Improving Wireless Federated Learning via Joint Downlink-Uplink Beamforming over Analog Transmission

Chong Zhang, Student Member, IEEE, Min Dong, Fellow, IEEE, Ben Liang, Fellow, IEEE, Ali Afana, Member, IEEE, and Yahia Ahmed

Abstract—Federated learning (FL) over wireless networks using analog transmission can efficiently utilize the communication resource but is susceptible to errors caused by noisy wireless links. In this paper, assuming a multi-antenna base station, we jointly design downlink-uplink beamforming to maximize FL training convergence over time-varying wireless channels. We derive the round-trip model updating equation and use it to analyze the FL training convergence to capture the effects of downlink and uplink beamforming and the local model training on the global model update. Aiming to maximize the FL training convergence rate, we propose a low-complexity joint downlinkuplink beamforming (JDUBF) algorithm, which adopts a greedy approach to decompose the multi-round joint optimization and convert it into per-round online joint optimization problems. The per-round problem is further decomposed into three subproblems over a block coordinate descent framework, where we show that each subproblem can be efficiently solved by projected gradient descent with fast closed-form updates. An efficient initialization method that leads to a closed-form initial point is also proposed to accelerate the convergence of JDUBF. Simulation demonstrates that JDUBF substantially outperforms the conventional separatelink beamforming design.

I. INTRODUCTION

Federated learning (FL) is a widely recognized method to train a machine learning model by multiple collaborating devices using their local training datasets [2]. A parameter server (PS) coordinates the devices for local model updates and aggregates these updates to perform a global model update. In the wireless environment, the PS is usually hosted by a base station (BS), and FL requires frequent exchange of a large amount of model parameters between the BS and many devices over the wireless links, stressing the limited communication resource, such as transmission bandwidth and power [3]. Furthermore, the fluctuation of the wireless links and noisy reception at the receivers introduce distortion, resulting in training errors that degrade the FL performance in both training accuracy and convergence rate. Thus, it is crucial to obtain efficient communication design for wireless FL.

Many existing communication-efficient wireless FL solutions use digital transmission-then-aggregation schemes for uplink acquisition of local parameters from the devices to BS [4]–[8]. The transmission and aggregation are designed separately in these schemes. Conventional digital transmission via orthogonal channels is used, which can consume large bandwidth and incur high latency when the number of devices is large.

To address this issue, analog transmission-and-aggregation schemes have been proposed and analyzed [9]–[13]. In these schemes, devices simultaneously transmit their local models via analog modulation over the shared multiple access channel,

Chong Zhang and Ben Liang are with the Department of ECE, University of Toronto, Canada (e-mail: {chongzhang, liang}@ece.utoronto.ca). Min Dong is with the Department of ECSE, Ontario Tech University, Canada (e-mail: min.dong@ontariotechu.ca). Ali Afana and Yahia Ahmed are with Ericsson Canada, Canada (e-mail: {ali.afana, yahia.ahmed}@ericsson.com). Part of this work was presented in [1].

achieving over-the-air aggregation of local models by superposition. Compared with the digital schemes, such analog schemes can significantly conserve the communication resource and reduce communication latency. The studies in [9]–[11] focus on uplink aggregation while assuming an error-free downlink. However, it is shown that the downlink transmission can be more vulnerable to communication error than the uplink [12]. Noisy downlink transmission for wireless FL has been studied in [13] by assuming an error-free uplink, where it is shown that, since the gradient descent training method in FL is noise resilient, analog transmission can be more efficient than digital transmission even for the downlink.

It is further recognized in [14]–[17] that, especially with analog transmission, downlink and uplink communication for model parameter exchange are coupled during the iterative FL training process. The noise and distortion propagate over the FL communication and computation iterations. This suggests that a joint downlink-uplink design is needed. However, the intertwined process brings significant challenges to tractable analysis and design optimization. The literature on joint downlink-uplink communication design for wireless FL is limited. A recent work has studied the effect of noisy downlink and uplink channels on the convergence of FL with non-independent and identically distributed local datasets [14], where a simple generic signalin-noise receiver model is used to facilitate the analysis. Analog designs have been proposed for noisy downlink and uplink in both single-cell [15], [16] and multi-cell [17] scenarios. However, these schemes only consider single-antenna BSs, and their solutions and convergence analysis cannot be applied to multi-antenna BSs, which are typical in practical wireless systems.

For multi-antenna BSs, beamforming is an essential transmission technique that can be applied to enhance communication quality and reduce noise in wireless FL. Receive beamforming for uplink analog over-the-air aggregation is considered in [18], [19]. Subsequently, various uplink beamforming designs have been proposed for analog schemes to improve the training performance of wireless FL [20]-[25]. These studies consider joint device transmit beamforming and BS receive beamforming. It is shown that carefully designed transmit and receive beamforming schemes can improve uplink over-the-air aggregation. However, the existing literature mostly focus on uplink beamforming design for FL. Our recent work [26] studies downlink beamforming in the context of multi-model FL. As far as we are aware, there is no existing work on joint downlinkuplink beamforming design, which is essential to optimize the overall learning performance.

Besides the analog over-the-air aggregation designs mentioned above, a few recent works have considered digital over-the-air aggregation [27]–[29], where coding schemes with scalar or vector quantization schemes are proposed. These studies focus on how to perform over-the-air computation using digital

modulation and do not consider multi-antenna beamforming. The problems considered there are beyond the scope of this paper.

A. Contributions

We consider wireless FL between a multi-antenna BS and collaborating devices, with noisy analog transmission over both the downlink and the uplink and over-the-air aggregation over the uplink. We jointly design downlink-uplink beamforming to optimize the communication process for FL training over time-varying wireless channels. Such a joint design is challenging as it requires the round-trip iterative model updating structure, and its dependency on the beamforming design can be highly complex. We derive the round-trip global model updating equation, and use it to analyze the training convergence to capture the effects of both communication and the local model training on the global model update. Based on this, we then propose a fast joint downlink-uplink beamforming (JDUBF) algorithm to maximize the FL learning performance. The main contributions are summarized below:

- We formulate the round-trip analog transmission and reception process with downlink and uplink beamforming. We utilize multicast beamforming [30]–[32] for downlink broadcasting of the global model update to devices, which is an efficient beamforming technique to send common data to multiple devices simultaneously. For the uplink, we address uplink beamforming for over-the-air aggregation in two cases based on the availability of channel state information (CSI) at the devices: 1) receive beamforming only, and 2) joint transmit-receive beamforming. Based on this downlink-uplink FL process, we derive the overall global model updating equation over each communication round, capturing the effects of transceiver beamforming and processing over noisy downlink-uplink and the local model training on the global model update.
- Aiming at maximizing the FL training convergence rate, we use the obtained global updating equation to formulate a problem of joint downlink-uplink beamforming and device power optimization, to minimize the expected global training loss after T communication rounds under transmit power constraints at the BS and the devices. For a tractable design, we derive upper bounds on the global training loss through FL training convergence analysis for the two uplink beamforming cases. The bounds show the impact of downlink-uplink beamforming and local device training on the convergence of the global model update, through a weighted sum of the inverse of received signal-to-noise ratios (SNRs) at the BS from all devices.
- Based on the upper bounds, we propose a low-complexity JDUBF algorithm for each of the two uplink beamforming cases: JDUBF-R for receive beamforming only, and JDUBF-TR for joint transmit-receive beamforming. The JDUBF algorithm adopts a greedy approach to decompose the T-round joint beamforming and power allocation optimization problem into separate per-round problems, each then solved in an online optimization manner based on the available global model update at the BS. In particular, we decompose each per-round joint optimization problem into three subproblems via the block coordinate descent (BCD) method [33]. We show that each subproblem can be efficiently solved by the projected gradient descent (PGD) algorithm [34] with fast closed-form gradient updates. To accelerate the

- convergence of the proposed methods, we also propose an efficient initialization method that uses closed-form initial points.
- Our simulation results under typical wireless network settings show that the proposed JDUBF algorithms substantially outperform the conventional separate beamforming design over each link, leading to faster training convergence for a wide range of configurations of the devices and the BS antennas. In particular, JDUBF-TR is shown to nearly attain the learning performance of ideal FL with error-free communication links. It outperforms JDUBF-R at the cost of a higher communication overhead for required information for enabling transmit beamforming among devices.

