum_{t=1}^T \mathbb{E} \left[\frac{(\mathbf{g}_t^j)^2}{\sqrt{t\mathbf{v}_t^j} + \sqrt{t}\tau} \right] \leq \mathbb{E} \left[\sqrt{T\mathbf{v}_T^j} + \sqrt{T}\tau \right]. \quad (11)$$

In order to reflect the effects caused by channel fading and interference in the final results, we consider an auxiliary factor. We denote the auxiliary factor as

$$\tilde{\mathbf{v}}_t = \tilde{\mathbf{v}}_{t-1} + (\nabla f(\mathbf{w}_t))^2, \quad (12)$$

so $\tilde{\mathbf{v}}_T^j$ means the j -th coordinate of $\tilde{\mathbf{v}}_T$. Then we introduce this factor into the final result

Theorem 1. Let assumptions 1 and 2 hold. After running Algorithm 1 for T iterations with initial learning rate η_1 , we have

$$\begin{aligned} & \min_{t \in \{1, \dots, T\}} \mathbb{E} [f(\mathbf{w}_t)] - f^* \\ & \leq \frac{G}{\sqrt{T}} \sum_{j=1}^d \mathbb{E} \left[\sqrt{(\mu_c^2 + \sigma_c^2) \tilde{\mathbf{v}}_T^j + \sigma_n + \tau} \right], \end{aligned} \quad (13)$$

where $f^* := f(\mathbf{w}^*)$ and $G = \frac{D^2}{2\eta_1} + \frac{\eta_1}{2}$.

Proof. We omit the proof due to space limitations. \square

ADOTA-FL's upper bound is provided by Theorem 1, which suggests that the algorithm converges. The upper bound indicates that channel fading and interference work together to affect the algorithm's convergence.

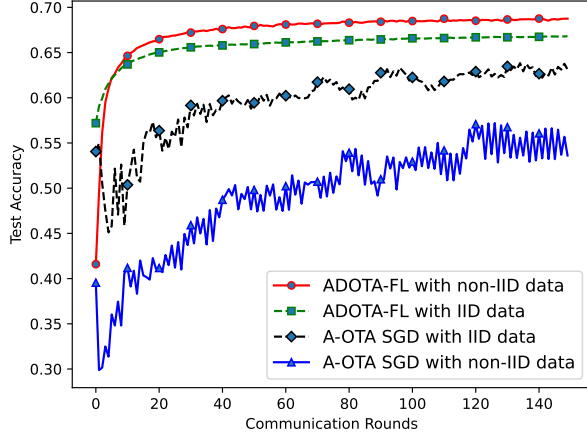


Fig. 2. Performance comparison for test accuracy on EMNIST with non-IID settings

V. NUMERICAL RESULTS

We conduct extensive experiments in this section to examine the performance and efficacy of the proposed ADOTA-FL. We explore the impacts of parameters pertaining to the algorithm (e.g., hyperparameters) and system (e.g., the number of participating clients in the network) on the performance of the A-OTA federated learning framework with diverse data settings and tasks. More details are given in the sequel.

A. Experimental setup

To evaluate the performance of the proposed framework, we conduct experiments with both convex and non-convex settings. Starting with the theoretical findings, we first carry out the use logistic regression model for multi-class classification tasks on EMNIST [19]. To further consolidate the findings, we carry out the experiments of image classification tasks on the CIFAR-10 [20] with ResNet-18 [21]. All the reported results are averaged with 3 trials. Unless otherwise specified, we use the total number of clients $N = 100$. Both the IID and non-IID data settings are evaluated, in which non-IID data are partitioned with widely used Dirichlet distribution with parameter $\alpha = 0.1$.

B. Performance evaluation

We compare the model performance of ADOTA-FL with conventional FedSGD in A-OTA setting (i.e., A-OTA SGD), with respect to the test accuracy and the training loss performance. We also provide a sensitivity analysis of the key system parameter and hyper-parameters. Overall, our ADOTA-FL method outperforms the original A-OTA method in terms of training loss across various models and datasets.

