

Over-the-Air Federated Learning: Rethinking Edge AI Through Signal Processing

Seyed Mohammad Azimi-Abarghouyi, *Member, IEEE*, Carlo Fischione, *Fellow, IEEE*, and Kaibin Huang, *Fellow, IEEE*

ABSTRACT

Over-the-Air Federated Learning (AirFL) is an emerging paradigm that tightly integrates wireless signal processing and distributed machine learning to enable scalable AI at the network edge. By leveraging the superposition property of wireless signals, AirFL performs communication and model aggregation of the learning process simultaneously, significantly reducing latency, bandwidth, and energy consumption. This article offers a tutorial treatment of AirFL, presenting a novel classification into three design approaches: CSIT-aware, blind, and weighted AirFL. We provide a comprehensive guide to theoretical foundations, performance analysis, complexity considerations, practical limitations, and prospective research directions.

I. INTRODUCTION

The convergence of edge AI and wireless communications is shaping the next generation of intelligent, networked systems. At the core of this paradigm lies Federated Learning (FL), a distributed machine learning framework in which devices such as smartphones, IoT sensors, and autonomous vehicles collaboratively train a global model without sharing their raw data [1]. The fundamental process involves each device performing local model updates based on its private data and then transmitting these updates to a central server, where aggregation takes place to refine the global model. This cycle is repeated across multiple rounds until convergence.

Originally introduced to preserve privacy and data locality in decentralized environments, FL has rapidly evolved to address broader challenges inherent to large-scale networked systems. However, its real-world deployment remains hindered by a critical limitation: communication overhead. In traditional FL setups,

S. M. Azimi-Abarghouyi and C. Fischione are with the School of Electrical Engineering and Computer Science, KTH Royal Institute of Technology, Stockholm, Sweden (Emails: {seyaa, carlofi}@kth.se). K. Huang is with the Department of Electrical and Electronic Engineering, The University of Hong Kong, Hong Kong SAR, China (Email: huangkb@hku.hk).

each device must reliably send high-dimensional model updates to the central server at every round. These frequent transmissions often occur over limited-bandwidth, energy-constrained wireless links, resulting in a significant mismatch between the communication-intensive nature of FL and the capacity of existing wireless networks—especially as the number of participating devices grows.

To overcome this bottleneck, Over-the-Air Computation (AirComp) has emerged as a transformative solution [2], [3]. Unlike conventional methods that recover individual messages, AirComp exploits the superposition property of wireless signals to compute desired functions—such as sums—directly over the multiple-access channel. AirComp should thus be viewed as one component within the broader problem of FL, where its function-aggregation capability is leveraged as a building block toward a larger objective. This forms the foundation of Over-the-Air Federated Learning (AirFL), a new class of FL systems where all devices transmit simultaneously, and the server receives a naturally aggregated signal. By collapsing communication and aggregation into a single step, AirFL significantly reduces latency, bandwidth consumption, and energy usage. This makes it a highly promising approach for achieving scalable, efficient FL in resource-constrained wireless environments. However, realizing this potential requires machine-learning-aware wireless signal-processing techniques—capable of handling statistical and computational heterogeneity across devices and of providing unbiased (or controllably biased) aggregation while ensuring privacy, fairness, and convergence.

The benefits of AirFL come with significant challenges. Unlike orthogonal schemes, AirFL is highly sensitive to channel impairments, synchronization errors, and hardware limitations. In particular, the accuracy of signal aggregation depends on the precise alignment of transmitted signals. One early approach to address this relies on accurate Channel State Information at the Transmitter (CSIT) and power control to compensate for channel gains and phase shifts. However, acquiring and maintaining CSIT in practical wireless environments—especially in large-scale or rapidly varying systems—remains a major challenge. These difficulties have sparked growing interest in methods that require only local CSIT [4]. Nonetheless, approaches based on global CSIT [5] continue to be explored due to their potential for achieving optimal performance.

More recently, CSIT-free approaches—most notably blind AirFL [6] and weighted AirFL [7]—have emerged as promising practical alternatives. These approaches eliminate the reliance on channel estimation and power control at transmitters through partial phase compensation, leveraging massive MIMO and high-dimensional processing at the server, or by applying adaptive weighting within aggregation. By offloading complexity from edge devices and avoiding fine synchronization and feedback overhead, these solutions reduce device complexity—albeit at the cost of higher processing demands on the server or a potential risk of aggregation bias (objective drift).

In this article, we categorize AirFL schemes into three main classes: CSIT-aware, blind, and weighted AirFL, each representing a distinct design philosophy centered on one of three key system parameters: transmission power control at the devices, equalization at the server, and aggregation weighting. The representative methods in each class reflect the corresponding design principle and collectively capture the dominant directions that have shaped the existing AirFL literature.

A. Related Work and Contributions of this Tutorial

Previous tutorials have typically focused on either FL [8], [9] or AirComp [10], [11]. In contrast, this tutorial presents a unified and integrated treatment of their combination, namely AirFL, and delivers it as a comprehensive and detailed lecture note. Unlike the general AirComp perspective, which typically remains agnostic to learning objectives and relies mainly on mean squared error (MSE) as the sole metric, this tutorial focuses on the communication-learning interface, calling for a fundamental rethinking of how to fully integrate learning-oriented objectives, metrics, and influencing factors. It covers theoretical foundations, performance analysis, signal processing techniques, system designs, practical challenges, and recent developments.

While there are a few brief overview papers on AirFL [12], [13], as well as broader surveys on distributed learning over wireless networks that touch on AirFL [14], [15], their discussions remain at a high level and primarily focus on the local CSIT-based approach. However, this approach lacks the specialized processing capabilities found in more recent and advanced methods, which are thoroughly explored and clearly classified in this tutorial.

We begin by revisiting the principles of AirComp and its role in aggregation in FL. We then examine the fundamentals of FL, highlighting the standard learning process and its integration with wireless systems. Next, we delve into the first class of CSIT-aware schemes, analyzing both local and global approaches, along with their convergence behavior and optimization frameworks. This is followed by an in-depth discussion of the second class, namely blind schemes, including convergence guarantees under heavy-tailed interference and the antenna requirements for accurate learning. We conclude with the third class: weighted schemes, examined under both homogeneous and heterogeneous computational settings. To systematically compare these approaches and highlight their trade-offs, we introduce a common system model and analytical framework that enables consistent evaluation and meaningful insights across all classes.

II. OVER-THE-AIR COMPUTATION FUNDAMENTALS

This section presents the fundamental principles of AirComp over the standard multiple access channel (MAC).

A. Signal Model over the MAC

Consider a wireless system comprising K single-antenna transmitting devices and a receiver with M antennas. This is the most widely adopted wireless model, as devices typically have limited hardware capabilities and physical constraints that restrict them to a single antenna. Each transmitter k sends a scaled version of its data symbol x_k using a complex precoding scalar

$$p_k = |p_k|e^{j\angle p_k}, \quad (1)$$

which controls both the transmission power and phase. The choice of p_k , referred to as *power control*, must satisfy either the average power constraint $\mathbb{E}[|p_k|^2] \leq P$ or the instantaneous constraint $|p_k|^2 \leq P$. The received signal vector $\mathbf{y} \in \mathbb{C}^M$ at the receiver is then given by

$$\mathbf{y} = \sum_{k=1}^K \mathbf{h}_k p_k x_k + \mathbf{z}. \quad (2)$$

Here, $\mathbf{h}_k = [h_{k,1}, \dots, h_{k,M}]^T$ denotes the channel vector from device k to the receiver, with each element

$$h_{k,m} = |h_{k,m}|e^{j\angle h_{k,m}}, \quad (3)$$

representing the channel coefficient from device k to antenna m , comprising channel gain $|h_{k,m}|$ and phase $\angle h_{k,m}$. The channels are assumed to remain constant during a given transmission but may vary across different transmissions. The vector $\mathbf{z} \sim \mathcal{CN}(0, \sigma_z^2 \mathbf{I})$ models complex additive white Gaussian noise (AWGN) at the receiver. On the other hand, to estimate a desired function of the transmitted data, the receiver applies an equalization vector $\mathbf{b} \in \mathbb{C}^M$, yielding the scalar output

$$f = \mathbf{b}^H \mathbf{y} = \mathbf{b}^H \left(\sum_{k=1}^K \mathbf{h}_k p_k x_k + \mathbf{z} \right). \quad (4)$$

B. Synchronization Aspects

The baseband MAC model in (2) assumes that the transmitted waveforms from all devices overlap coherently at the receiver so that their joint contribution can be represented by complex scalar channel coefficients $\{h_{k,m}\}$. To ensure the validity of this model within an AirComp block, three distinct forms of synchronization are relevant:

1) Clock synchronization (frame and symbol timing): All devices must begin transmission within a common timing window such that the receiver samples their symbols at consistent time instants. Small

residual timing offsets within this window can be absorbed into the effective channel coefficients, provided they do not introduce significant inter-symbol interference.

2) Time alignment of received signals (waveform overlap): For superposition-based computation to hold, all transmitted symbols must overlap at the receiver within the same symbol interval. This alignment depends on both clock synchronization and the propagation delays of the channels. When the residual delay spread is small relative to the symbol duration, its effect is adequately captured by the complex gains $\{h_{k,m}\}$ in (2); otherwise, receive-side equalization or OFDM with cyclic prefix may be required.

3) Carrier-frequency and carrier-phase synchronization: The carrier frequencies of all devices must be sufficiently close to that of the receiver so that residual carrier-frequency offsets and phase noise evolve slowly over one AirComp block. Schemes that employ phase-sensitive power control require stricter, *fine* synchronization—namely, very small carrier-frequency offsets and minimal phase drift. In contrast, phase-agnostic schemes operate under only *coarse* synchronization, as their residual timing and carrier offsets can be absorbed into the random channel coefficients and handled statistically or through equalization.