B. Organization and Notation

The rest of this paper is organized as follows. Section II presents the system model for wireless FL. Section III describes the downlink-uplink analog communication process for FL. In Section IV, we formulate the joint downlink-uplink beamforming optimization problem and develop the upper bounds on the global training loss via FL training convergence analysis. Section V presents the low-complexity JDUBF-R and JDUBF-TR algorithms for the two uplink beamforming cases. Simulation results are shown in Section VI, followed by the conclusion Section VII.

Notation: Hermitian and transpose are denoted as $(\cdot)^H$ and $(\cdot)^T$, respectively. Real and imaginary parts of a complex number are respectively denoted as $\mathfrak{Re}\{\cdot\}$ and $\mathfrak{Im}\{\cdot\}$. The Euclidean norm of a vector is denoted as $\|\cdot\|$, the identity matrix is denoted as \mathbf{I} , and $\mathbf{x} \sim \mathcal{CN}(\mathbf{0}, \mathbf{C})$ denotes a complex Gaussian random vector \mathbf{x} with zero mean and covariance \mathbf{C} .

II. SYSTEM MODEL

A. FL System

We consider a wireless FL system consisting of a server and K devices that collaboratively train a machine learning (ML) model. Let $\mathcal{K} = \{1,\ldots,K\}$ denote the set of devices. Each device $k \in \mathcal{K}$ holds a local training dataset with S_k samples, denoted by $S_k = \{(\mathbf{s}_{k,i}, v_{k,i}) : 1 \leq i \leq S_k\}$, where $\mathbf{s}_{k,i} \in \mathbb{R}^{b \times 1}$ is the i-th data feature vector, and $v_{k,i}$ is the label for this data sample. Let $\theta \in \mathbb{R}^{D \times 1}$ denote the parameter vector of the ML model with D parameters, which predicts the true labels of data feature vectors. Using their respective local training datasets, the devices collaboratively train the global model θ at the server, while keeping their local datasets private. The local training loss function representing the training error at device k is defined as

$$F_k(\boldsymbol{\theta}) = \frac{1}{S_k} \sum_{i=1}^{S_k} L(\boldsymbol{\theta}; \mathbf{s}_{k,i}, v_{k,i})$$

where $L(\cdot)$ is the sample-wise training loss function associated with each data sample. The global training loss function is given by the weighted sum of the local loss functions over all K devices:

$$F(\boldsymbol{\theta}) = \sum_{k=1}^{K} \frac{S_k}{S} F_k(\boldsymbol{\theta}) \tag{1}$$

where $S \triangleq \sum_{k=1}^{K} S_k$ is the total number of training samples from all devices. The learning objective is to find the optimal global model θ^* that minimizes $F(\theta)$.

The devices communicate with the server over the downlink and uplink wireless channels to exchange the model update information iteratively for model training. The training procedure in each communication round $t=0,1,\ldots$ is as follows:

- Downlink broadcast: The server sends the parameter vector of the current global model θ_t to all K devices via the downlink channels.
- Local model update: Each device k uses its dataset S_k to perform local training independently based on the received global model θ_t . Specifically, device k divides S_k into J smaller mini-batches and then performs J iterative local updates using the mini-batches to generate the updated local model $\theta_{k,t}^J$.
- Uplink aggregation: The K devices send their updated local models $\{\theta_{k,t}^J\}_{k\in\mathcal{K}}$ to the server via the uplink channels. The server aggregates $\theta_{k,t}^J$'s to generate an updated global model θ_{t+1} for the next communication round t+1.

B. Wireless Communication Model

Consider a wireless communication system where the server is hosted by a BS equipped with N antennas, and each device has a single antenna. The BS applies downlink transmit beamforming to send the global model update θ_t to the K devices and receive beamforming to the received signals from the K devices to update the global model.

We consider analog communication to transmit the global and local models in the downlink and uplink, respectively. In particular, the BS and the devices send the respective values of θ_t and $\{\theta_{k,t}^J\}_{k\in\mathcal{K}}$ directly under their transmit power budgets. To use the communication bandwidth efficiently for uplink model aggregation, we consider over-the-air computation via analog aggregation over the uplink multiple access channel. Specifically, the devices send their local models $\theta_{k,t}^J$'s to the BS simultaneously over the same frequency band. The BS receives the signal that is the superposition of $\theta_{k,t}^J$'s over the air. We assume that the control and signaling channels of the system still use digital transmissions and are perfect.

Due to the wireless downlink and uplink channels, the received model updates at the devices and the BS are the distorted noisy versions of $\boldsymbol{\theta}_t$ and $\{\boldsymbol{\theta}_{k,t}^J\}_{k\in\mathcal{K}}$, respectively. These errors in the model updates further propagate over the subsequent communication rounds, degrading the learning performance for FL model training. In this paper, we focus on this communication aspect of FL model training. Our goal is to jointly design downlink and uplink beamforming to maximize the learning performance of FL over wireless transmission.

III. DOWNLINK-UPLINK ANALOG TRANSMISSION FOR FL

To study the impact of the non-ideal wireless communication on the learning performance of FL, in this section, we formulate the detailed transmission and reception process with downlink and uplink beamforming in the three stages of a communication round for FL model update as described in Section II-A.

A. Downlink Broadcast of Global Model Update

At the start of communication round t, the BS has the current global model, denoted by $\boldsymbol{\theta}_t = [\theta_{1,t}, \dots, \theta_{D,t}]^\mathsf{T}$. For efficient transmission, we represent $\boldsymbol{\theta}_t$ using a complex signal vector, with real and imaginary parts respectively containing the first and second halves of the elements in $\boldsymbol{\theta}_t$. Specifically, we reexpress $\boldsymbol{\theta}_t = [(\tilde{\boldsymbol{\theta}}_t^{\mathrm{re}})^\mathsf{T}, (\tilde{\boldsymbol{\theta}}_t^{\mathrm{im}})^\mathsf{T}]^\mathsf{T}$, where $\tilde{\boldsymbol{\theta}}_t^{\mathrm{re}} \triangleq [\theta_{1,t}, \dots, \theta_{\frac{D}{2},t}]^\mathsf{T}$,

and $\tilde{\theta}_t^{\text{im}} \triangleq [\theta_{\frac{D}{2}+1,t},\dots,\theta_{D,t}]^\mathsf{T}$. Let $\tilde{\theta}_t$ denote the equivalent complex vector representation of θ_t , given by $\tilde{\theta}_t = \tilde{\theta}_t^{\text{re}} + j\tilde{\theta}_t^{\text{im}} \in \mathbb{C}^{\frac{D}{2} \times 1}$.

Let $\mathbf{h}_{k,t} \in \mathbb{C}^{N \times 1}$ denote the downlink channel vector from the BS to device $k \in \mathcal{K}$ in communication round t. We assume $\{\mathbf{h}_{k,t}\}_{k \in \mathcal{K}}$ remain unchanged during the downlink transmission in round t and are known perfectly at the BS. For sending the common global model to all K devices via beamforming, we consider multicast beamforming [30]–[32], which is an efficient transmission technique to send common signals to multiple devices simultaneously. Specifically, let $\mathbf{w}_t^{\mathrm{dl}} \in \mathbb{C}^{N \times 1}$ be the downlink multicast beamforming vector in communication round t. The BS uses $\mathbf{w}_t^{\mathrm{dl}}$ to send the common $\tilde{\boldsymbol{\theta}}_t$ to all K devices in $\frac{D}{2}$ channel uses. Each device $k \in \mathcal{K}$ receives a complex signal vector, given by

$$\mathbf{u}_{k,t} = (\mathbf{w}_t^{ ext{dl}})^\mathsf{H} \mathbf{h}_{k,t} ilde{ heta}_t + \mathbf{n}_{k,t}^{ ext{dl}}$$

where $\mathbf{n}_{k,t}^{\mathrm{dl}} \in \mathbb{C}^{\frac{D}{2} \times 1}$ is the receiver additive white Gaussian noise (AWGN) vector with i.i.d. elements that are zero mean with variance σ_{d}^2 . The beamforming vector $\mathbf{w}_t^{\mathrm{dl}}$ is subject to the BS transmit power budget. Let DP^{dl} be the total transmit power budget at the BS for sending the global model θ_t in $\frac{D}{2}$ channel uses, where $2P^{\mathrm{dl}}$ represents the average transmit power budget per channel use for sending two elements of θ_t using a complex signal. Then, for transmitting $\tilde{\theta}_t$, $\mathbf{w}_t^{\mathrm{dl}}$ is subject to the transmit power constraint $\|\mathbf{w}_t^{\mathrm{dl}}\|^2 \|\tilde{\theta}_t\|^2 \leq DP^{\mathrm{dl}}$. The BS also sends the scaling factor $\frac{\mathbf{h}_{k,t}\mathbf{w}_t^t}{|(\mathbf{w}_t^{\mathrm{dl}})|^{\mathrm{H}}\mathbf{h}_{k,t}|^2}$ to device k via the downlink signaling channel to facilitate the receiver processing. Device k post-processes the received signal $\mathbf{u}_{k,t}$ using this received scaling factor and obtains the estimate of the complex-valued global model update as