Convex Settings: We first evaluate our method with convex settings. Fig. 2 illustrates the performance of ADOTA-FL as well as the performance comparison with A-OTA SGD on both IID and non-IID EMNIST, using the same system configurations. As shown in the figure, ADOTA-FL is convergent in the

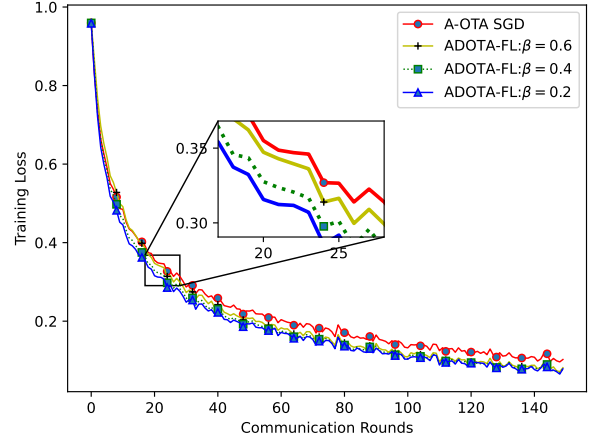


Fig. 3. Performance comparison for training loss on CIFAR-10 with non-IID settings.

context of A-OTA FL scenarios and consistently outperforms A-OTA SGD in both IID and non-IID settings. Furthermore, ADOTA-FL shows superiority in terms of stability, in the presence of the same channel fading, interference, and data heterogeneity in A-OTA FL systems.

Non-convex Settings: Fig. 3 compares the training loss performances with non-IID CIFAR-10 under A-OTA SGD [13] and ADOTA-FL, where the impacts of different controlling factors β are also investigated. It can be observed that ADOTA-FL attains a faster convergence rate than the A-OTA SGD for a varying value of β , verifying the effectiveness of employing adaptive optimization methods in the framework of A-OTA FL systems. Moreover, we notice that, as the β value decreases, the convergence speed of ADOTA-FL increases. This observation is consistent with what we expected because β is related to historical information, and the larger the β , the more information will be retained during the model update process, which will result in slower convergence but more stable convergence overall. Note that the results obtained under IID settings present a similar pattern to the non-IID case.

Performance with different system scale: In Fig. 4, we depict the test accuracy of the non-IID setting under ADOTA-FL as a function of communication rounds under different values of client number N . This figure unveils a phenomenon unique to A-OTA FL systems, that an increase in the number of participating clients can improve the system's performance [13]. This is because, using A-OTA computing, all the clients can simultaneously upload their locally computed gradients in each communication round. In contrast to the conventionally adopted digital communication-based FL that requires scheduling [3], it significantly boosts up the amount of information aggregated in every global iteration, thus enhancing the generalization performance. The aggregation of the imperfect gradients from more local clients could mitigate the negative impacts of channel interference and fading.

Performance with different data settings: We also in-

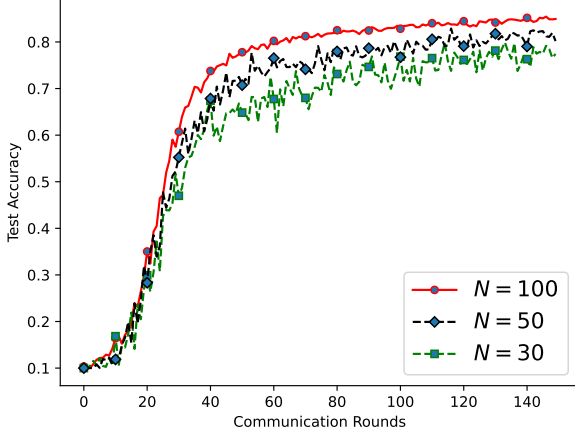


Fig. 4. Performance comparison for test accuracy with the different total number of clients N on CIFAR-10 with non-IID data.