Throughout the rest of this article, we refer to *fine synchronization* as symbol-level alignment together with tight carrier-frequency and carrier-phase synchronization, and to *coarse synchronization* as frame-level alignment with moderately stable carriers whose residual offsets are absorbed into the effective channels.

C. Nomographic Function Computation

More generally, AirComp aims to compute a function of the form

$$f = \Phi \left(\sum_{k=1}^K \mathbf{h}_k \Psi_k(x_k) + \mathbf{z} \right), \quad (5)$$

where $\Psi_k(\cdot)$ is the pre-processing function applied at device k and $\Phi(\cdot)$ is a post-processing function applied at the receiver. The objective of AirComp is to make (5) as close as possible to a nomographic target function F of the form

$$F = \Theta \left(\sum_{k=1}^K \Lambda_k(x_k) \right), \quad (6)$$

where $\Lambda_k(\cdot)$ and $\Theta(\cdot)$ are known functions. Functions such as (6) are highly expressive and appears in many practical applications across signal processing, distributed computation, and machine learning. Examples include:

- Weighted Sum: $\sum_{k=1}^K \alpha_k x_k$
- Majority Vote: $\text{sign} \left(\sum_{k=1}^K \text{sign}(x_k) \right)$
- Polynomial: $\sum_{k=1}^K a_k x_k^{k-1}$

- P-norm: $\left(\sum_{k=1}^K |x_k|^p\right)^{1/p}$

D. Performance Metric and Optimization Objective

A common metric to evaluate the accuracy of AirComp is the mean squared error (MSE) between the received function f and the desired function value F as

$$\text{MSE} = \mathbb{E} [|f - F|^2] = \mathbb{E} \left[\left| \Phi \left(\sum_{k=1}^K \mathbf{h}_k \Psi_k(x_k) + \mathbf{z} \right) - \Theta \left(\sum_{k=1}^K \Lambda_k(x_k) \right) \right|^2 \right]. \quad (7)$$

The goal is to design the pre-processing functions Ψ_k and post-processing function Φ to minimize this error as $\min_{\Phi, \Psi_k} \text{MSE}$. This optimization typically depends on the availability of channel state information (CSI), i.e., knowledge of \mathbf{h}_k for all k in (7). Depending on the design approach, CSI may be required at the transmitter side (CSIT) to compute pre-processing functions, or at the receiver side (CSIR) for post-processing. In most systems, CSI is first estimated at the receiver—serving as a central node—by collecting training sequences from all transmitters, thereby acquiring CSIR. To enable CSIT, the receiver must then feed this information back to the transmitters. However, acquiring CSIT introduces significant signaling overhead and depends on the channel remaining stable; otherwise, the CSIT may become outdated by the time it is used for transmission pre-processing. Once CSIT is available, each transmitter can apply the appropriate pre-processing and transmit its data. In contrast, if the design relies solely on CSIR without pre-processing, devices can blindly transmit their data immediately after sending training sequences, without waiting for feedback. In this case, the receiver performs all necessary post-processing using the available CSIR.

The existing studies on AirComp and its applications focus on linear pre-processing and post-processing functions, specifically power control and equalization, as discussed in Subsection II-A [4]–[7], [16], [17].

E. Challenges and Practical Considerations

AirComp presents several practical challenges:

- It requires the synchronization aspects in Subsection II-B to the extent that the baseband MAC model in (2) remains valid; any residual misalignment appears as additional distortion.
- Perfect computation is generally unachievable due to noise and the fact that channel coefficients can take arbitrary values, leading to unavoidable residual errors.
- It entails several hardware constraints, including the need for accurate channel estimation and analog modulation to freely shape the transmitted signals. Although the original design is based on analog communication, recent advancements have proposed digital communication as a viable

alternative [18], [19]. They take advantage of digital chips already embedded in many existing devices. Nevertheless, this aspect is beyond the scope of this article.

- It is inherently restricted to computing certain classes of functions, with prior research primarily concentrating on summation.

Despite these limitations, AirComp offers significant benefits:

- It drastically reduces communication overhead by enabling direct function computation through simultaneous transmission within a single resource block, eliminating the need for full data recovery and thereby achieving highly reduced bandwidth requirements, lower energy consumption, and minimized latency.
- It scales efficiently with the number of devices, making it suitable for large-scale systems.
- It opens up the potential for creating a "wireless distributed computing platform"—a kind of *computer over the air*.

III. FEDERATED LEARNING FUNDAMENTALS

Device k possesses its own local (private) dataset \mathcal{D}_k . The learning model is parametrized by the model vector $\mathbf{w} = [w_1, \dots, w_s]^\top \in \mathbb{R}^{s \times 1}$, where s is the model size. Then, the local loss function of the model vector \mathbf{w} on \mathcal{D}_k is

$$F_k(\mathbf{w}) = \frac{1}{D_k} \sum_{\xi_i \in \mathcal{D}_k} \ell(\mathbf{w}, \xi_i), \quad (8)$$

where $D_k = |\mathcal{D}_k|$ is the size of the dataset and the function $\ell(\mathbf{w}, \xi_i)$ represents the sample-wise loss, measuring the prediction error of \mathbf{w} on the sample ξ_i . Following this, the global loss function applied to all distributed datasets, denoted as $\cup_{k=1}^K \mathcal{D}_k$, is

$$F(\mathbf{w}) = \frac{1}{\sum_{k=1}^K D_k} \sum_{k=1}^K D_k F_k(\mathbf{w}). \quad (9)$$

The goal of the training procedure is to discover an optimal parameter vector \mathbf{w} that minimizes $F(\mathbf{w})$, expressed as

$$\mathbf{w}^* = \min_{\mathbf{w}} F(\mathbf{w}). \quad (10)$$

A widely used FL algorithm for solving (10) is FedAvg [1], which serves as the foundational component of most AirFL schemes, as follows. Consider a particular round $t \in \{0, \dots, T-1\}$, where T denotes the number of rounds. In this round, each device k first updates its own learning model via τ local training steps, each based on a randomly sampled mini-batch ξ_k^i with size B drawn from \mathcal{D}_k , as

$$\mathbf{w}_k^{t,i+1} = \mathbf{w}_k^{t,i} - \mu \nabla F_k(\mathbf{w}_k^{t,i}, \xi_k^{t,i}), \forall i \in \{0, \dots, \tau-1\}, \quad (11)$$

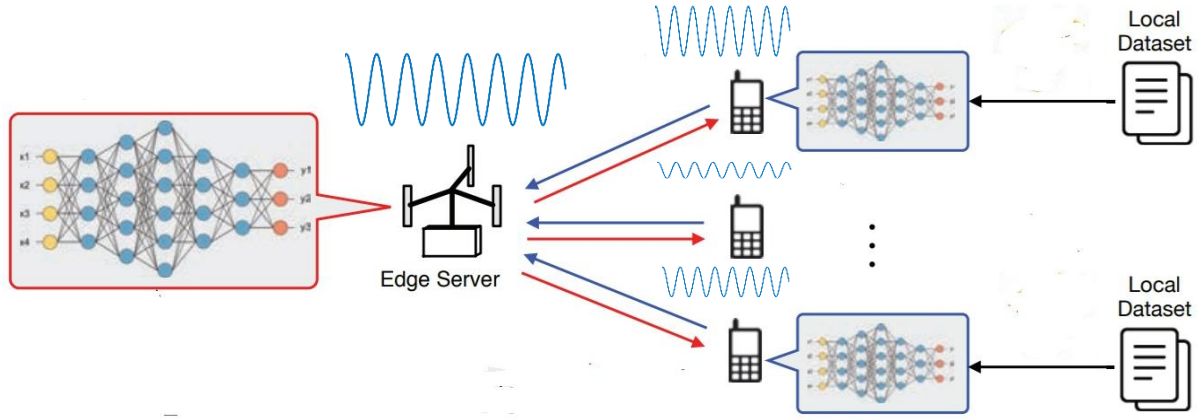


Fig. 1: Standard AirFL setup.

where μ is the learning rate. Then, each device k uploads the local model $\mathbf{w}_k^t = \mathbf{w}_k^{t,\tau}$ to the server for aggregation. As the ideal aggregation, the global gradient can be obtained as an average of model parameters from all the devices with equal weighting as

$$\mathbf{w}_G^{t+1} = \frac{1}{K} \sum_{k=1}^K \mathbf{w}_k^t. \quad (12)$$

Next, the server broadcasts the obtained global model \mathbf{w}_G^{t+1} to the devices, based on which each device k updates its initial state for the next round as $\mathbf{w}_k^{t+1,0} = \mathbf{w}_G^{t+1}, \forall k$.

Thus, a prominent application of AirComp lies in FL, known as over-the-air FL (AirFL) with a setup shown in Fig. 1, where the objective is to compute the following linear nomographic function:

$$F = \frac{1}{K} \sum_{k=1}^K x_k. \quad (13)$$

This direct wireless aggregation approach significantly reduces both latency and bandwidth consumption compared to traditional orthogonal transmission schemes, such as TDMA and FDMA, which require a separate resource block for each device. In these schemes, the server must wait to receive individual model updates from all devices before performing the aggregation.

The following sections present several well-recognized approaches for enabling AirFL, each focusing on a distinct key design element.

IV. CSIT-AWARE OVER-THE-AIR FEDERATED LEARNING

This approach, referred to as CSIT AirFL, requires full channel knowledge for its design—not only CSIR, but also CSIT. Two variants arise depending on the scope of CSIT: local CSIT, where each device is aware only of its own channel and independently aligns its transmission in a decentralized and typically

suboptimal manner; and global CSIT, where each device has access to the complete set of channel information across the system. With this comprehensive knowledge, the design becomes centralized, potentially achieving optimal performance. In both cases, CSIT is leveraged to emphasize *transmission power control* as a key design element.