$$\hat{\tilde{\boldsymbol{\theta}}}_{k,t} = \frac{\mathbf{h}_{k,t}^{\mathsf{H}} \mathbf{w}_{t}^{\mathsf{dl}}}{|(\mathbf{w}_{t}^{\mathsf{dl}})^{\mathsf{H}} \mathbf{h}_{k,t}|^{2}} \mathbf{u}_{k,t} = \tilde{\boldsymbol{\theta}}_{t} + \tilde{\mathbf{n}}_{k,t}^{\mathsf{dl}}$$
(2)

where $\tilde{\mathbf{n}}_{k,t}^{\text{dl}} \triangleq \frac{\mathbf{h}_{k,t}^{\text{H}} \mathbf{w}_{t}^{\text{dl}}}{|(\mathbf{w}_{t}^{\text{dl}})^{\text{H}} \mathbf{h}_{k,t}|^{2}} \mathbf{n}_{k,t}^{\text{dl}}$ is the post-processed noise vector at device k. By the equivalence of real and complex signal representations between $\boldsymbol{\theta}_{t}$ and $\tilde{\boldsymbol{\theta}}_{t}$, device k obtains the estimate of the global model $\boldsymbol{\theta}_{t}$, denoted by $\hat{\boldsymbol{\theta}}_{k,t}$, given by

$$\hat{\boldsymbol{\theta}}_{k,t} = \left[\Re \left(\hat{\tilde{\boldsymbol{\theta}}}_{k,t} \right)^{\mathsf{T}}, \Im \left(\hat{\tilde{\boldsymbol{\theta}}}_{k,t} \right)^{\mathsf{T}} \right]^{\mathsf{T}} = \boldsymbol{\theta}_t + \hat{\mathbf{n}}_{k,t}^{\mathsf{dl}}$$
(3)

where $\hat{\mathbf{n}}_{k,t}^{\text{dl}} \triangleq [\mathfrak{Re}\{\tilde{\mathbf{n}}_{k,t}^{\text{dl}}\}^\mathsf{T}, \mathfrak{Im}\{\tilde{\mathbf{n}}_{k,t}^{\text{dl}}\}^\mathsf{T}]^\mathsf{T}$ is the real-valued equivalent post-processed noise vector.

B. Local Model Update

Device k performs local model training based on $\hat{\theta}_{k,t}$ in (3) using its local dataset \mathcal{S}_k . We assume each device adopts the widely used mini-batch stochastic gradient descent (SGD) algorithm to perform the local training for minimizing the local training loss function $F_k(\theta)$ [35]. It uses a subset of the training dataset to compute the gradient update at each iteration and achieves a favorable tradeoff between computational efficiency and convergence rate. Assume J mini-batch SGD iterations are used at each device for its local model update. Let $\theta_{k,t}^{\tau}$ be the local model update by device k at iteration $\tau \in \{0,\ldots,J-1\}$, with $\theta_{k,t}^0 = \hat{\theta}_{k,t}$, and let $\mathcal{B}_{k,t}^{\tau}$ denote the mini-batch, i.e., a subset of \mathcal{S}_k , at iteration τ . Then, the local model update is given by

$$\boldsymbol{\theta}_{k,t}^{\tau+1} = \boldsymbol{\theta}_{k,t}^{\tau} - \eta_t \nabla F_k(\boldsymbol{\theta}_{k,t}^{\tau}; \mathcal{B}_{k,t}^{\tau})$$

$$= \boldsymbol{\theta}_{k,t}^{\tau} - \frac{\eta_t}{|\mathcal{B}_{k,t}^{\tau}|} \sum_{(\mathbf{s},v) \in \mathcal{B}_{k,t}^{\tau}} \nabla L(\boldsymbol{\theta}_{k,t}^{\tau}; \mathbf{s}, v)$$
(4)

where η_t is the learning rate in communication round t, and ∇F_k and ∇L are the gradient functions with respect to (w.r.t.) $\theta_{k,t}^{\tau}$. After J iterations, device k obtains the updated local model $\theta_{k,t}^{J}$.

C. Uplink Aggregation of Local Model Updates

The devices send their local model updates $\{\boldsymbol{\theta}_{k,t}^{J}\}_{k\in\mathcal{K}}$ to the BS using over-the-air aggregation to generate the global model update $\boldsymbol{\theta}_{t+1}$ for the next communication round t+1. For efficient transmission, similar to the downlink, we represent $\boldsymbol{\theta}_{k,t}^{J}$ using a complex vector with its real and imaginary parts containing half of the elements in $\boldsymbol{\theta}_{k,t}^{J}$, respectively. Let $\boldsymbol{\theta}_{k,t}^{J} = [(\tilde{\boldsymbol{\theta}}_{k,t}^{J,\mathrm{re}})^{\mathsf{T}}, (\tilde{\boldsymbol{\theta}}_{k,t}^{J,\mathrm{im}})^{\mathsf{T}}]^{\mathsf{T}}$, where $\tilde{\boldsymbol{\theta}}_{k,t}^{J,\mathrm{re}} \triangleq [\boldsymbol{\theta}_{k1,t}^{J}, \ldots, \boldsymbol{\theta}_{k\frac{D}{2},t}^{J}]^{\mathsf{T}}$ and $\tilde{\boldsymbol{\theta}}_{k,t}^{J,\mathrm{im}} \triangleq [\boldsymbol{\theta}_{k(\frac{D}{2}+1),t}^{J}, \ldots, \boldsymbol{\theta}_{kD,t}^{J}]^{\mathsf{T}}$. The equivalent complex vector representation of $\boldsymbol{\theta}_{k,t}^{J}$ is thus given by $\tilde{\boldsymbol{\theta}}_{k,t}^{J} = \tilde{\boldsymbol{\theta}}_{k,t}^{J,\mathrm{re}} + j\tilde{\boldsymbol{\theta}}_{k,t}^{J,\mathrm{im}} \in \mathbb{C}^{\frac{D}{2} \times 1}$. Device k transmits $\tilde{\boldsymbol{\theta}}_{k,t}^{J}$ to the BS using a total of $\frac{D}{2}$ channel uses.

of $\frac{\tilde{D}}{2}$ channel uses. Let $\mathbf{g}_{k,t} \in \mathbb{C}^{N \times 1}$ denote the uplink channel vector from device $k \in \mathcal{K}$ to the BS in communication round t. We assume $\{\mathbf{g}_{k,t}\}_{k \in \mathcal{K}}$ remain unchanged during the uplink transmission in round t and are known perfectly at the BS. Let $\tilde{\theta}_{kl,t}^J$ be the l-th element in $\tilde{\theta}_{k,t}^J$. The devices send $\tilde{\theta}_{kl,t}^J$'s simultaneously in channel use l. As a result, the received signal vector at the BS for channel use l, denoted by $\mathbf{v}_{l,t}$, is given by

$$\mathbf{v}_{l,t} = \sum_{l=1}^{K} \mathbf{g}_{k,t} a_{k,t} \tilde{\theta}_{kl,t}^{J} + \mathbf{u}_{l,t}^{\mathrm{ul}}$$

where $a_{k,t} \in \mathbb{C}$ is the transmit beamforming weight at device k, and $\mathbf{u}_{l,t}^{\mathrm{ul}} \in \mathbb{C}^{N \times 1}$ is the BS receiver AWGN vector with i.i.d. elements that are zero mean with variance σ_{u}^2 .