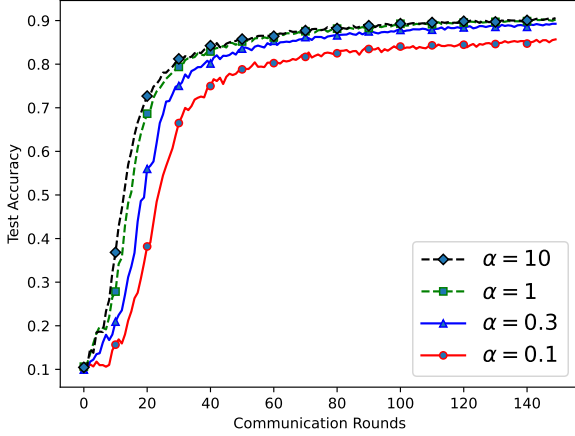


Fig. 5. Performance comparison for test accuracy with different α on CIFAR-10 with non-IID data.

investigated the performance of the ADOTA-FL with different degrees of data heterogeneity. To specify the degree of discrepancy of the data distribution, we use different concentration parameters α of Dirichlet distribution, where a smaller α would lead to more heterogeneous local data. Fig. 5 shows the impacts of different α on training performance. Our proposed ADOTA-FL presents stable performance with diverse heterogeneous data settings.

VI. CONCLUSION

In this paper, we proposed ADOTA-FL, which uses an adaptive optimization method to improve the convergence performance of OTA computing-based FL systems. We derived the convergence rate of ADOTA-FL, accounting for effects from both algorithmic and system perspectives. We also carried out several experiments to demonstrate the efficacy of the proposed method and explored the effects of different system

parameters on the ADOTA-FL performance. The robustness performance would be further explored in our future work.

REFERENCES

- [1] D. Ramage S. Hampson H. B. McMahan, E. Moore and B. A. Y. Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Proc. 20th Int. Conf. Artif. Intell. Stat. (AISTATS)*, Fort Lauderdale, FL, USA, Apr. 2017, pp. 1273–1282.
- [2] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith, "Federated learning: Challenges, methods, and future directions," *IEEE Signal Process. Mag.*, vol. 37, no. 3, pp. 50–60, 2020.
- [3] H. H. Yang, Z. Liu, T. Q. S. Quek, and H. V. Poor, "Scheduling policies for federated learning in wireless networks," *IEEE Trans. Commun.*, vol. 68, no. 1, pp. 317–333, 2020.
- [4] J. Park, S. Samarakoon, M. Bennis, and M. Debbah, "Wireless network intelligence at the edge," *Proc. IEEE*, vol. 107, no. 11, pp. 2204–2239, Nov. 2019.
- [5] K. Yang, T. Jiang, Y. Shi, and Z. Ding, "Federated learning via over-the-air computation," *IEEE Trans. Wireless Commun.*, vol. 19, no. 3, pp. 2022–2035, 2020.
- [6] G. Zhu, J. Xu, and K. Huang, "Over-the-air computing for 6g - turning air into a computer," *CoRR*, vol. abs/2009.02181, 2020.
- [7] M. M. Amiri and D. Gündüz, "Machine learning at the wireless edge: Distributed stochastic gradient descent over-the-air," *IEEE Trans. Signal Process.*, vol. 68, pp. 2155–2169, Mar. 2020.
- [8] H. Guo, Y. Zhu, H. Ma, V. K. N. Lau, K. Huang, X. Li, H. Nong, and M. Zhou, "Over-the-air aggregation for federated learning: Waveform superposition and prototype validation," *J. of Commun. and Inf. Networks*, vol. 6, no. 4, pp. 429–442, 2021.
- [9] G. Zhu, Y. Wang, and K. Huang, "Broadband analog aggregation for low-latency federated edge learning," *IEEE Trans. Wireless Commun.*, vol. 19, no. 1, pp. 491–506, 2020.
- [10] H. H. Yang, Z. Chen, T. Q. S. Quek, and H. V. Poor, "Revisiting analog over-the-air machine learning: The blessing and curse of interference," *IEEE J. Sel. Topics Signal Process.*, vol. 16, no. 3, pp. 406–419, 2021.
- [11] Z. Zhang, G. Zhu, R. Wang, V. K. N. Lau, and K. Huang, "Turning channel noise into an accelerator for over-the-air principal component analysis," *IEEE Trans. Wireless Commun.*, vol. 21, no. 10, pp. 7926–7941, 2022.
- [12] G. Zhu, J. Xu, K. Huang, and S. Cui, "Over-the-air computing for wireless data aggregation in massive iot," *IEEE Wireless Commun.*, vol. 28, no. 4, pp. 57–65, 2021.
- [13] T. Sery and K. Cohen, "On analog gradient descent learning over multiple access fading channels," *IEEE Trans. Signal Process.*, vol. 68, pp. 2897–2911, 2020.
- [14] J. Duchi, E. Hazan, and Y. Singer, "Adaptive subgradient methods for online learning and stochastic optimization," *J. Mach. Learn. Res.*, vol. 12, no. 61, pp. 2121–2159, 2011.
- [15] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *Int. Conf. Learn. Represent.*, Y. Bengio and Y. LeCun, Eds., 2015.
- [16] S. J. Reddi, Z. Charles, M. Zaheer, Z. Garrett, K. Rush, J. Konečný, S. Kumar, and H. B. McMahan, "Adaptive federated optimization," in *Int. Conf. Learn. Represent.*, 2020.
- [17] R. Ward, X. Wu, and L. Bottou, "Adagrad stepsizes: Sharp convergence over nonconvex landscapes," *J. Mach. Learn. Res.*, vol. 21, no. 1, pp. 9047–9076, 2020.
- [18] S. Mehta, C. Paunwala, and B. Vaidya, "Cnn based traffic sign classification using adam optimizer," in *Proc. IEEE Int. Conf. on Commun. Systems*, IEEE, 2019, pp. 1293–1298.
- [19] G. Cohen, S. Afshar, J. Tapson, and A. van Schaik, "Emnist: Extending mnist to handwritten letters," in *Proc. Int. Jt. Conf. Neural Netw.*, 2017, pp. 2921–2926.
- [20] A. Krizhevsky and G. Hinton, "Learning multiple layers of features from tiny images," Tech. Rep. 0, University of Toronto, Toronto, Ontario, 2009.
- [21] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 770–778.