A. Local CSIT AirFL

Since this approach relies on local CSIT, it is inherently limited to single-antenna server scenarios and has not been extended to settings involving multi-antenna servers. Originally introduced in [16], this approach has since been widely adopted in subsequent AirFL studies, including [20]–[23], where it has been extended to support digital modulations through one-bit aggregation [20], large-scale cellular deployments [21], hierarchical multi-cluster networks [23], and principal-component-based learning [22]. It comprises the following key components:

Power Control: The following power control at each device k employs direct channel compensation to exclusively align transmissions with the aggregation function (13).

$$p_k \propto \frac{1}{h_k} = \frac{1}{|h_k|} e^{-j\angle h_k}, \quad (14)$$

where h_k is the channel between device k and the server. However, this results in extremely high transmit power when the channel gain $|h_k|$ is small. To address this, *truncated power control*, where only devices with sufficiently strong channels transmit, is used as

$$p_k = \begin{cases} \frac{\sqrt{\rho}}{|h_k|} e^{-j\angle h_k}, & \text{if } |h_k|^2 \geq \theta, \\ 0, & \text{otherwise,} \end{cases} \quad (15)$$

where ρ is a denormalizing factor and θ denotes a threshold. Accordingly, the active device set, which includes the transmitting devices, is defined as

$$\mathcal{S} = \{k \mid |h_k|^2 \geq \theta\}. \quad (16)$$

Thus, *device selection* is inherent in this approach.

Aggregation: The single-antenna server estimates the aggregation function as

$$\mathbf{w}_G = \frac{\mathbf{y}}{\sqrt{\rho}|\mathcal{S}|}. \quad (17)$$

The result is unbiased and can be written as

$$\mathbf{w}_G = \underbrace{\frac{1}{|\mathcal{S}|} \sum_{k \in \mathcal{S}} \mathbf{w}_k}_{\text{Desired Aggregation}} + \underbrace{\frac{\mathbf{z}}{\sqrt{\rho}|\mathcal{S}|}}_{\text{Error from AWGN}}. \quad (18)$$

The factor ρ controls the error in recovering the aggregation $\frac{1}{|\mathcal{S}|} \sum_{k \in \mathcal{S}} \mathbf{w}_k$, and maximizing it leads to improved performance.

Denormalizing: Assuming Rayleigh fading, i.e., $|h_k|^2 \sim \exp(1)$, the factor ρ is chosen under the power constraint P as

$$\mathbb{E}[|p_k|^2] = \rho \mathbb{E}\left[\frac{1}{|h_k|^2} \mid |h_k|^2 \geq \theta\right] = \rho \text{Ei}(\theta) = P. \quad (19)$$

Thus,

$$\rho = \frac{P}{\text{Ei}(\theta)}, \quad (20)$$

where $\text{Ei}(x) = \int_x^\infty \frac{e^{-t}}{t} dt$ is the exponential integral function.

Limitations:

- It requires fine synchronization as defined in Subsection II-B: symbol-level time alignment and tight carrier-frequency/phase synchronization, since the phase-sensitive channel inversion in (15) is highly sensitive to residual offsets.
- It lacks spatial diversity gains due to single-antenna server.
- It requires careful device selection and consequently limits participation, being highly dependent on channel conditions, which may introduce bias and adversely affect learning fairness; in extreme cases, no device may participate, leading to learning failure.
- It overlooks the learning aspect in its design.

B. Global CSIT AirFL

This approach adopts a centralized design and leverages an arbitrary number of antennas at the server. While it was originally developed in [5], subsequent AirFL studies [24]–[26] have continued to build on its framework—either by reducing its complexity [24], [26] or by extending it to other wireless settings, such as systems incorporating intelligent reflecting surface (IRS) [25]. The key components of this approach are as follows:

Normalization: Each device normalizes its model vector as

$$\bar{\mathbf{w}}_k = \frac{\mathbf{w}_k - \eta_k \mathbf{1}}{\sigma_k}, \quad (21)$$

where

$$\eta_k = \frac{1}{s} \sum_{i=1}^s w_{k,i}, \quad \sigma_k^2 = \frac{1}{s} \sum_{i=1}^s (w_{k,i} - \eta_k)^2,$$

and transmits $p_k \bar{\mathbf{w}}_k$. The pair (η_k, σ_k) is also shared with the server.

Normalizing the model parameters provides two key benefits for the subsequent aggregation estimator. First, zero-mean entries ensure that the estimator is unbiased. Second, unit-variance entries make the power of the estimation error independent of the specific values of the model parameters. This normalization ultimately enables MSE-based design.

Aggregation: The multiple-antenna server computes its estimate based on the received signal \mathbf{y} as

$$\mathbf{w}_G = \bar{\sigma} \frac{\mathbf{b}^H \mathbf{y}}{\sqrt{\rho} |\mathcal{S}|} + \bar{\eta} \mathbf{1}, \quad (22)$$

where $|\mathcal{S}|$ is the number of active devices in this approach, ρ is the denormalizing factor, and

$$\bar{\sigma} = \frac{1}{|\mathcal{S}|} \sum_{k \in \mathcal{S}} \sigma_k, \quad \bar{\eta} = \frac{1}{|\mathcal{S}|} \sum_{k \in \mathcal{S}} \eta_k,$$

represent the averages over the active devices, ensuring an unbiased estimation. Then, (22) can be written as

$$\mathbf{w}_G = \underbrace{\frac{1}{|\mathcal{S}|} \sum_{k \in \mathcal{S}} \mathbf{w}_k}_{\text{Desired Aggregation}} + \underbrace{\frac{\bar{\sigma}}{|\mathcal{S}|} \sum_{k \in \mathcal{S}} \left(\frac{\mathbf{b}^H \mathbf{h}_k p_k}{\sqrt{\rho}} - 1 \right) \bar{\mathbf{w}}_k}_{\text{Error from Channel Misalignment with Aggregation Weights}} + \underbrace{\frac{\bar{\sigma} \mathbf{b}^H \mathbf{z}}{\sqrt{\rho} |\mathcal{S}|}}_{\text{Error from AWGN}}. \quad (23)$$

Power Control and Denormalizing: To minimize MSE caused by the error terms in (23), the optimal p_k is [5]

$$p_k = \sqrt{\rho} \frac{(\mathbf{b}^H \mathbf{h}_k)^H}{|\mathbf{b}^H \mathbf{h}_k|^2}. \quad (24)$$

With power constraint $|p_k|^2 \leq P$, the optimal ρ is [5]

$$\rho = \min_{k \in \mathcal{S}} P |\mathbf{b}^H \mathbf{h}_k|^2, \quad (25)$$

which results to

$$\text{MSE} = \frac{\sigma_z^2}{P} \max_{k \in \mathcal{S}} \frac{\|\mathbf{b}\|^2}{|\mathbf{b}^H \mathbf{h}_k|^2}. \quad (26)$$

If the sole objective were to compute the most accurate aggregation, one could simply minimize the MSE with respect to the remaining unknowns—the set of active devices \mathcal{S} and the equalization vector \mathbf{b} . However, the primary objective is to enable learning; thus, the focus shifts to evaluating the convergence performance of the learning process.

Convergence Analysis: The following key widely-used assumptions are made to facilitate the analysis.

Assumption 1 (Lipschitz-Continuous Gradient): The gradient of the loss function $F(\mathbf{w})$, as represented in (9), is characterized by Lipschitz continuity with a non-negative constant $L > 0$. This implies that for any pair of model vectors \mathbf{w}_1 and \mathbf{w}_2 , we have

$$F(\mathbf{w}_2) \leq F(\mathbf{w}_1) + \nabla F(\mathbf{w}_1)^T (\mathbf{w}_2 - \mathbf{w}_1) + \frac{L}{2} \|\mathbf{w}_2 - \mathbf{w}_1\|^2. \quad (27)$$

Assumption 2 (Gradient Variance Bound): The local stochastic gradient estimate for device k at \mathbf{w}_k , using a mini-batch ξ_k with size B , is an unbiased estimate of the ground-truth gradient $\nabla F(\mathbf{w}_k)$ with bounded variance

$$\mathbb{E} [\|\nabla F_k(\mathbf{w}_k, \xi_k) - \nabla F(\mathbf{w}_k)\|^2] \leq \frac{\sigma_g^2}{B}. \quad (28)$$

Convergence of Global CSIT AirFL

Let $1 - \frac{L^2\mu^2}{2}\tau(\tau-1) - L\mu\tau \geq 0$. Under Assumptions 1 and 2, the convergence rate is bounded as [7]

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} [\|\nabla F(\mathbf{w}_G^t)\|^2] \leq \frac{2(F(\mathbf{w}_G^0) - F(\mathbf{w}^*))}{\mu\tau T} + \frac{L^2\mu^2(\tau-1)}{2} \frac{\sigma_g^2}{B} + \frac{L}{\mu\tau T} \sum_{t=0}^{T-1} \mathcal{I}_t, \quad (29)$$

where

$$\mathcal{I}_t = \mu^2 \frac{\sigma_g^2}{B} \frac{\tau}{|\mathcal{S}_t|} + \text{MSE}_t. \quad (30)$$

The bound in (29) shows that the average squared gradient along the trajectory decomposes into three contributions. The first term, $\frac{2(F(\mathbf{w}_G^0) - F(\mathbf{w}^*))}{\mu\tau T}$, is a *transient* term: it decays as $1/T$ and captures how far the algorithm starts from the optimum. The second term, $\frac{L^2\mu^2(\tau-1)}{2} \frac{\sigma_g^2}{B}$, is an *irreducible noise floor* due to stochastic gradients: larger mini-batches B and smaller step sizes μ reduce this term, while too many local steps τ amplify it. The third term, $\frac{L}{\mu\tau T} \sum_{t=0}^{T-1} \mathcal{I}_t$, captures the *impact of wireless imperfections*: \mathcal{I}_t grows when fewer devices participate (small $|\mathcal{S}_t|$) or when the MSE is large. Thus, good convergence requires (i) enough rounds T , (ii) a carefully chosen learning rate and number of local steps, and (iii) system designs that keep the MSE small while involving as many devices as possible.