The BS applies receive beamforming to the received signal $\mathbf{v}_{l,t}$ over N antennas, for $l=1,\ldots,\frac{D}{2}$. Let $\mathbf{w}_t^{\mathrm{ul}} \in \mathbb{C}^{N\times 1}$ be the unit-norm receive beamforming vector at the BS in communication round t, with $\|\mathbf{w}_t^{\mathrm{ul}}\|^2=1$. Then, the post-processed received signal vector over all $\frac{D}{2}$ channel uses is given by

$$\mathbf{z}_{t} = \left((\mathbf{w}_{t}^{\text{ul}})^{\mathsf{H}} \sum_{k=1}^{K} \mathbf{g}_{k,t} a_{k,t} \right) \tilde{\boldsymbol{\theta}}_{k,t}^{J} + \mathbf{n}_{t}^{\text{ul}} = \sum_{k=1}^{K} \alpha_{k,t}^{\text{ul}} \tilde{\boldsymbol{\theta}}_{k,t}^{J} + \mathbf{n}_{t}^{\text{ul}}$$
(5)

where $\alpha_{k,t}^{\text{ul}} \triangleq (\mathbf{w}_t^{\text{ul}})^{\mathsf{H}} \mathbf{g}_{k,t} a_{k,t}$ represents the effective uplink channel from device k to the BS after applying device transmit and BS receive beamforming, and $\mathbf{n}_t^{\text{ul}} \in \mathbb{C}^{\frac{D}{2} \times 1}$ is the post-processed receiver noise with the l-th element being $(\mathbf{w}_t^{\text{ul}})^{\mathsf{H}} \mathbf{u}_{l,t}^{\text{ul}}$, for $l = 1, \ldots, \frac{D}{2}$. The BS further scales \mathbf{z}_t in (5) to obtain the complex equivalent global model update for the next round t+1:

$$\tilde{\boldsymbol{\theta}}_{t+1} = \frac{\mathbf{z}_t}{\sum_{k=1}^{K} \alpha_{k,t}^{\text{ul}}} = \sum_{k=1}^{K} \rho_{k,t} \tilde{\boldsymbol{\theta}}_{k,t}^{J} + \tilde{\mathbf{n}}_t^{\text{ul}}$$
(6)

where $\rho_{k,t} \triangleq \frac{\alpha_{k,t}^{\text{ul}}}{\sum_{j=1}^{K} \alpha_{j,t}^{\text{ul}}}$ represents the effective weight of device k with $\sum_{k=1}^{K} \rho_{k,t} = 1$, and $\tilde{\mathbf{n}}_t^{\text{ul}} \triangleq \frac{\mathbf{n}_t^{\text{ul}}}{\sum_{k=1}^{K} \alpha_{k,t}^{\text{ul}}}$ is the post-processed noise vector. The effective weight $\rho_{k,t}$ indicates the relative significance of device k's local model update in the global model update. It reflects the overall uplink processing

effect including the device uplink transmission and BS receiver processing.

For transmitting the local model update $\theta_{k,t}^J$ at device k, let $\tilde{\theta}_{k,t}^J$ be its equivalent complex representation and $\Delta \tilde{\theta}_{k,t} \triangleq \tilde{\theta}_{k,t}^J - \tilde{\theta}_{k,t}^0$ the corresponding model change after the local training. Since $\theta_{k,t}^0 = \hat{\theta}_{k,t}$, we have $\tilde{\theta}_{k,t}^0 = \hat{\theta}_{k,t}$. Then, from (2), we express the global model $\tilde{\theta}_{t+1}$ in (6) in terms of $\tilde{\theta}_t$ as

$$\tilde{\boldsymbol{\theta}}_{t+1} = \tilde{\boldsymbol{\theta}}_t + \sum_{k=1}^K \rho_{k,t} \Delta \tilde{\boldsymbol{\theta}}_{k,t} + \sum_{k=1}^K \rho_{k,t} \tilde{\mathbf{n}}_{k,t}^{\text{dl}} + \tilde{\mathbf{n}}_t^{\text{ul}}.$$
 (7)

Finally, the real-valued global model update θ_{t+1} is recovered from $\tilde{\theta}_{t+1}$ as

$$\boldsymbol{\theta}_{t+1} = [\mathfrak{Re}\{\tilde{\boldsymbol{\theta}}_{t+1}\}^\mathsf{T}, \, \mathfrak{Im}\{\tilde{\boldsymbol{\theta}}_{t+1}\}^\mathsf{T}]^\mathsf{T}.$$

Note that the effective uplink channel $\alpha_{k,t}^{\rm ul}$ in (5) depends on the specific transmit/receive beamforming design of the uplink transmission. In the following, we consider two cases: 1) receive beamforming only; 2) joint transmit-receive beamforming.

1) Receive Beamforming Only: We first consider the case where the BS does not send uplink CSI to the devices in order to keep the communication overhead low. In this case, each device has no knowledge of CSI and only applies power scaling to its local model update for transmission, and the BS applies receive beamforming to post-process the received model parameters. Specifically, device k sets $a_{k,t} = \sqrt{p_{k,t}}$, where $p_{k,t}$ denotes the transmit power scaling factor for device k in round t. After applying receive beamforming vector $\mathbf{w}_t^{\rm ul}$ at the BS, the resulting effective uplink channel $\alpha_{k,t}^{\rm ul}$ in (5) is given by

$$\alpha_{k,t}^{\text{ul}} = \sqrt{p_{k,t}} (\mathbf{w}_t^{\text{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}, \ k \in \mathcal{K}.$$
 (8)

Note that $\alpha_{k,t}^{\rm ul}$ is complex-valued, as there is no transmit beamforming applied at the devices, and thus, their signals received at the BS cannot be perfectly phase aligned. This may significantly limit the quality of the aggregated local model updates in (5).

2) Joint Transmit-Receive Beamforming: We also consider the case when CSI is available at the devices. In this case, we can jointly design the device transmit beamforming weights $\{a_{k,t}\}_{k\in\mathcal{K}}$ and the BS receive beamforming vector \mathbf{w}_t^{ul} to phasealign the effective uplink channels. Specifically, the transmit beamforming weight $a_{k,t}$ at device k is given by k

$$a_{k,t} = \sqrt{p_{k,t}} \frac{\mathbf{g}_{k,t}^{\mathsf{H}} \mathbf{w}_{t}^{\mathsf{ul}}}{|(\mathbf{w}_{t}^{\mathsf{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}|}.$$
 (9)

The effective uplink channel between device k and the BS in this case is given by

$$\alpha_{k,t}^{\text{ul}} = (\mathbf{w}_t^{\text{ul}})^{\mathsf{H}} \mathbf{g}_{k,t} a_{k,t} = \sqrt{p_{k,t}} |(\mathbf{w}_t^{\text{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}|, \ k \in \mathcal{K}.$$
 (10)

All effective uplink channels are phase-aligned as the result of joint transmit-receive beamforming.

Finally, we note that each device is subject to transmit power budget. Similar to the downlink, let DP_k^{ul} be the total transmit power budget at device k for sending each local model update $\theta_{k,t}^J$ in $\frac{D}{2}$ channel uses. Then, for transmitting $\tilde{\theta}_{k,t}^J$, the transmit power constraint is $p_{k,t} \|\tilde{\theta}_{k,t}^J\|^2 \leq DP_k^{\mathrm{ul}}$.

 $^1 \text{The value of } a_{k,t}$ in (9) can be obtained at device k by letting the BS send the scalar $\frac{\mathbf{g}_{k,t}^{\mathsf{H}} \mathbf{w}_t^{\mathsf{ul}}}{|(\mathbf{w}_t^{\mathsf{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}|}$ to the device via the downlink signaling channel. This requires additional communication overhead in each communication round.

Remark 1. The global model updating equation in (7) is derived from the entire round-trip FL procedure, including downlink transmission via multicast beamforming, local model updates at devices via mini-batch SGD, and uplink over-the-air aggregation with a given transmit/receive beamforming scheme. The second term in (7) is a weighted average of local model changes at all devices. It represents the aggregated local model change from K devices as the result of uplink transmission. The third term in (7) is the downlink receiver noise from all K devices aggregated at the BS via uplink transmission. The fourth term in (7) is the BS receiver noise from uplink transmission. Overall, the global model updating equation (7) is a noisy version of the aggregated local model updates. It shows how local model updates contribute to the global model update over the noisy communication channel and transmitter and receiver multi-antenna processing.