APPENDIX

A. Proof of Lemma 1

We prove the following inequality by induction:

$$\sum_{t=1}^T \mathbb{E} \left[\frac{(g_j^t)^2}{\sqrt{tv_j^t} + \sqrt{t}\epsilon} \right] \leq \mathbb{E} \left[\sqrt{Tv_j^T} + \sqrt{T}\epsilon^T \right] \quad (14)$$

For $t = 1$

$$\begin{aligned} \mathbb{E} \left[\frac{(g_j^1)^2}{\sqrt{v_j^1} + \epsilon} \right] &= \mathbb{E} \left[\frac{(g_j^1)^2}{\sqrt{(g_j^1)^2} + \epsilon} \right] \\ &\leq \mathbb{E} \left[\sqrt{(g_j^1)^2} + \epsilon \right] \leq 2\mathbb{E} \left[\sqrt{(g_j^1)^2} + \epsilon \right] \end{aligned} \quad (15)$$

Suppose that the conclusion holds when $t = T - 1$, i.e., for any $j \in [d]$,

$$\sum_{t=1}^{T-1} \mathbb{E} \left[\frac{(g_j^t)^2}{\sqrt{tv_j^t} + \sqrt{t}\epsilon} \right] \leq \mathbb{E} \left[\sqrt{(T-1)v_j^{T-1}} + \sqrt{T-1}\epsilon \right]. \quad (16)$$

In addition, combined with the fact that $v_j^T = v_j^{T-1} + (g_j^T)^2$ and $\sqrt{T}\epsilon \geq \sqrt{T-1}\epsilon^{T-1}$, we have

$$\begin{aligned} \sqrt{(T-1)v_j^{T-1}} + \sqrt{T-1}\epsilon &\leq \\ \sqrt{(T-1)v_j^T - (T-1)(g_j^T)^2} + \sqrt{T}\epsilon \end{aligned} \quad (17)$$