Device Selection and Equalization: One approach to optimizing the convergence rate in (29) is to maximize device participation while keeping the MSE below a specified threshold θ , as

$$\max_{\mathcal{S}, \mathbf{b}} |\mathcal{S}|, \quad (31)$$

subject to

$$\max_{k \in \mathcal{S}} \frac{\|\mathbf{b}\|^2}{|\mathbf{b}^H \mathbf{h}_k|^2} \leq \theta.$$

This is a mixed combinatorial and non-convex optimization problem. Two main approaches have been proposed to solve (31).

- **Difference-of-Convex (DC) Programming Approach [5]:** This approach tackles (31) as a sparse and low-rank optimization problem. It employs DC programming to induce sparsity in device selection and enforce a rank-one constraint essential for effective equalization. The main advantage of this

approach lies in its theoretical rigor and guaranteed global convergence. However, it comes with a high computational complexity of order $\mathcal{O}\left(K((M^2 + K)^3 + M^6)\right)$, which can be a limitation in large-scale or time-sensitive deployments.

- *Matching Pursuit (MP) Based Greedy Scheduling Approach [26]*: This approach reformulates (31) as a sparse support selection problem and solves it using a greedy, iterative matching pursuit algorithm inspired by compressive sensing. By iteratively selecting devices that contribute least to the MSE, it achieves a near-optimal solution with significantly lower computational burden. Its complexity scales polynomially with the number of devices and antennas as $\mathcal{O}(K^2 M^2)$, making it efficient and scalable. However, the approach is heuristic in nature and lacks the theoretical optimality guarantees of DC programming.

The above optimization flow of Global CSIT AirFL can be represented by three sequential stages:



Limitations:

- It requires global CSIT at each device, which is difficult to acquire and maintain in practice due to the overhead of obtaining comprehensive channel knowledge.
- It requires fine synchronization, ensuring symbol-level time alignment and tight carrier-frequency/phase synchronization so that the coherent combining implied by (24) and the equalizer \mathbf{b} remains valid.
- It has high computational complexity due to the joint optimization (31).
- Power control may exceed device capabilities, particularly in low-power IoT applications.
- Owing to transmission power constraints, device selection becomes a critical factor, and fairness can be adversely affected. Nevertheless, as this scheme prioritizes maximizing the number of participating devices in (31), it seeks to maintain high learning performance despite these constraints.

V. BLIND OVER-THE-AIR FEDERATED LEARNING

This approach eliminates the need for CSIT by allowing devices to transmit model updates at a constant power level, regardless of their channel conditions, a feature commonly referred to as *blindness*. Notably, this strategy is in complete harmony with traditional wireless system designs, where constant-power, channel-agnostic transmission is widely used. In contrast, CSIT-aware AirFL introduces additional complexity due to the need for accurate CSIT and dynamic power control. Beyond simplifying system requirements, blind transmission offers several practical advantages: it maintains average energy consumption regardless of channel variations, avoids expanding the dynamic range of transmitted signals—thereby simplifying hardware and reducing costs—and eliminates the risk of performance degradation from

channel estimation errors that can lead to unpredictable signal distortions at the receiver. An additional advantage is that all devices contribute to the aggregation process, thereby maximizing the inclusion of locally trained models in the global learning update.

Two main variants exist within this approach: fully blind, which operates without any CSIT, and partial-phase-aware blind, where each device has access to a limited channel phase information only. The fully blind approach, which can rival the performance of CSIT-aware approach, emphasizes post-processing *equalization* as its central design element.

A. Fully Blind AirFL

The server, equipped with a large number of antennas, leverages high dimensional statistical properties to aggregate effectively. This approach was initially proposed in [6] and subsequently examined in greater depth in [27]–[30], providing deeper analyses [30] and several extensions to more practical and complex scenarios, including hardware-impaired transceivers [27], hierarchical multi-cluster architectures [28], and time-varying fading environments [29]. The main process consists of the following key components:

Signal Assumption: No pre-processing is performed at the devices. Assuming a constant transmit power P and only coarse frame-level synchronization among devices (cf. Subsection II-B), the signal received at the m -th antenna of the server can be expressed, based on (2), as

$$\mathbf{y}_m = \sum_{k=1}^K \sqrt{P} h_{k,m} \mathbf{w}_k + \mathbf{z}_m. \quad (32)$$

Here, the channel coefficients are further assumed to be i.i.d. and Gaussian distributed, with

$$\mathbb{E}[h_{k,m}] = 0, \quad \mathbb{E}[|h_{k,m}|^2] = \sigma_h^2, \forall k, m. \quad (33)$$

Equalization and Aggregation: The server applies a special high-dimensional equalization over all antennas as

$$\mathbf{b} = \left[\left(\sum_{k=1}^K h_{k,1} \right)^H, \dots, \left(\sum_{k=1}^K h_{k,M} \right)^H \right], \quad (34)$$

and estimates as

$$\mathbf{w}_G = \frac{\sum_{m=1}^M \left(\sum_{k=1}^K h_{k,m} \right)^H \mathbf{y}_m}{KM\sigma_h^2}. \quad (35)$$

Aggregation Behavior as $M \rightarrow \infty$: We now analyze how the aggregation converges to the true average as the number of antennas grows.

$$\begin{aligned}
 & \lim_{M \rightarrow \infty} \frac{\sum_{m=1}^M \left(\sum_{k=1}^K h_{k,m} \right)^H \left(\sum_{k=1}^K h_{k,m} \mathbf{w}_k + \mathbf{z}_m \right)}{KM\sigma_h^2} \\
 &= \lim_{M \rightarrow \infty} \underbrace{\frac{1}{K} \sum_{k=1}^K \left(\frac{\sum_{m=1}^M |h_{k,m}|^2}{M\sigma_h^2} \right)}_{\text{Desired Aggregation}} \mathbf{w}_k \\
 & \quad + \underbrace{\frac{\sum \sum_{k' \neq k} \sum_{m=1}^M (h_{k,m})^H h_{k',m} \mathbf{w}_{k'}}{KM\sigma_h^2}}_{\text{Error from Interference}} + \underbrace{\frac{\sum_{m=1}^M \left(\sum_{k=1}^K h_{k,m} \right)^H \mathbf{z}_m}{KM\sigma_h^2}}_{\text{Error from AWGN}} \\
 &= \frac{1}{K} \sum_{k=1}^K \mathbf{w}_k.
 \end{aligned} \tag{36}$$

Here, the result relies on the large-number theory as

$$\frac{\sum_{m=1}^M |h_{k,m}|^2}{M} \rightarrow \sigma_h^2, \tag{37}$$

which is a result of the i.i.d. channel assumption. The cross-terms due to interference and noise contributions vanish as $M \rightarrow \infty$.

Number of Antennas for Convergence: A probabilistic lower bound on the number of antennas M required to ensure a bounded estimation error in (35) is established. Specifically, the following theorem shows that, for sufficiently large M , the estimation error remains bounded with high probability, thereby guaranteeing convergence.

Minimum Number of Antennas for Convergence

The absolute and expected estimation error in (35) are bounded as [30]

$$\left\| \mathbf{w}_G - \frac{1}{K} \sum_{k=1}^K \mathbf{w}_k \right\| \leq \frac{\epsilon}{K}, \tag{38}$$

$$\mathbb{E} \left[\left\| \mathbf{w}_G - \frac{1}{K} \sum_{k=1}^K \mathbf{w}_k \right\| \right] \leq \frac{4\gamma_n}{\sqrt{M}c_n} (\sqrt{\pi} + \ln(6K)), \tag{39}$$

where $c_n = 1/\gamma_n + \sigma_h/\sigma_z$ and γ_n is a positive constant. Moreover, to satisfy the above with probability at least $1 - \delta$, the number of antennas M must fulfill

$$M \geq \frac{8\gamma_n^2 K^2}{\epsilon^2 c_n^2} \ln \left(\frac{6K}{\delta} \right). \tag{40}$$

The error $\|\mathbf{w}_G - \frac{1}{K} \sum_k \mathbf{w}_k\|$ scales roughly as $1/\sqrt{M}$: more antennas average out interference and noise more effectively. The parameter ϵ controls how tight the approximation must be (smaller ϵ means

stricter accuracy and thus larger M), while δ specifies the confidence level (smaller δ means the guarantee must hold with higher probability, again requiring more antennas). The dependence $M \propto K^2 \ln(6K/\delta)$ shows that, as the number of devices grows, the receiver needs more antennas to maintain the same aggregation quality. In short, this bound formalizes the *massive-MIMO averaging* effect: with sufficiently many antennas, blind AirFL can approximate exact averaging arbitrarily well, but the required array size grows with both network size and desired reliability.

Limitations:

- Requiring a large number of antennas poses practical challenges, is costly, and may be infeasible in many real-world applications.
- The primary assumption underlying this approach is the presence of i.i.d. channel conditions; however, this assumption does not always hold in practical wireless environments.
- No optimization or adaptation is incorporated in this design, rendering it suboptimal; however, this simplification has been made to preserve the tractability and analytical clarity of the design.
- Finite antenna arrays result in residual errors due to interference and noise, leading to biased aggregation that can substantially degrade learning performance and fairness.

B. Partial-Phase-Aware Blind AirFL

This blind approach is primarily applicable in the single-antenna setting and assumes that each device k has access to a partial estimate of its channel phase $\angle h_k$, denoted by $\angle_p h_k$, with an accuracy within the range $[0, \frac{\pi}{2})$. Using this partial information, each device performs *quadrant phase compensation*, adjusting its transmission phase to ensure that all effective channels are coherently aligned with a positive sign at the server. This technique accommodates a wide inaccuracy range—covering one-quarter of the entire phase space—and prevents channel-induced sign inversions in the aggregated signal. This requires coarse clock and time alignment together with moderate carrier-frequency/phase stability so that the residual phase error remains within the compensated quadrant during one aggregation block (cf. Subsection II-B).