IV. JOINT DOWNLINK-UPLINK BEAMFORMING FOR WIRELESS FL

A. Joint Downlink-Uplink Formulation

Our objective in this paper is to design the communication aspect of FL, to maximize the model training convergence rate. In particular, we consider the expected global loss function after T communication rounds to measure for training convergence rate. Let $\mathbf{p}_t \triangleq [p_{1,t},\ldots,p_{K,t}]^\mathsf{T}$ contain the uplink transmit power scaling factors of all K devices in round t. Let $\mathbf{w}^{\mathrm{dl}} \triangleq [(\mathbf{w}_0^{\mathrm{dl}})^{\mathrm{H}},\ldots,(\mathbf{w}_{T-1}^{\mathrm{dl}})^{\mathrm{H}}]^{\mathrm{H}} \in \mathbb{C}^{TN \times 1}$, $\mathbf{w}^{\mathrm{ul}} \triangleq [(\mathbf{w}_0^{\mathrm{ul}})^{\mathrm{H}},\ldots,(\mathbf{w}_{T-1}^{\mathrm{ul}})^{\mathrm{H}}]^{\mathrm{H}} \in \mathbb{C}^{TN \times 1}$, and $\mathbf{p} \triangleq [\mathbf{p}_0^\mathsf{T},\ldots,\mathbf{p}_{T-1}^\mathsf{T}]^\mathsf{T} \in \mathbb{R}^{TK \times 1}$ respectively denote the stacked downlink multicast beamforming vectors, the uplink receive beamforming vectors, and the device power scaling vectors over the entire T communication rounds. We aim to find the joint beamforming and power control solution $(\mathbf{w}^{\mathrm{dl}},\mathbf{w}^{\mathrm{ul}},\mathbf{p})$ to maximize the expected global loss after T communication rounds, which is formulated as

$$\mathcal{P}_{o}: \min_{\mathbf{w}^{\mathrm{dl}}, \mathbf{w}^{\mathrm{ul}}, \mathbf{p}} \mathbb{E}[F(\boldsymbol{\theta}_{T})]$$
s.t. $\|\mathbf{w}_{t}^{\mathrm{dl}}\|^{2} \|\boldsymbol{\theta}_{t}\|^{2} \leq DP^{\mathrm{dl}}, \quad t \in \mathcal{T},$ (11)
$$p_{k,t} \|\boldsymbol{\theta}_{k,t}^{J}\|^{2} \leq DP_{k}^{\mathrm{ul}}, \quad k \in \mathcal{K}, t \in \mathcal{T},$$
 (12)

$$\|\mathbf{y}_{k,t}^{\mathrm{ul}}\|_{k,t}^{2} \leq DI_{k}, \quad k \in \mathbb{N}, t \in \mathbb{N},$$

$$\|\mathbf{y}_{k}^{\mathrm{ul}}\|_{k,t}^{2} = 1, \quad t \in \mathcal{T}$$

$$(12)$$

where $\mathcal{T} = \{0, \dots, T-1\}$, and $\mathbb{E}[\cdot]$ is the expectation taken w.r.t. the receiver noise at the devices and the BS and the minibatch sampling for local training at each device. The constraints in (11) and (12) are the transmit power constraints at the BS and each device k, respectively.

Problem \mathcal{P}_o is a finite-horizon stochastic optimization problem. The expected global loss in the objective function is difficult to evaluate. Furthermore, the beamforming in both downlink and uplink is a one-to-many beamforming design, which is equivalent to a multicast beamforming problem that is nonconvex and NP-hard. To tackle this challenging problem, we first analyze the convergence rate of the global loss function and derive a more tractable upper bound on $\mathbb{E}[F(\theta_T)]$ as a function of the downlink and uplink beamforming design parameters $(\mathbf{w}^{\mathrm{dl}}, \mathbf{w}^{\mathrm{ul}}, \mathbf{p})$. We then develop our fast joint downlink-uplink beamforming algorithm to minimize this upper bound.

B. Convergence Analysis on Global Training Loss

The FL learning objective is to find the optimal global model θ^* that minimizes the global training loss function $F(\theta)$.

Let F^{\star} denote the minimum global loss obtained by θ^{\star} . The expected optimality gap between the global loss under θ_T and the minimum global loss after T communication rounds is $\mathbb{E}[F(\theta_T)] - F^{\star}$. To examine $\mathbb{E}[F(\theta_T)]$, we can equivalently analyze $\mathbb{E}[F(\theta_T)] - F^{\star}$ based on the global model updating equation in (7).

We first make the following three assumptions on the local loss functions, the SGD, and the difference between the global and weighted average of the local loss gradients. These assumptions are commonly adopted for the convergence analysis of the FL model training [13], [15], [17].

Assumption 1. The local loss functions $F_k(\cdot)$'s are differentiable and are L-smooth: $F_k(\mathbf{y}) \leq F_k(\mathbf{x}) + (\mathbf{y} - \mathbf{x})^\mathsf{T} \nabla F_k(\mathbf{x}) + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2$, $\forall k \in \mathcal{K}$, $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^{D \times 1}$. Also, $F_k(\cdot)$'s are λ -strongly convex: $F_k(\mathbf{y}) \geq F_k(\mathbf{x}) + (\mathbf{y} - \mathbf{x})^\mathsf{T} \nabla F_k(\mathbf{x}) + \frac{\lambda}{2} \|\mathbf{y} - \mathbf{x}\|^2$, $\forall k \in \mathcal{K}$, $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^{D \times 1}$.

Assumption 2. The mini-batch SGD is unbiased: $\mathbb{E}_{\mathcal{B}}[\nabla F_k(\boldsymbol{\theta}_{k,t}^{\tau}; \mathcal{B}_{k,t}^{\tau})] = \nabla F_k(\boldsymbol{\theta}_{k,t}^{\tau}), \ \forall \ k \in \mathcal{K}, \ \forall \ t \in \mathcal{T}, \\ \forall \ \tau. \ \text{The variance of the mini-batch stochastic local loss gradient is bounded by } \mu > 0 \text{: For } \forall \ k \in \mathcal{K}, \ \forall \ t \in \mathcal{T}, \ \forall \ \tau, \\ \mathbb{E}[\|\nabla F_k(\boldsymbol{\theta}_{k,t}^{\tau}; \mathcal{B}_{k,t}^{\tau}) - \nabla F_k(\boldsymbol{\theta}_{k,t}^{\tau})\|^2] \leq \mu.$

Assumption 3. The gradient divergence is bounded by $\delta > 0$: For $\forall k \in \mathcal{K}, \forall t \in \mathcal{T}, \mathbb{E}[\|\nabla F(\boldsymbol{\theta}_t) - \sum_{k=1}^K \phi_k \nabla F_k(\boldsymbol{\theta}_t)\|^2] \leq \delta$, where $\phi_k \in \mathbb{R}$, and $\sum_{k=1}^K \phi_k = 1$.

In addition, we have the following assumption on the gradient of the global loss and the change in the local model update.

Assumption 4. The gradient of the global loss and the change in the local model update are finite: For some $\zeta > 0$, $\nu > 0$, $\forall \ k \in \mathcal{K}, \ \forall \ t \in \mathcal{T}, \ \|\nabla F(\boldsymbol{\theta}_t)\| \leq \zeta, \ \|\Delta \boldsymbol{\theta}_{k,t}\| \leq \nu.$

We point out that $\|\nabla F(\theta_t)\|$ and $\|\Delta \theta_{k,t}\|$ being finite is typical for FL systems, and thus, Assumption 4 generally holds.

From the global model updating equation in (7), we note that weight $\rho_{k,t}$ may be complex- or real-valued. It depends on the uplink beamforming case considered, as shown in (8) and (10) of Sections III-C1 and III-C2, respectively. In the following convergence rate analysis, we first assume $\rho_{k,t}$ takes the general form as a complex-valued scalar, and then discuss the special case when $\rho_{k,t}$ is real-valued.