Using $\sqrt{x-c} \leq \sqrt{x} - \frac{c}{2\sqrt{x}}$ ($x > c$), we have

$$\begin{aligned} &\sqrt{(T-1)v_j^T - (T-1)(g_j^T)^2} + \sqrt{T}\epsilon \\ &\leq \sqrt{(T-1)v_j^T} - \frac{(T-1)(g_j^T)^2}{2\sqrt{(T-1)v_j^T}} + \sqrt{T}\epsilon \\ &\leq \sqrt{Tv_j^T} - \frac{(T-1)(g_j^T)^2}{2(\sqrt{Tv_j^T} + \sqrt{T}\epsilon)} + \sqrt{T}\epsilon \end{aligned} \quad (18)$$

Hence

$$\begin{aligned} &\sum_{t=1}^T \mathbb{E} \left[\frac{(g_j^t)^2}{\sqrt{tv_j^t} + \sqrt{t}\epsilon} \right] \\ &\leq \mathbb{E} \left[\sqrt{Tv_j^T} - \frac{(T-1)(g_j^T)^2}{2(\sqrt{Tv_j^T} + \sqrt{T}\epsilon)} + \sqrt{T}\epsilon \right] \\ &\quad + \mathbb{E} \left[\frac{(g_j^T)^2}{\sqrt{Tv_j^T} + \sqrt{T}\epsilon} \right] \\ &\leq \mathbb{E} \left[\sqrt{Tv_j^T} + \sqrt{T}\epsilon \right] + \left(1 - \frac{T-1}{2} \right) \mathbb{E} \left[\frac{(g_j^T)^2}{\sqrt{Tv_j^T} + \sqrt{T}\epsilon} \right] \\ &\leq \mathbb{E} \left[\sqrt{Tv_j^T} + \sqrt{T}\epsilon \right] \end{aligned} \quad (19)$$

B. Proof of Theorem 1

We employ the norm as

$$\|\cdot\|_M = \sqrt{\langle \cdot, M \cdot \rangle}$$

for symmetric and positive definite $M \in \mathbb{R}^d$. If $M \in \mathbb{R}^d$ and has entries being non-negative, then diagonalize it.

To simplify the notation, we denote $\sqrt{\nu_t} + \varepsilon_t$ as A_t , and $\eta^t = \frac{\eta}{\sqrt{t}}$, which means η^t is a decay factor

Based on the updating rule, we have

$$\begin{aligned} &\|w^{t+1} - w^*\|_{A^t}^2 \\ &= \|w^t - \eta^t (A^t)^{-1} g^t - w^*\|_{A^t}^2 \\ &= \|w^t - w^*\|_{A^t}^2 + \|\eta^t (A^t)^{-1} g^t\|_{A^t}^2 - 2\langle w^t - w^*, \eta^t g^t \rangle \\ &= \|w^t - w^*\|_{A^t}^2 - 2\eta^t \langle g^t, w^t - w^* \rangle + (\eta^t)^2 \langle g^t, (A^t)^{-1} g^t \rangle \end{aligned} \quad (20)$$

Rearranging terms gives

$$\begin{aligned} \langle g^t, w^t - w^* \rangle &= \frac{\|w^t - w^*\|_{A^t}^2 - \|w^{t+1} - w^*\|_{A^t}^2}{2\eta^t} \\ &\quad + \frac{\eta^t}{2} \langle g^t, (A^t)^{-1} g^t \rangle. \end{aligned} \quad (21)$$

Take expectation on both sides conditioned on w^t

$$\begin{aligned} \langle \nabla F(w^t), w^t - w^* \rangle &= \frac{\mathbb{E}_t [\|w^t - w^*\|_{A^t}^2] - \mathbb{E}_t [\|w^{t+1} - w^*\|_{A^t}^2]}{2\eta^t} \\ &\quad + \frac{\eta^t}{2} \mathbb{E}_t [\langle g^t, (A^t)^{-1} g^t \rangle], \end{aligned} \quad (22)$$

where we have used the fact that N is a zero-mean Gaussian variable independent of g^t, w^t . Taking expectation on both sides and using the convexity of $F(\cdot)$:

$$\begin{aligned} &\mathbb{E} [F(w^t)] - F(w^*) \\ &\leq \frac{\mathbb{E} [\|w^t - w^*\|_{A^t}^2] - \mathbb{E} [\|w^{t+1} - w^*\|_{A^t}^2]}{2\eta^t} \\ &\quad + \frac{\eta^t}{2} \mathbb{E} [\langle g^t, (A^t)^{-1} g^t \rangle] \end{aligned} \quad (23)$$