No additional pre-processing is required at the device side, and due to the inherent limitations of the single-antenna scenario, post-processing at the receiver is not applied either. Initially proposed in [31], later enhanced with momentum-based acceleration in [32], and further extended to account for heavy-tailed interference effects in [33], this scheme demonstrates that AirFL can achieve convergence even without sophisticated signal processing. Although performance is degraded, this method serves as a useful benchmark representing the lower-bound performance of AirFL systems under minimal system complexity.

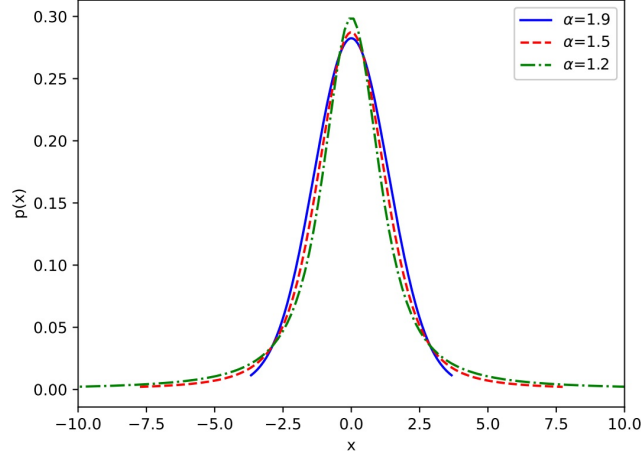


Fig. 2: Probability density function of the α -stable distribution under different values of the tail index [33].

Signal Model: A more general model is considered here, incorporating interference ξ in place of AWGN \mathbf{z} . Accordingly, the real part of the received signal at the single-antenna server is expressed as

$$\mathbf{y}_r = \sum_{k=1}^K \sqrt{P} h_{r,k} \mathbf{w}_k + \xi, \quad (41)$$

where $h_{r,k} = h_k \cos(\angle h_k - \angle_p h_k)$ is the real part of channel h_k after phase compensation.

Two critical assumptions are made:

- All channels $\{h_{r,k}\}$ are i.i.d. random variables with mean ω and variance σ^2 .
- The interference vector ξ has i.i.d. entries following a symmetric α -stable distribution.

It has been well established, both theoretically and empirically, that electromagnetic interference in wireless systems often follows a heavy-tailed distribution [34], [35]. To capture this behavior, the interference can be modeled using the symmetric α -stable distribution, defined below.

Definition: A random variable ξ is said to follow a symmetric α -stable distribution if its characteristic function is given by

$$\mathbb{E} \left[e^{j\omega\xi} \right] = \exp(-\delta^\alpha |\omega|^\alpha), \quad (42)$$

where $\delta > 0$ is the scale parameter and $\alpha \in (0, 2]$ is the tail index. As shown in Fig. 2, smaller values of α correspond to heavier tails. If $\alpha = 1$, the distribution reduces to Cauchy and when $\alpha = 2$, it reduces to Gaussian.

Aggregation: The aggregated signal without any post-processing is estimated as

$$\mathbf{w}_G = \frac{\mathbf{y}}{K}, \quad (43)$$

which is equal to

$$\mathbf{w}_G = \underbrace{\frac{1}{K} \sum_{k=1}^K h_{r,k} \mathbf{w}_k}_{\text{Imposed Aggregation}} + \frac{\boldsymbol{\xi}}{K}. \quad (44)$$

Here, the desired aggregation function is not explicitly defined; instead, the actual outcome, including all its imperfections, is regarded as the imposed aggregation function.

Convergence Analysis: Surprisingly, even without any form of signal processing (except phase compensation), the partial-phase-aware blind scheme can still converge despite the presence of heavy-tailed interference. To support the analysis, the following assumptions are introduced in addition to Assumptions 1 and 2:

Assumption 3 (Strong Convexity): The loss function $F(\mathbf{w})$ is γ -strongly convex, i.e., for any \mathbf{w}_1 and \mathbf{w}_2 , it satisfies

$$F(\mathbf{w}_2) \geq F(\mathbf{w}_1) + \nabla F(\mathbf{w}_1)^\top (\mathbf{w}_2 - \mathbf{w}_1) + \frac{\gamma}{2} \|\mathbf{w}_2 - \mathbf{w}_1\|^2. \quad (45)$$

Assumption 4 (Positive Definiteness of the Hessian): For any given vector \mathbf{w} , the Hessian matrix of $F(\mathbf{w})$, i.e., $\nabla^2 F(\mathbf{w})$, is α -positive definite. A symmetric matrix \mathbf{Q} is said to be α -positive definite if

$$\mathbf{v}^\top \mathbf{Q} \mathbf{v}^{(\alpha-1)} > 0, \quad \forall \mathbf{v} \text{ with } \|\mathbf{v}\|_\alpha^\alpha > 1,$$

where the operator $\mathbf{v}^{(\alpha-1)}$ denotes the *signed-power vector* defined element-wise as

$$\mathbf{v}^{(\alpha-1)} = [\text{sgn}(v_1)|v_1|^{\alpha-1}, \text{sgn}(v_2)|v_2|^{\alpha-1}, \dots, \text{sgn}(v_d)|v_d|^{\alpha-1}]^\top,$$

and the α -norm is defined as

$$\|\mathbf{v}\|_\alpha^\alpha = \left(\sum_{i=1}^d |v_i|^\alpha \right)^{1/\alpha}.$$

Assumption 5 (Bounded Gradient): The gradient at each device k is bounded as

$$\|\nabla F_k(\mathbf{w})\| \leq G, \quad (46)$$

for all k and all model parameters \mathbf{w} .

Assumption 6 (Bounded Interference): The α -moment of the interference is bounded as

$$\mathbb{E} [\|\boldsymbol{\xi}\|_\alpha^\alpha] \leq C, \quad (47)$$

where $C > 0$ is a constant.

Convergence of Partial-Phase-Aware Blind AirFL

If the learning rate in (11) is set to diminish as $\mu_t = \theta/t$ where $\theta > \frac{\alpha-1}{\omega L}$, then the scheme under Assumptions 1–6 converges as [33]

$$\mathbb{E} [\|\mathbf{w}_G^t - \mathbf{w}^*\|_\alpha] \leq \frac{4\theta^\alpha \left(C + \sigma^\alpha G^\alpha s^{1-\frac{1}{\alpha}} / K^{\alpha/2} \right)}{\omega\theta L - \alpha + 1} \cdot \frac{1}{t^{\alpha-1}}. \quad (48)$$

The bound in (48) shows that the distance to the optimum (in the α -norm) decays like $1/t^{\alpha-1}$, where α is the tail index of the interference. When interference is light-tailed (e.g., Gaussian, $\alpha = 2$), the decay behaves like $1/t$, which is relatively fast. As the interference becomes heavier-tailed (smaller α), the convergence slows down; in the extreme, very impulsive interference (small α) causes a much slower decay. The constant in front captures the combined effect of the interference strength C , the channel randomness σ , the bound on local gradients G , the model dimension s , and the number of devices K : more devices and better averaging (larger K) help suppress the randomness, while stronger/heavier-tailed interference inflates the constant. Overall, the quality of the interference environment directly controls how fast the learning error shrinks over time.

Limitations:

- The design is inherently rigid, offering no flexibility for optimization or adaptation.
- Quadrant phase compensation requires intermediate synchronization: coarse clock/time alignment and carrier-frequency/phase stability sufficient to keep the residual phase error within one quadrant.
- It heavily relies on restrictive channel, interference, and learning assumptions.
- Since the aggregation weights are simply the channel coefficients, this leads to biased aggregation and fairness issues, resulting in degraded performance.

VI. WEIGHTED OVER-THE-AIR FEDERATED LEARNING

Weighted AirFL, also known as WAFeL, offers a flexible approach that removes the need for CSIT, power control, and large antenna arrays by introducing and emphasizing a new design dimension: *aggregation weights*. This enables operation even with a single-antenna server. This approach was originally proposed in [7], [36], encompassing both homogeneous and heterogeneous computational settings, and was subsequently extended in [37] to regression problems, with the following key components:

Weighted Aggregation: Other AirFL approaches, like traditional FL, use fixed equal-weight aggregation as in (13). This simple aggregation is grounded on ideal communication conditions, yet remains adopted despite interference, noise, and other wireless impairments. In such cases, AirComp aims to approximate the fixed aggregation function as closely as possible, with the accuracy of this approximation

depending on the effectiveness of the system design. To overcome this, WAFeL proposes a general weighted aggregation as

$$\mathbf{w}_G^t = \sum_{k=1}^K \alpha_k^t \mathbf{w}_k^t, \quad (49)$$

where $\alpha_k^t \geq 0$ is the weight corresponding to device k , and the weight vector $\boldsymbol{\alpha}_t = [\alpha_1^t, \dots, \alpha_K^t]^\top$, such that $\mathbf{1}^\top \boldsymbol{\alpha}_t = 1$. Here, rather than forcing one to approximate the other, both AirComp and the aggregation function are allowed to adapt toward each other, enabling more accurate and robust computation with reduced system requirements.

Quadrant Phase Compensation and Normalization: The WAFeL approach can be applied regardless of the number of antennas at the server. However, in the single-antenna case, it requires quadrant phase compensation. After normalization, as described in (21), the transmitted signal with constant power P is given by

$$\mathbf{x}_k = \sqrt{P} e^{-j\angle_p h_k} \bar{\mathbf{w}}_k, \quad (50)$$

enabling blind transmission.