Based on Assumptions 1–4, we now analyze the convergence rate of the expected global loss $\mathbb{E}[F(\theta_t)]$ for each round t and provide an upper bound on $\mathbb{E}[F(\theta_T)] - F^\star$ after T rounds. Using the global model update obtained in (7), we first bound the expected change of the global loss function in two consecutive rounds as

$$\mathbb{E}[F(\boldsymbol{\theta}_{t+1}) - F(\boldsymbol{\theta}_{t})] = \sum_{k=1}^{K} \frac{S_{k}}{S} \mathbb{E}[F_{k}(\boldsymbol{\theta}_{t+1}) - F_{k}(\boldsymbol{\theta}_{t})]$$

$$\stackrel{(a)}{\leq} \mathbb{E}[(\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_{t})^{\mathsf{T}} \nabla F(\boldsymbol{\theta}_{t})] + \frac{L}{2} \mathbb{E}[\|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_{t}\|^{2}] \qquad (14)$$

$$\stackrel{(b)}{=} \mathfrak{Re}\{\mathbb{E}[(\tilde{\boldsymbol{\theta}}_{t+1} - \tilde{\boldsymbol{\theta}}_{t})^{\mathsf{H}} \nabla \tilde{F}(\boldsymbol{\theta}_{t})]\} + \frac{L}{2} \mathbb{E}[\|\tilde{\boldsymbol{\theta}}_{t+1} - \tilde{\boldsymbol{\theta}}_{t}\|^{2}]$$

$$\stackrel{(c)}{=} \mathfrak{Re}\{\mathbb{E}\left[\left(\sum_{k=1}^{K} \rho_{k,t} \Delta \tilde{\boldsymbol{\theta}}_{k,t} + \sum_{k=1}^{K} \rho_{k,t} \tilde{\mathbf{n}}_{k,t}^{\mathsf{dl}} + \tilde{\mathbf{n}}_{t}^{\mathsf{ul}}\right)^{\mathsf{H}} \nabla \tilde{F}(\boldsymbol{\theta}_{t})\right]\}$$

$$+ \frac{L}{2} \mathbb{E}\left[\left\|\sum_{k=1}^{K} \rho_{k,t} \Delta \tilde{\boldsymbol{\theta}}_{k,t} + \sum_{k=1}^{K} \rho_{k,t} \tilde{\mathbf{n}}_{k,t}^{\mathsf{dl}} + \tilde{\mathbf{n}}_{t}^{\mathsf{ul}}\right\|^{2}\right] \qquad (15)$$

where (a) follows the L-smoothness of $F_k(\cdot)$'s in Assumption 1

and the fact that $\nabla F(\theta) = \sum_{k=1}^K \frac{S_k}{S} \nabla F_k(\theta)$ from (1), (b) is based on an equivalent expression of (14) by using the equivalent complex representation $\tilde{\theta}_t$ of θ_t , where $\nabla \tilde{F}(\theta_t)$ denotes the equivalent complex representation of the global loss gradient $\nabla F(\theta_t)$ in round t, and (c) is obtained following the global model update in (7). The upper bound in (15) shows the effect of noisy channels and beamforming processing at both downlink and uplink on the loss function. Let $A_{1,t}$ denote the first term and $B_{1,t}$ denote the expectation part $\mathbb{E}[\cdot]$ in the second term in (15), respectively. They are functions of the aggregated local model change, the joint downlink-uplink transmission processing, and the receiver noise at the BS and devices at round t. Below, we provide upper bounds for $A_{1,t}$ and $B_{1,t}$, respectively.

We first derive the upper bound for $A_{1,t}$. For $\Delta \tilde{\theta}_{k,t}$, let $\Delta \theta_{k,t} \triangleq \theta_{k,t}^J - \theta_{k,t}^0$ denote the corresponding real-valued local model change after the local training at device k in round t, where $\Delta \theta_{k,t} = [\mathfrak{Re}\{\Delta \tilde{\theta}_{k,t}\}^\mathsf{T}, \mathfrak{Im}\{\Delta \tilde{\theta}_{k,t}\}^\mathsf{T}]^\mathsf{T}$. Define $\Delta \bar{\theta}_{k,t} = [-\mathfrak{Im}\{\Delta \tilde{\theta}_{k,t}\}^\mathsf{T}, \mathfrak{Re}\{\Delta \tilde{\theta}_{k,t}\}^\mathsf{T}]^\mathsf{T}$. Then, we have

$$A_{1,t} \stackrel{(a)}{=} \mathfrak{Re} \left\{ \mathbb{E} \left[\left(\sum_{k=1}^{K} \rho_{k,t} \Delta \tilde{\boldsymbol{\theta}}_{k,t} \right)^{\mathsf{H}} \nabla \tilde{F}(\boldsymbol{\theta}_{t}) \right] \right\}$$

$$= \mathbb{E} \left[\left(\sum_{k=1}^{K} \mathfrak{Re} \{ \rho_{k,t} \} \Delta \boldsymbol{\theta}_{k,t} \right)^{\mathsf{T}} \nabla F(\boldsymbol{\theta}_{t}) \right]$$

$$+ \mathbb{E} \left[\left(\sum_{k=1}^{K} \mathfrak{Im} \{ \rho_{k,t} \} \Delta \bar{\boldsymbol{\theta}}_{k,t} \right)^{\mathsf{T}} \nabla F(\boldsymbol{\theta}_{t}) \right]$$

$$(17)$$

where (a) is because the receiver noise vectors $\tilde{\mathbf{n}}_{k,t}^{\mathrm{dl}}$'s at the devices and $\tilde{\mathbf{n}}_t^{\mathrm{ul}}$ at the BS are zero mean and independent of $\nabla \tilde{F}(\boldsymbol{\theta}_t)$, and (17) is an equivalent expression of (16) by using the corresponding real-valued parameters. Based on Assumptions 1–4, we obtain an upper bound on $A_{1,t}$ in Lemma 1.

Lemma 1. Consider the FL system described in Section III. Let $\beta_t^{\mathrm{re}} \triangleq \sum_{k=1}^K |\mathfrak{Re}\{\rho_{k,t}\}|$ and $\beta_t^{\mathrm{im}} \triangleq \sum_{k=1}^K |\mathfrak{Im}\{\rho_{k,t}\}|$. Let $Q_t \triangleq 1 - 4\eta_t^2 J^2 L^2$, and assume $\eta_t J < \frac{1}{2L}$, $\forall \ t \in \mathcal{T}$. Based on Assumptions 1–4, $A_{1,t}$ is upper bounded as

$$A_{1,t} \leq 2\eta_{t}J\left(\frac{1-Q_{t}}{Q_{t}}(\beta_{t}^{\text{re}})^{2} - \frac{1}{4}\right)\mathbb{E}\left[\|\nabla F(\boldsymbol{\theta}_{t})\|^{2}\right] + \frac{D(1-Q_{t})}{4\eta_{t}JQ_{t}}\beta_{t}^{\text{re}}\sigma_{d}^{2}\sum_{k=1}^{K}\frac{|\Re \mathfrak{e}\{\rho_{k,t}\}|}{|(\mathbf{w}_{t}^{\text{dl}})^{\mathsf{H}}\mathbf{h}_{k,t}|^{2}} + \frac{\eta_{t}J(1-Q_{t})}{2Q_{t}}(4\delta + \mu)(\beta_{t}^{\text{re}})^{2} + \nu\zeta\beta_{t}^{\text{im}} + \eta_{t}J\delta.$$
 (18)

Proof: See Appendix A.

Note that by the definition of $\rho_{k,t}$ below (6), we have $\beta_t^{\rm re} \geq |\Re {\rm e}\{\sum_{k=1}^K \rho_{k,t}\}| = 1$, and also $\beta_t^{\rm im} \geq 0$. For the special case of $\rho_{k,t}$ being real-valued (i.e., in (10) under the joint transmit-receive beamforming scheme), we have $\beta_t^{\rm re} = 1$ and $\beta_t^{\rm im} = 0$. Also, for the bound in (18), η_t and J are the parameters set in the SGD for the local model update at each device, L, μ, δ are parameters specified in Assumptions 1–3, and ν, ζ are parameters specified in Assumption 4.

For $B_{1,t}$, since the receiver noise at the BS is zero mean and independent of $\sum_{k=1}^K \rho_{k,t}(\Delta \tilde{\boldsymbol{\theta}}_{k,t} + \tilde{\mathbf{n}}_{k,t}^{\mathrm{dl}})$, we have

$$B_{1,t} = \mathbb{E}\left[\left\|\sum_{k=1}^{K} \rho_{k,t} (\Delta \tilde{\boldsymbol{\theta}}_{k,t} + \tilde{\mathbf{n}}_{k,t}^{\mathrm{dl}})\right\|^{2}\right] + \mathbb{E}[\|\tilde{\mathbf{n}}_{t}^{\mathrm{ul}}\|^{2}]$$

$$= \mathbb{E}\left[\left\|\sum_{k=1}^{K} \rho_{k,t} \left(\Delta \tilde{\boldsymbol{\theta}}_{k,t} + \tilde{\mathbf{n}}_{k,t}^{\text{dl}}\right)\right\|^{2}\right] + \frac{D\sigma_{\mathbf{u}}^{2}}{2\left|\sum_{k=1}^{K} \alpha_{k,t}^{\text{ul}}\right|^{2}}.$$
 (19)

By Assumptions 1–4, we upper bound $B_{1,t}$ in Lemma 2.