Applying telescope sum, we have

$$\begin{aligned} &\sum_{t=1}^T (\mathbb{E} [F(w^t)] - F(w^*)) \leq \frac{\|w^1 - w^*\|_{A^1}^2}{2\eta_1} \\ &\quad + \sum_{t=2}^T \left(\frac{\mathbb{E} [\|w^t - w^*\|_{A^t}^2]}{2\eta^t} - \frac{\mathbb{E} [\|w^t - w^*\|_{A^{t-1}}^2]}{2\eta^{t-1}} \right) \\ &\quad + \sum_{t=1}^T \frac{\eta^t}{2} \mathbb{E} [\langle g^t, (A^t)^{-1} g^t \rangle] \end{aligned} \quad (24)$$

Using $\eta^t = \frac{\eta}{\sqrt{t}}$, we now have

$$\begin{aligned}
& \sum_{t=1}^T (\mathbb{E}[F(w^t)] - F(w^*)) \\
& \leq \underbrace{\frac{\|w^1 - w^*\|_{A^1}^2}{2\eta} + \sum_{t=2}^T \frac{\mathbb{E}[\|w^t - w^*\|_{\sqrt{t}A^t - \sqrt{t-1}A^{t-1}}^2]}{2\eta}}_{T_1} + \\
& \quad \underbrace{\sum_{t=1}^T \frac{\eta}{2\sqrt{t}} \mathbb{E}[\langle g^t, (A^t)^{-1} g^t \rangle]}_{T_2}
\end{aligned} \tag{25}$$

We first bound T_1 . Based on the relations between v^t and v^{t-1} , apparently for any j, t we have

$$\sqrt{t}A_j^t \geq \sqrt{t-1}A_j^{t-1}. \tag{26}$$

Hence,

$$\begin{aligned}
& \mathbb{E} \left[\sum_{t=2}^T \|w^t - w^*\|_{\sqrt{t}A^t - \sqrt{t-1}A^{t-1}}^2 \right] \\
& = \mathbb{E} \left[\sum_{t=2}^T \sum_{j=1}^d (w_j^t - w_j^*)^2 \left(\sqrt{tv_j^t} + \sqrt{t}\epsilon - \sqrt{(t-1)v_j^{t-1}} - \sqrt{t-1}\epsilon \right) \right] \\
& = \mathbb{E} \left[\sum_{j=1}^d \sum_{t=2}^T (w_j^t - w_j^*)^2 \left(\sqrt{tv_j^t} + \sqrt{t}\epsilon - \sqrt{(t-1)v_j^{t-1}} - \sqrt{t-1}\epsilon \right) \right] \\
& \leq \mathbb{E} \left[\sum_{j=1}^d D^2 \sum_{t=2}^T \left(\sqrt{tv_j^t} + \sqrt{t}\epsilon - \sqrt{(t-1)v_j^{t-1}} - \sqrt{t-1}\epsilon \right) \right] \\
& = \mathbb{E} \left[\sum_{j=1}^d D^2 \left(\sqrt{Tv_j^T} + \sqrt{T}\epsilon - \sqrt{v_j^1} - \epsilon \right) \right]
\end{aligned} \tag{27}$$

We next bound T_2 . According to Lemma 1, we have

$$\sum_{t=1}^T \mathbb{E} \left[\frac{(g_j^t)^2}{\sqrt{tv_j^t} + \sqrt{t}\epsilon} \right] \leq \mathbb{E} \left[\sqrt{Tv_j^T} + \sqrt{T}\epsilon \right] \tag{28}$$

so we can bound T_2 as follows.

$$\begin{aligned}
T_2 & = \mathbb{E} \left[\sum_{t=1}^T \frac{\eta^t}{2} \sum_{j=1}^d \frac{(g_j^t)^2}{\sqrt{tv_j^t} + \epsilon} \right] = \frac{\eta}{2} \mathbb{E} \left[\sum_{t=1}^T \sum_{j=1}^d \frac{(g_j^t)^2}{\sqrt{tv_j^t} + \sqrt{t}\epsilon} \right] \\
& \leq \frac{\eta}{2} \sum_{j=1}^d \mathbb{E} \left[\sqrt{Tv_j^T} + \sqrt{T}\epsilon \right].
\end{aligned} \tag{29}$$