Signal Model: Arranging the real and imaginary components of the received signal into a vector yields the following real-valued representation at the single-antenna server:

$$\mathbf{Y} = \sqrt{P} \mathbf{H} \bar{\mathbf{W}} + \mathbf{Z}, \quad (51)$$

where $\bar{\mathbf{W}} = [\bar{\mathbf{w}}_1, \dots, \bar{\mathbf{w}}_K]^\top$, and

$$\mathbf{Y} = \begin{bmatrix} \Re\{\mathbf{y}^\top\} \\ \Im\{\mathbf{y}^\top\} \end{bmatrix}, \quad \mathbf{H} = \begin{bmatrix} \Re\{\mathbf{h}^\top\} \\ \Im\{\mathbf{h}^\top\} \end{bmatrix}, \quad \mathbf{Z} = \begin{bmatrix} \Re\{\mathbf{z}^\top\} \\ \Im\{\mathbf{z}^\top\} \end{bmatrix},$$

with $\mathbf{h} = [h_1 e^{j(\angle h_1 - \angle_p h_1)}, \dots, h_K e^{j(\angle h_K - \angle_p h_K)}]^\top$.

Aggregation: The server employs an equalization vector $\mathbf{b} \in \mathbb{R}^{2 \times 1}$ to estimate the weighted aggregation as

$$\mathbf{w}_G^\top = \frac{1}{\sqrt{P}} \mathbf{b}^\top \mathbf{Y} + \sum_{k=1}^K \alpha_k \eta_k \mathbf{1}^\top, \quad (52)$$

which can be written as

$$\mathbf{w}_G^\top = \underbrace{\sum_{k=1}^K \alpha_k \mathbf{w}_k^\top}_{\text{Desired Aggregation}} + \underbrace{\left(\mathbf{b}^\top \mathbf{H} - (\boldsymbol{\alpha} \odot \boldsymbol{\sigma})^\top \right) \bar{\mathbf{W}}}_{\text{Error from Channel Misalignment with Aggregation Weights}} + \underbrace{\frac{1}{\sqrt{P}} \mathbf{b}^\top \mathbf{Z}}_{\text{Error from AWGN}}. \quad (53)$$

where $\boldsymbol{\sigma} = [\sigma_1, \dots, \sigma_K]^\top$ and $(\boldsymbol{\alpha} \odot \boldsymbol{\sigma})^\top = [\alpha_1 \sigma_1, \dots, \alpha_K \sigma_K]$. Due to the normalized transmission, this estimation is unbiased.

Equalization: The optimal equalizer that minimizes MSE arising from the error terms in (53) is obtained as [7]

$$\mathbf{b}_{\text{opt}}^\top = (\boldsymbol{\alpha} \odot \boldsymbol{\sigma})^\top \mathbf{H}^\top \left(\frac{\sigma_z^2}{P} \mathbf{I}_2 + \mathbf{H}\mathbf{H}^\top \right)^{-1}. \quad (54)$$

This equalization leads to

$$\text{MSE}(\boldsymbol{\alpha}) = s \boldsymbol{\alpha}^\top \text{diag}(\boldsymbol{\sigma}) \left(\mathbf{I}_K + \frac{P}{\sigma_z^2} \mathbf{H}^\top \mathbf{H} \right)^{-1} \text{diag}(\boldsymbol{\sigma}) \boldsymbol{\alpha}. \quad (55)$$

The final step involves selecting the aggregation weight vector $\boldsymbol{\alpha}$, which can be optimized to enhance convergence performance.

Convergence Analysis: The convergence rate of WAFeL is given in the next.

Convergence of WAFeL

Let $1 - \frac{L^2 \mu^2}{2} \tau(\tau - 1) - L\mu\tau \geq 0$ and $\boldsymbol{\alpha}_t$ be the weight vector at round t . Then the convergence rate under Assumptions 1 and 2 is bounded as [36]

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} [\|\nabla F(\mathbf{w}_G^t)\|^2] \leq \frac{2(F(\mathbf{w}_G^0) - F(\mathbf{w}^*))}{\mu\tau T} + L^2 \mu^2 \frac{\tau - 1}{2} \frac{\sigma_g^2}{B} + \frac{L}{\mu\tau T} \sum_{t=0}^{T-1} \mathcal{I}_t(\boldsymbol{\alpha}_t), \quad (56)$$

where

$$\mathcal{I}_t(\boldsymbol{\alpha}_t) = \mu^2 \frac{\sigma_g^2}{B} \tau \|\boldsymbol{\alpha}_t\|^2 + \text{MSE}_t(\boldsymbol{\alpha}_t). \quad (57)$$

The convergence bound in (56) has the same qualitative structure as in the global CSIT case, but with an explicit dependence on the aggregation weights $\boldsymbol{\alpha}_t$. The first term, $\frac{2(F(\mathbf{w}_G^0) - F(\mathbf{w}^*))}{\mu\tau T}$, again captures the transient phase and decays as $1/T$. The second term, $\frac{L^2 \mu^2 (\tau - 1)}{2} \frac{\sigma_g^2}{B}$, is the stochastic gradient noise floor. The third term, $\frac{L}{\mu\tau T} \sum_t \mathcal{I}_t(\boldsymbol{\alpha}_t)$, links convergence directly to the choice of weights and the wireless channel. Here, $\mathcal{I}_t(\boldsymbol{\alpha}_t)$ has two parts: $\mu^2 \frac{\sigma_g^2}{B} \tau \|\boldsymbol{\alpha}_t\|^2$ reflects how unbalanced weights amplify gradient noise (large $\|\boldsymbol{\alpha}_t\|$ means a few devices dominate the update), while $\text{MSE}_t(\boldsymbol{\alpha}_t)$ quantifies the aggregation distortion due to AirComp. By optimizing $\boldsymbol{\alpha}_t$ to keep both $\|\boldsymbol{\alpha}_t\|^2$ and $\text{MSE}_t(\boldsymbol{\alpha}_t)$ small, WAFeL jointly balances *learning fairness* and *wireless reliability*, potentially leading to faster and more stable convergence than fixed, non-optimized weights.

Aggregation Weight Selection: The convergence rate in (56) can be optimized by minimizing the norm of the weight vector while ensuring that the MSE remains below a threshold θ , as

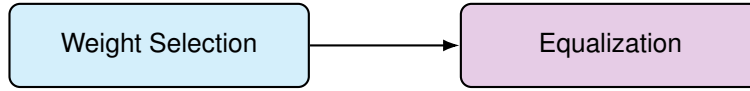
$$\boldsymbol{\alpha}^t = \arg \min_{\boldsymbol{\alpha} \geq 0} \|\boldsymbol{\alpha}\|^2, \quad (58)$$

subject to

$$\begin{aligned} \alpha^\top \text{diag}(\sigma_t) \left(\mathbf{I}_K + \frac{P}{\sigma_z^2} \mathbf{H}_t^\top \mathbf{H}_t \right)^{-1} \text{diag}(\sigma_t) \alpha &\leq \theta, \\ \mathbf{1}^\top \alpha &= 1. \end{aligned}$$

This is a convex optimization problem and can be efficiently solved using standard solvers. A low-complexity iterative algorithm with computational order $\mathcal{O}(K^3)$ for solving (58) is presented in [7]. This approach enables participation from all devices while adapting the aggregation weights to automatically account for both communication and learning aspects—the norm of the weight vector captures the learning contribution, whereas the MSE term reflects the communication quality.

Therefore, the above optimization flow of WAFeL can be represented by two sequential stages:



Soft vs Hard Selection: The continuous flexibility in selecting aggregation weights in (58)—referred to as the *soft* device selection approach—offers lower complexity and greater potential for optimization. In contrast, the *hard* device selection approaches used in local CSIT AirFL (16) and global CSIT AirFL (31) are limited to a small set of discrete choices. In particular, the global CSIT formulation involves a high-complexity integer optimization problem, further increasing the computational burden. Unlike hard selection, which forces a binary decision on whether a device is selected or not, the soft approach allows each device to contribute proportionally to its potential—even if that contribution is small—ensuring that all devices participate in the learning process.

Extension to Computationally Heterogeneous Devices: In heterogeneous environments, each device k possesses a unique hardware capability, typically characterized by its computing speed f_k (e.g., CPU or GPU). To meet a common computation deadline across all devices, the local batch size B_k is adjusted proportionally to $1/f_k$, effectively capturing computational heterogeneity. Under this setting, the heterogeneity-aware FedAvg aggregation is given by

$$\mathbf{w}_G = \frac{1}{B_{\text{tot}}} \sum_{k=1}^K B_k \mathbf{w}_k, \quad (59)$$

where $B_{\text{tot}} = \sum_{k=1}^K B_k$, and the aggregation weights are compactly expressed as the vector $\mathbf{b}_w = \left[\frac{B_1}{B_{\text{tot}}}, \dots, \frac{B_K}{B_{\text{tot}}} \right]^\top$.

The inequality introduced below is fundamental for quantifying computational heterogeneity through the batch sizes B_k , and forms the basis for the following convergence result.

Assumption 7 (Heterogeneity-Aware Gradient Variance Bound): The local stochastic gradient estimate for device k at \mathbf{w}_k , using a mini-batch ξ_k with size B_k , is an unbiased estimate of the ground-truth gradient $\nabla F(\mathbf{w}_k)$ with bounded variance

$$\mathbb{E} [\|\nabla F_k(\mathbf{w}_k, \xi_k) - \nabla F(\mathbf{w}_k)\|^2] \leq \frac{\sigma_g^2}{B_k}. \quad (60)$$

As the batch size increases, the computed gradient becomes more accurate.