Lemma 2. Consider the FL system described in Section III. Assume $\eta_t J < \frac{1}{2L}, \ \forall \ t \in \mathcal{T}$. Based on Assumptions 1–4, $B_{1,t}$ is upper bounded as

$$B_{1,t} \leq \frac{2}{L^{2}} \left(\frac{1 - Q_{t}}{Q_{t}} \right) (\beta_{t}^{\text{re}})^{2} \mathbb{E}[\|\nabla F(\boldsymbol{\theta}_{t})\|^{2}] + \frac{D\sigma_{u}^{2}}{2\left|\sum_{k=1}^{K} \alpha_{k,t}^{\text{ul}}\right|^{2}} + D\sigma_{d}^{2} \left(\frac{1 - Q_{t}}{Q_{t}} \beta_{t}^{\text{re}} \sum_{k=1}^{K} \frac{|\Re \{\rho_{k,t}\}|}{|(\mathbf{w}_{t}^{\text{dl}})^{\mathsf{H}} \mathbf{h}_{k,t}|^{2}} + \sum_{k=1}^{K} \frac{|\rho_{k,t}|^{2}}{|(\mathbf{w}_{t}^{\text{dl}})^{\mathsf{H}} \mathbf{h}_{k,t}|^{2}} \right) + \frac{1 - Q_{t}}{2L^{2}Q_{t}} \left(\left(1 - Q_{t} + \frac{Q_{t}}{J} \right) \mu + 4\delta \right) (\beta_{t}^{\text{re}})^{2} + 2\nu^{2} (\beta_{t}^{\text{im}})^{2}. \quad (20)$$

Proof: See Appendix B.

We now analyze the expected gap $\mathbb{E}[F(\theta_T)] - F^*$ after T communication rounds. From (15), the expected optimality gap at round t+1 is bounded as

$$\mathbb{E}[F(\boldsymbol{\theta}_{t+1})] - F^* \le \mathbb{E}[F(\boldsymbol{\theta}_t)] - F^* + A_{1,t} + \frac{L}{2}B_{1,t} \quad (21)$$

where the RHS of (21) is further upper bounded using (18) and (20). Summing up both sides over $t \in \mathcal{T}$ and rearranging the terms, we obtain the upper bound on $\mathbb{E}[F(\theta_T)] - F^*$, which is stated in Proposition 1 below.

Proposition 1. Consider the FL system described in Section III. Let $V_t \triangleq \frac{1-Q_t+\sqrt{1-Q_t}}{Q_t}$, and assume $\frac{1}{2L(4(\beta_t^{\text{re}})^2+1)} \leq \eta_t J < \frac{1}{2L}$, $\forall \ t \in \mathcal{T}$. Based on Assumptions 1–4, the expected gap $\mathbb{E}[F(\boldsymbol{\theta}_T)]-F^\star$ after T communication rounds is upper bounded by

$$\mathbb{E}[F(\boldsymbol{\theta}_{T})] - F^{*} \leq \Gamma \prod_{t=0}^{T-1} G_{t} + \Lambda + \sum_{t=0}^{T-2} H(\mathbf{w}_{t}^{\text{dl}}, \mathbf{w}_{t}^{\text{ul}}, \mathbf{p}_{t}) \prod_{s=t+1}^{T-1} G_{s} + H(\mathbf{w}_{T-1}^{\text{dl}}, \mathbf{w}_{T-1}^{\text{ul}}, \mathbf{p}_{T-1})$$
(22)

where $\Gamma \triangleq \mathbb{E}[F(\boldsymbol{\theta}_0)] - F^{\star}$, $\Lambda \triangleq \sum_{t=0}^{T-2} C_t \left(\prod_{s=t+1}^{T-1} G_s \right) + C_{T-1}$ with

$$G_{t} \triangleq \frac{1 - Q_{t}}{4\eta_{t}J\lambda} \left(4V_{t}(\beta_{t}^{\text{re}})^{2} - 1\right) + 1, \tag{23}$$

$$C_{t} \triangleq \frac{1 - Q_{t}}{4L} \left(V_{t}(4\delta + \mu) + 4\delta + \frac{\mu}{J}\right) (\beta_{t}^{\text{re}})^{2} + \eta_{t}J\delta$$

$$+ L\nu^{2} \left(\beta_{t}^{\text{im}} + \frac{\zeta}{L\nu}\right) \beta_{t}^{\text{im}},$$

and

$$H(\mathbf{w}_{t}^{\text{dl}}, \mathbf{w}_{t}^{\text{ul}}, \mathbf{p}_{t}) \triangleq \frac{LD}{2} \left(V_{t} \beta_{t}^{\text{re}} \sigma_{d}^{2} \sum_{k=1}^{K} \frac{|\mathfrak{Re}\{\rho_{k,t}\}|}{|(\mathbf{w}_{t}^{\text{dl}})^{\text{H}} \mathbf{h}_{k,t}|^{2}} + \sigma_{d}^{2} \sum_{k=1}^{K} \frac{|\rho_{k,t}|^{2}}{|(\mathbf{w}_{t}^{\text{dl}})^{\text{H}} \mathbf{h}_{k,t}|^{2}} + \frac{\sigma_{u}^{2}}{2|\sum_{k=1}^{K} \alpha_{k}^{\text{ul}}|^{2}} \right). \tag{24}$$

Proof: See Appendix C.

The upper bound on $\mathbb{E}[F(\theta_T)] - F^*$ in (22) shows how the downlink-uplink transmission and the local device training impact the convergence of the global model update. In particular, the first term for Γ shows the effect of the initial starting point θ_0 on the convergence. The second term Λ is a weighted sum of C_t 's over T rounds. Each term for round t in the sum accounts

for the gradient difference of the local loss function from the optimal global loss F^* by using the mini-batch SGD at each device, and the uplink aggregation effect of the local gradient updates via beamforming. The third and fourth terms represent a weighted sum of $H(\mathbf{w}_t^{\text{dl}}, \mathbf{w}_t^{\text{ul}}, \mathbf{p}_t)$ over T rounds. The expression of $H(\mathbf{w}_t^{\text{dl}}, \mathbf{w}_t^{\text{ul}}, \mathbf{p}_t)$ in (24) implicitly depends on \mathbf{w}_t^{ul} and \mathbf{p}_t , since both $\rho_{k,t}$ and $\alpha_{k,t}^{\text{ul}}$ are functions of \mathbf{w}_t^{ul} and \mathbf{p}_t .

In particular, we note that $H(\mathbf{w}_t^{\text{dl}}, \mathbf{w}_t^{\text{ul}}, \mathbf{p}_t)$ in (24) is in the form of a weighted sum of the inverse of SNRs (*i.e.*, the noise-to-signal ratio). Two types of SNRs are captured in $H(\mathbf{w}_t^{\text{dl}}, \mathbf{w}_t^{\text{ul}}, \mathbf{p}_t)$: 1) For the terms with σ_d^2 , the inverse of each term in the summation is related to the effective SNR of the downlink multicast (due to downlink beamforming and receiver processing), evaluated at the BS receiver (after uplink beamforming and receiver processing); 2) For the term with σ_u^2 , its inverse is the effective SNR of the uplink local model aggregation due to the device transmission scheme and receive beamforming and processing.

1) Special case for joint transmit-receive beamforming: Proposition 1 considers general weights $\rho_{k,t}$'s for an arbitrary uplink beamforming case. If we consider the joint transmit-receive beamforming scheme in Section III-C2, where the uplink transmission can be phased aligned for local model aggregation as in (10), $\rho_{k,t}$ becomes a positive real-valued weight. In this case, we have $\beta_t^{\rm re}=1$ and $\beta_t^{\rm im}=0$, and $|\Re \{\rho_{k,t}\}|=\rho_{k,t}$. As a result, the bounds in Lemmas 1 and 2 and Proposition 1 are simplified as follows.