Noting that

$$\frac{\|w^1 - w^*\|_{A^1}^2}{2\eta} \leq \frac{D^2}{2\eta} \sum_{j=1}^d \left(\sqrt{v_j^1} + \epsilon \right), \tag{30}$$

combined with the bounds of T_1, T_2 yields

$$\sum_{t=1}^T (\mathbb{E}[F(w^t)] - F(w^*)) \leq \left(\frac{D^2}{2\eta} + \frac{\eta}{2} \right) \sum_{j=1}^d \mathbb{E} \left[\sqrt{Tv_j^T} + \sqrt{T}\epsilon \right], \tag{31}$$

which implies that

$$\begin{aligned}
\min_{t \in [T]} \mathbb{E}[F(w^t)] - F(w^*) & \leq \left(\frac{D^2}{2\eta} + \frac{\eta}{2} \right) \frac{1}{T} \sum_{j=1}^d \mathbb{E} \left[\sqrt{Tv_j^T} + \sqrt{T}\epsilon \right] \\
& \leq \left(\frac{D^2}{2\eta} + \frac{\eta}{2} \right) \frac{1}{\sqrt{T}} \sum_{j=1}^d \mathbb{E} \left[\sqrt{v_j^T} + \epsilon \right]
\end{aligned} \tag{32}$$

Further note that

$$\begin{aligned}
v_T & = v_{T-1} + g_T^2 \\
& = v_{T-2} + g_{T-1}^2 + g_T^2 \quad \dots
\end{aligned} \tag{33}$$

Since

$$g_k = \frac{1}{N} \sum_{n=1}^N h_{n,k} \cdot \nabla g_n(\omega_k) + \xi_k \tag{34}$$

We average out randomness from fading and noise from

$$\begin{aligned}
& \mathbb{E} \left[\sqrt{Tv_j^T} \right] \\
& \text{since } f(x) = \sqrt{x} \text{ is convex, we have}
\end{aligned}$$

$$\mathbb{E}_{h,\xi} \left[\sqrt{Tv_T} \right] \leq \mathbb{E} \left[\sqrt{T\mathbb{E}_{h,\xi}(v_T)} \right] \tag{35}$$

So

$$\begin{aligned}
\mathbb{E} \left[(g_k^j)^2 \right] & = \mathbb{E} \left[\left(\frac{1}{N} \sum_{n=1}^N h_{n,k} (\nabla f_n(w_k))_j + \xi_k^j \right)^2 \right] \\
& \leq (\mu_c^2 + \sigma_c^2) \left(\frac{1}{N} \sum_{n=1}^N \nabla F(w_k)_j^2 \right) + \sigma_n^2 \\
& = (\mu_c^2 + \sigma_c^2) (\nabla F(w_k))_j^2 + \sigma_n^2
\end{aligned} \tag{36}$$

We denote an auxiliary factor as

$$\tilde{v}_t = \tilde{v}_{t-1} + (\nabla f(w_t))^2 \tag{37}$$

where $f(\omega) = \frac{1}{N} \sum_{n=1}^N f_n(w)$. Then

$$\begin{aligned}
& \min_{t \in [T]} \mathbb{E}[F(w^t)] - F(w^*) \\
& \leq \left(\frac{D^2}{2\eta} + \frac{\eta}{2} \right) \frac{1}{T} \sum_{j=1}^d \mathbb{E} \left[\sqrt{Tv_j^T} + \sqrt{T}\epsilon \right] \\
& \leq \left(\frac{D^2}{2\eta} + \frac{\eta}{2} \right) \frac{1}{\sqrt{T}} \sum_{j=1}^d \mathbb{E} \left[\sqrt{(\mu_c^2 + \sigma_c^2) \tilde{v}_j^T} + \epsilon \right]
\end{aligned} \tag{38}$$