Convergence of Heterogeneity-Aware WAFel

Let $1 - \frac{L^2\eta^2}{2}\tau(\tau-1) - L\eta\tau \geq 0$ and α_t as the weight vector for each round t , then the convergence rate under Assumptions 1 and 7 is bounded as [7]

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} [\|\nabla F(\mathbf{w}_G^t)\|^2] \leq \frac{2(F(\mathbf{w}_G^0) - F^*)}{\eta\tau T} + \frac{L}{\eta\tau T} \sum_{t=0}^{T-1} \mathcal{I}_t(\alpha_t), \quad (61)$$

where

$$\mathcal{I}_t(\alpha_t) = L\eta^3 \frac{\tau(\tau-1)}{2} \sigma_g^2 \alpha_t^\top \mathbf{b}_s + \eta^2 \sigma_g^2 \tau \alpha_t^\top \text{diag}\{\mathbf{b}_s\} \alpha_t + \text{MSE}_t(\alpha_t), \quad (62)$$

where the vector $\mathbf{b}_s = [\frac{1}{B_1}, \dots, \frac{1}{B_K}]^\top$ is formed from the inverse of batch sizes of all the devices.

In the heterogeneous setting, the asymmetry shows up explicitly in $\mathcal{I}_t(\alpha_t)$ through the vector \mathbf{b}_s . Devices with smaller B_k produce noisier gradients (larger $1/B_k$), which makes them more dangerous to overweight. The first two terms in (62) describe how the combination of local steps, gradient noise, and the chosen weights interacts with this heterogeneity: large weights on low-batch devices can substantially increase $\mathcal{I}_t(\alpha_t)$ and thus slow convergence. The third term, $\text{MSE}_t(\alpha_t)$, again captures the AirComp distortion. The optimization problem in (63) therefore seeks weights that down-weight unreliable (small-batch) devices just enough to control the noise, while also respecting an MSE constraint that ensures acceptable wireless aggregation quality. Intuitively, the bound shows that WAFel can *re-balance* the influence of fast and slow devices: it lets more reliable, high-batch devices have a stronger say in the global update, without completely ignoring weaker devices, thereby achieving convergence even under strong computational heterogeneity.

Therefore, the weight selection problem to optimize the convergence rate can be characterized as

$$\alpha_t = \arg \min_{\alpha \geq 0} \alpha^\top \text{diag}\{\mathbf{b}_s\} \alpha, \quad (63)$$

subject to

$$\begin{aligned} \boldsymbol{\alpha}^\top \text{diag}(\boldsymbol{\sigma}_t) \left(\mathbf{I}_K + \frac{P}{\sigma_z^2} \mathbf{H}_t^\top \mathbf{H}_t \right)^{-1} \text{diag}(\boldsymbol{\sigma}_t) \boldsymbol{\alpha} &\leq \theta, \\ \mathbf{1}^\top \boldsymbol{\alpha} &= 1, \end{aligned}$$

which is a convex problem. Here, the term $\boldsymbol{\alpha}_t^\top \mathbf{b}_s$ in (62) is neglected in comparison to $\boldsymbol{\alpha}_t^\top \text{diag}\{\mathbf{b}_s\} \boldsymbol{\alpha}_t$ for several reasons. First, this simplification improves tractability. Second, both terms attain their minimum when $\boldsymbol{\alpha} = \mathbf{b}_w$, meaning that minimizing the quadratic form $\boldsymbol{\alpha}^\top \text{diag}\{\mathbf{b}_s\} \boldsymbol{\alpha}$ inherently suppresses the linear term $\boldsymbol{\alpha}^\top \mathbf{b}_s$. Moreover, the quadratic term aligns with the widespread use of ℓ_2 -norm-based objectives in optimization, offering smoother behavior compared to the ℓ_1 -norm.

Limitations:

- In the single-antenna setting, it employs quadrant phase compensation, which requires intermediate synchronization as in the partial-phase-aware blind scheme.
- The potential risk of aggregation bias and fairness degradation increases as the norm of the selected weight vector deviates from that of the equal-weight case.
- When aggregation weights are strictly dictated by the learning process—as in applications with high heterogeneity—such an approach may fail due to the lack of flexibility needed for effective adaptation.

VII. COMPARISON BY EXPERIMENTS

An experimental comparison of different approaches is presented in [7], where the learning task involves classification on the standard MNIST dataset using a convolutional neural network (CNN) as the model. The dataset is distributed non-i.i.d. across devices, with each device holding samples from only two classes and varying quantities of data. Performance is measured in terms of learning accuracy on the test set as a function of the global iteration count t . Each result represents the average of 20 independent realizations, accounting for the randomness introduced by Gaussian wireless channels. The experiments are conducted with $K = 30$, $\tau = 3$, and $\text{SNR} = \frac{P}{\sigma_z^2} = 10$.

The evaluated schemes include WAFEL [7], a weighted aggregation scheme; BAA [4], a local CSIT-aware scheme; and GBMA [31], a partial-phase-aware blind scheme, all under the single-antenna server setting to ensure a fair comparison. Additionally, an ideal benchmark using orthogonal transmission is considered, which assumes error-free communication without interference but requires at least $K = 30$ times more resource blocks.

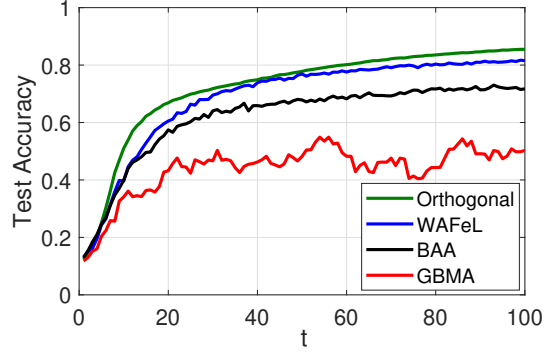


Fig. 3: Performance comparison [7]

While BAA entails higher complexity due to its reliance on perfect CSIT and power control, GBMA operates without any optimization or processing. In contrast, WAFEL incorporates an optimization algorithm with complexity $\mathcal{O}(K^3)$.

Experimental results in Fig. 3 show that WAFEL outperforms both BAA and GBMA. For instance, at $t = 100$, WAFEL achieves approximately 15% and 30% higher accuracy than BAA and GBMA, respectively, and closely approaches the ideal orthogonal baseline. Notably, despite relying only on partial phase compensation, WAFEL outperforms BAA, which has perfect channel gain and phase knowledge. This superiority is attributed to the use of optimized adaptive weights in the aggregation process rather than fixed, unoptimized ones.

A comparative overview is provided in Table I, drawing on the discussions in the preceding sections.

TABLE I: Comparison

Approach	Channel Condition	Fine Synchronization	Channel Estimation	Number of Antennas	Resource	Optimization Complexity Order
Orthogonal	Any	✗	CSIR	≥ 1	K	0
CSIT-Aware	Any	✓	CSIT and CSIR	≥ 1	1	$\mathcal{O}\left(K((M^2 + K)^3 + M^6)\right)$
Blind	i.i.d.	✗	CSIR	∞	1	0
Weighted (WAFEL)	Any	✗	CSIR	≥ 1	1	$\mathcal{O}(K^3)$

VIII. CONCLUSIONS AND FUTURE DIRECTIONS

These lecture notes provide a comprehensive tutorial of Over-the-Air Federated Learning (AirFL), an emerging paradigm that integrates wireless communication, signal processing, and machine learning to enable efficient distributed AI across edge devices. AirFL leverages the waveform superposition property

of wireless signals to perform model aggregation directly over the air, offering significant reductions in communication latency, energy consumption, and bandwidth requirements. We categorized AirFL schemes into three major classes: (i) CSIT-aware approaches; (ii) blind approaches; and (iii) weighted approaches. Each class was analyzed in terms of system complexity, resource requirements, and learning convergence, providing both theoretical foundations and practical implications.

Despite substantial progress, several open problems and promising research directions remain:

- *Local CSIT AirFL with Multiple-Antenna Server*: While local CSIT-based schemes offer a practical middle ground, they are limited to single-antenna servers and heuristic, learning-free designs. Future research may investigate advanced strategies for decentralized yet optimal designs that also incorporate the learning aspect. In particular, even with only local CSIT, multi-antenna servers can be effectively utilized.
- *Global CSIT AirFL with Lower Complexity*: Current global CSIT-based methods incur significant computational and signaling overhead, limiting their scalability. Future work should focus on developing lightweight approximation techniques and efficient device selection algorithms that reduce complexity while preserving convergence guarantees. Incorporating learning-based device selection strategies may further enhance performance by adapting to dynamic network conditions. Moreover, the assumption of full channel knowledge at all devices and centralized coordination is impractical for many applications, highlighting the need to relax these requirements for more feasible and decentralized implementations.
- *Fully Blind AirFL with Optimal Equalization and Limited Number of Antennas*: Blind AirFL schemes rely on asymptotic assumptions, such as large antenna arrays and i.i.d. channel conditions. A key direction for future research is to develop equalization and aggregation techniques that remain effective with a limited number of antennas and can accommodate realistic channel environments.
- *Partial-Phase-Aware Blind AirFL under Generalized Conditions*: Existing analyses of partial-phase-aware blind AirFL rely on restrictive assumptions such as i.i.d. fading and symmetric interference. Future research can relax these assumptions by investigating convergence behavior under correlated, frequency-selective, or time-varying channels, as well as asynchronous or mobility-induced phase variations. Such extensions would provide theoretical guarantees on how AirFL behaves and converges even without explicit processing, offering deeper insight into its potential given its inherently low complexity.
- *WAFEL with Multiple-Antenna Servers*: Extending WAFEL to multi-antenna systems and eliminating its reliance on quadrant phase compensation through a fully blind design open up promising directions

for future research.

- *AirFL with Heterogeneity*: Most AirFL schemes assume homogeneous devices with similar computational capabilities and i.i.d. data. However, real-world FL often involves both computational and statistical heterogeneity, presenting significant challenges. Future research should explore the impact of these heterogeneities and develop heterogeneity-aware AirFL designs that are robust and adaptive to diverse device and data conditions.
- *AirFL for Multi-Task and Personalized Learning*: Conventional AirFL aims to train a single global model, which may not generalize well across diverse devices with distinct data distributions or task objectives. In contrast, multi-task and personalized learning seek to provide each device with a model adapted to its local data or task. To enable personalization in AirFL, aggregation mechanisms must go beyond simple averaging and integrate device-specific model preferences, task identifiers, or metadata. Additionally, meta-learning frameworks and regularization-based personalization techniques can be adapted for AirFL to support model customization at the edge.
- *Secure and Private AirFL*: The open nature of wireless channels in AirFL introduces significant security and privacy risks. Adversaries may eavesdrop on transmissions, manipulate the aggregation process, or infer sensitive information from individual devices, underscoring the need for more robust aggregation functions beyond simple averaging, such as geometric medians. Furthermore, because AirFL operates in the analog domain, conventional digital cryptographic techniques are not directly applicable.
- *Server-less AirFL*: Most AirFL schemes rely on a centralized server to coordinate model aggregation. However, in many emerging applications—such as disaster recovery, rural connectivity, and autonomous multi-agent systems—a central server may be unavailable, unreliable, or undesirable due to latency, cost, or privacy concerns. Future research can explore server-less AirFL, where learning is conducted in a fully decentralized fashion. In such settings, devices would collaboratively train models by communicating directly with their neighbors using local wireless links, leveraging over-the-air computation for peer-to-peer aggregation.
- *AirFL in Other Wireless Scenarios*: Although primarily studied in the basic multiple-access scenario, AirFL can be adapted to a variety of advanced wireless environments. Potential directions include:
 - **AirFL with Mobile Devices**: Addressing challenges related to mobility, handovers, and dynamic topology.
 - **AirFL with RIS, Cell-Free MIMO, or ISAC**: Incorporating intelligent reflecting surfaces (RIS), distributed antenna systems (Cell-Free MIMO), or integrated sensing and communication

(ISAC) to enhance aggregation quality, coverage, and situational awareness in AirFL systems. Overall, AirFL represents a promising direction for realizing wireless distributed intelligence. With continued advances in federated learning and wireless system design, it has the potential to enable reliable, scalable, and resource-efficient AI training across next-generation networks.

REFERENCES

- [1] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Proc. AISTATS*, 2017, pp. 1273–1282.
- [2] B. Nazer and M. Gastpar, "Computation over multiple-access channels," *IEEE Trans. Inf. Theory*, vol. 53, no. 10, pp. 3498–3516, Oct. 2007.
- [3] M. Goldenbaum, H. Boche, and S. Stanczak, "Harnessing interference for analog function computation in wireless sensor networks," *IEEE Trans. Signal Proc.*, vol. 61, no. 20, pp. 4893–4906, 2013.
- [4] 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, Jan. 2020.
- [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, March 2020.
- [6] M. Mohammadi Amiri, T. M. Duman, D. Gunduz, S. R. Kulkarni, and H. V. Poor, "Blind federated edge learning," *IEEE Trans. Wireless Commun.*, vol. 20, no. 8, pp. 5129–5143, Aug. 2021.
- [7] S. M. Azimi-Abarghouyi and L. Tassiulas, "Over-the-air federated learning via weighted aggregation," *IEEE Trans. Wireless Commun.*, vol. 23, no. 12, pp. 18240–18253, Dec. 2024.
- [8] T. Li, A. K. Sahu, A. Talwalkar, V. Smith, "Federated learning: Challenges, methods, and future directions," *IEEE Signal Process. Mag.*, vol. 37, no. 3, pp. 50–60, May 2020.
- [9] T. Gafni, N. Shlezinger, K. Cohen, Y. C. Eldar, and H. V. Poor, "Federated learning: A signal processing perspective," *IEEE Signal Process. Mag.*, vol. 39, no. 3, pp. 14–41, May 2022.
- [10] A. Sahin and R. Yang, "A survey on over-the-air computation," *IEEE Commun. Surv. Tutor.*, vol. 25, no. 3, pp. 1877–1908, 2023.
- [11] A. Perez-Neira, M. Martinez-Gost, A. Sahin, S. Razavikia, C. Fischione, and K. Huang, "Waveforms for computing over the air," *IEEE Signal Process. Mag.*, under review.
- [12] X. Cao, Z. Lyu, G. Zhu, J. Xu, L. Xu, and S. Cui, "An overview on over-the-air federated edge learning," *IEEE Wirel. Commun.*, vol. 31, no. 3, June 2024.
- [13] Z. Chen, H. H. Yang, and T. Q. S. Quek, "Edge intelligence over the air: Two faces of interference in federated learning," *IEEE Commun. Mag.*, vol. 61, no. 12, pp. 62–68, Dec. 2023.
- [14] 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," *IEEE J. Sel. Areas Commun.*, pp. 1–26, 2021.
- [15] H. Hellstrom, J. M. Barros da Silva Jr., M. M. Amiri, M. Chen, V. Fodor, H. V. Poor, and C. Fischione, "Wireless for machine learning: A survey," *Foundations and Trends in Signal Processing*, vol. 15, no. 4, pp. 290–399, 2022.
- [16] X. Cao, G. Zhu, J. Xu, K. Huang, "Optimized power control for over-the-air computation in fading channels," *IEEE Trans. Wireless Commun.*, vol. 19, no. 11, pp. 7498–7513, Nov. 2020.
- [17] X. Cao, G. Zhu, J. Xu, and K. Huang, "Cooperative interference management for over-the-air computation networks," *IEEE Trans. Wireless Commun.*, vol. 20, no. 4, pp. 2634–2651, Apr. 2020.

- [18] S. M. Azimi-Abarghouyi and L. R. Varshney, "Compute-update federated learning: A lattice coding approach," *IEEE Trans. Signal Process.*, vol. 72, pp. 5213-5227, Nov. 2024.
- [19] S. Razavikia, J. M. Barros da Silva Jr., and C. Fischione, "ChannelComp: A general method for computation by communications," *IEEE Trans. Commun.*, vol. 72, no. 2, pp. 692-706, Feb. 2024.
- [20] G. Zhu, Y. Du, D. Gunduz, and K. Huang, "One-bit over-the-air aggregation for communication-efficient federated edge learning: Design and convergence analysis," *IEEE Trans. Wireless Commun.*, vol. 20, no. 3, pp. 2120-2135, Mar. 2021.
- [21] Z. Lin, X. Li, V. K. N. Lau, Y. Gong, and K. Huang, "Deploying federated learning in large-scale cellular networks: Spatial convergence analysis," *IEEE Trans. Wireless Commun.*, vol. 21, no. 3, pp. 1542-1556, Mar. 2022.
- [22] Z. Zhang, G. Zhu, R. Wang, V. K. 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, Oct. 2021.
- [23] S. M. Azimi-Abarghouyi and V. Fodor, "Scalable hierarchical over-the-air federated learning," *IEEE Trans. Wireless Commun.*, vol. 23, no. 8, pp. 8480-8496, Aug. 2024.
- [24] M. Kim, A. L. Swindlehurst, and D. Park, "Beamforming vector design and device selection in over-the-air federated learning," *IEEE Trans. Wireless Commun.*, vol. 22, no. 11, pp. 7464-7477, Nov. 2023.
- [25] Z. Wang, J. Qiu, Y. Zhou, Y. Shi, L. Fu, W. Chen, and K. B. Letaief, "Federated learning via intelligent reflecting surface," *IEEE Trans. Wireless Commun.*, vol. 21, no. 2, pp. 808-822, Feb. 2022.
- [26] A. Bereyhi, A. Vagollari, S. Asaad, R. R. Muller, W. Gerstacker, and H. V. Poor, "Device scheduling in over-the-air federated learning via matching pursuit," *IEEE Trans. Signal Process.*, vol. 71, pp. 2188-2203, June 2023.
- [27] B. Tegin and T. M. Duman, "Blind federated learning at the wireless edge with low-resolution ADC and DAC," *IEEE Trans. Wireless Commun.*, vol. 20, no. 12, pp. 7786-7798, Dec. 2021.
- [28] O. Aygun, M. Kazemi, D. Gunduz, T. M. Duman, "Over-the-air federated edge learning with hierarchical clustering," *IEEE Trans. Wireless Commun.*, vol. 23, no. 12, pp. 17856-17871, Dec. 2024.
- [29] B. Tegin and T. M. Duman, "Federated learning with over-the-air aggregation over time-varying channels," *IEEE Trans. Wireless Commun.*, vol. 22, no. 8, pp. 5671-5684, Aug. 2023.
- [30] S. Razavikia, J. M. Barros da Silva Jr., and C. Fischione, "Blind federated learning via over-the-air q-QAM," *IEEE Trans. Wireless Commun.*, vol. 23, no. 12, pp. 19570-19586, Dec. 2024.
- [31] 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.
- [32] R. Paul, Y. Friedman, and K. Cohen, "Accelerated gradient descent learning over multiple access fading channels," *IEEE J. Sel. Areas Commun.*, vol. 40, no. 2, pp. 532-547, Feb. 2022.
- [33] 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, Apr. 2022.
- [34] L. Clavier, T. Pedersen, I. Rodriguez, M. Lauridsen, and M. Egan, "Experimental evidence for heavy tailed interference in the IoT," *IEEE Commun. Lett.*, vol. 25, no. 3, pp. 692-695, Mar. 2021.
- [35] D. Middleton, "Statistical-physical models of electromagnetic interference," *IEEE Trans. Electromagn. Compat.*, no. 3, pp. 106-127, Aug. 1977.
- [36] S. M. Azimi-Abarghouyi, L. Tassiulas, and C. Fischione, "Weighted over-the-air federated learning," *IEEE ICMLCN*, Barcelona, Spain, May 2025.
- [37] K. Sato and K. Ishibashi, "Adaptively weighted averaging over-the-air computation and its application to distributed Gaussian process regression," *IEEE Trans. Cogn. Commun. Netw.*, early access.