Lemma 3. For $\rho_{k,t}$, $\forall k \in \mathcal{K}$, $\forall t \in \mathcal{T}$, being positive and real-valued, and for $\eta_t J < \frac{1}{2L}$, $\forall t \in \mathcal{T}$, the upper bound on $A_{1,t}$ in (18) is given by

$$A_{1,t} \leq 2\eta_t J \left(\frac{1 - Q_t}{Q_t} - \frac{1}{4} \right) \mathbb{E} \left[\|\nabla F(\boldsymbol{\theta}_t)\|^2 \right]$$

$$+ \frac{D(1 - Q_t)}{4\eta_t J Q_t} \sigma_d^2 \sum_{k=1}^K \frac{\rho_{k,t}}{|(\mathbf{w}_t^{\text{dl}})^{\mathsf{H}} \mathbf{h}_{k,t}|^2}$$

$$+ \frac{\eta_t J (1 - Q_t)}{2Q_t} (4\delta + \mu) + \eta_t J \delta.$$

Lemma 4. For $\rho_{k,t}$, $\forall k \in \mathcal{K}$, $\forall t \in \mathcal{T}$, being positive and real-valued, and for $\eta_t J < \frac{1}{2L}$, $\forall t \in \mathcal{T}$, the upper bound on $B_{1,t}$ in (20) is given by

$$\begin{split} B_{1,t} \leq & \frac{2}{L^2} \bigg(\frac{1 - Q_t}{Q_t} \bigg) \mathbb{E}[\|\nabla F(\boldsymbol{\theta}_t)\|^2] + \frac{D\sigma_{\mathrm{u}}^2}{2 \big| \sum_{k=1}^K \alpha_{k,t}^{\mathrm{ull}} \big|^2} \\ & + D\sigma_{\mathrm{d}}^2 \bigg(\frac{1 - Q_t}{Q_t} \sum_{k=1}^K \frac{\rho_{k,t}}{|(\mathbf{w}_t^{\mathrm{dl}})^{\mathrm{H}} \mathbf{h}_{k,t}|^2} + \sum_{k=1}^K \frac{\rho_{k,t}^2}{|(\mathbf{w}_t^{\mathrm{dl}})^{\mathrm{H}} \mathbf{h}_{k,t}|^2} \bigg) \\ & + \frac{1 - Q_t}{2L^2Q_t} \bigg(\bigg(1 - Q_t + \frac{Q_t}{L} \bigg) \mu + 4\delta \bigg). \end{split}$$

Corollary 1. For $\rho_{k,t}$, $\forall \ k \in \mathcal{K}$, $\forall \ t \in \mathcal{T}$, being positive and real-valued, and for $\frac{1}{10L} \leq \eta_t J < \frac{1}{2L}$, $\forall \ t \in \mathcal{T}$, the upper bound on the expected gap $\mathbb{E}[F(\boldsymbol{\theta}_T)] - F^\star$ in (22) is given by

$$\mathbb{E}[F(\boldsymbol{\theta}_{T})] - F^{*} \leq \Gamma \prod_{t=0}^{T-1} G_{t} + \Lambda + \sum_{t=0}^{T-2} H(\mathbf{w}_{t}^{\text{dl}}, \mathbf{w}_{t}^{\text{ul}}, \mathbf{p}_{t}) \prod_{s=t+1}^{T-1} G_{s} + H(\mathbf{w}_{T-1}^{\text{dl}}, \mathbf{w}_{T-1}^{\text{ul}}, \mathbf{p}_{T-1})$$
(25) where $G_{t} = \frac{1 - Q_{t}}{4n_{t}J\lambda} (4V_{t} - 1) + 1, C_{t} = \frac{1 - Q_{t}}{4L} (V_{t}(4\delta + \mu) + 4\delta + \mu)$

 $\left(\frac{\mu}{I}\right) + \eta_t J \delta$, and

$$H(\mathbf{w}_{t}^{\text{dl}}, \mathbf{w}_{t}^{\text{ul}}, \mathbf{p}_{t}) = \frac{LD}{2} \left(V_{t} \sigma_{d}^{2} \sum_{k=1}^{K} \frac{\rho_{k,t}}{|(\mathbf{w}_{t}^{\text{dl}})^{\text{H}} \mathbf{h}_{k,t}|^{2}} + \sigma_{d}^{2} \sum_{k=1}^{K} \frac{\rho_{k,t}^{2}}{|(\mathbf{w}_{t}^{\text{dl}})^{\text{H}} \mathbf{h}_{k,t}|^{2}} + \frac{\sigma_{u}^{2}}{2(\sum_{k=1}^{K} \alpha_{k,t}^{\text{ul}})^{2}} \right).$$
(26)

Remark 2. There are different upper bounds on the expected gap $\mathbb{E}[F(\theta_T)] - F^*$ that have been obtained in the FL literature [13], [15]. However, they are based on either idealized or simplified communication models without considering multi-antenna processing. Specifically, assuming an error-free uplink, [13] provides an upper bound for noisy downlink transmission using a single-antenna BS. The upper bound in [15] is obtained by considering joint noisy downlink-uplink transmission with a single-antenna BS. In contrast, the upper bounds we obtain for the FL convergence rate in (22) and (25) take into account the transmit and receive beamforming and processing over noisy downlink and uplink, which represent more realistic communication model for the practical multi-antenna systems than those in [13], [15].

Note that the upper bounds on $\mathbb{E}[F(\theta_T)] - F^\star$ given in Proposition 1 and Corollary 1 are both in a more tractable form than $\mathbb{E}[F(\theta_T)]$ in the original optimization problem \mathcal{P}_o , which can be explored for the joint beamforming optimization design. In the next section, we propose the joint downlink-uplink beamforming algorithms to minimize these two upper bounds directly.

V. JOINT DOWNLINK-UPLINK BEAMFORMING DESIGN

Instead of \mathcal{P}_o , we now optimize the joint downlink-uplink beamforming to minimize the expected optimality gap $\mathbb{E}[F(\theta_T)] - F^\star$. To do so, we consider the two uplink beamforming cases discussed in Sections III-C1 and III-C2 for local model aggregation.

For the upper bound in (22), only $H(\mathbf{w}_t^{\text{dl}}, \mathbf{w}_t^{\text{ul}}, \mathbf{p}_t)$ depends on the beamforming design. It is a weighted sum of the inverse of SNRs for downlink multicast and uplink aggregation, as discussed below Proposition 1. In particular,

- For the uplink receive beamforming only case: $H(\mathbf{w}_t^{\text{dl}}, \mathbf{w}_t^{\text{ul}}, \mathbf{p}_t)$ is given in (24). Substituting $\alpha_{k,t}^{\text{ul}} = \sqrt{p_{k,t}} (\mathbf{w}_t^{\text{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}$ from (8) into $H(\mathbf{w}_t^{\text{dl}}, \mathbf{w}_t^{\text{ul}}, \mathbf{p}_t)$, we arrive at the expression in (27), denoted by $H^{\mathsf{R}}(\mathbf{w}_t^{\text{dl}}, \mathbf{w}_t^{\text{ul}}, \mathbf{p}_t)$.
- For the uplink joint transmit-receive beamforming case: $H(\mathbf{w}_t^{\mathrm{dl}}, \mathbf{w}_t^{\mathrm{ul}}, \mathbf{p}_t)$ is given in (26). Substituting $\alpha_{k,t}^{\mathrm{ul}} = \sqrt{p_{k,t}} |(\mathbf{w}_t^{\mathrm{ul}})^{\mathrm{H}} \mathbf{g}_{k,t}|$ from (10) into $H(\mathbf{w}_t^{\mathrm{dl}}, \mathbf{w}_t^{\mathrm{ul}}, \mathbf{p}_t)$, we arrive at the expression in (28), denoted by $H^{\mathrm{TR}}(\mathbf{w}_t^{\mathrm{dl}}, \mathbf{w}_t^{\mathrm{ul}}, \mathbf{p}_t)$.

Define

$$\Psi(\mathbf{w}^{\text{dl}}, \mathbf{w}^{\text{ul}}, \mathbf{p})
\triangleq \sum_{t=0}^{T-2} H(\mathbf{w}_{t}^{\text{dl}}, \mathbf{w}_{t}^{\text{ul}}, \mathbf{p}_{t}) \prod_{s=t+1}^{T-1} G_{s} + H(\mathbf{w}_{T-1}^{\text{dl}}, \mathbf{w}_{T-1}^{\text{ul}}, \mathbf{p}_{T-1}), \quad (29)$$

which contains $H(\mathbf{w}_t^{\text{dl}}, \mathbf{w}_t^{\text{ul}}, \mathbf{p}_t)$'s in T rounds. For both the upper bounds on $\mathbb{E}[F(\theta_T)] - F^*$ in (22) and (25), omitting the first two constant terms, we arrive at the joint optimization problem:

$$\mathcal{P}_1 : \min_{\mathbf{w}^{dl}, \mathbf{w}^{ul}, \mathbf{p}} \quad \Psi(\mathbf{w}^{dl}, \mathbf{w}^{ul}, \mathbf{p})$$
s.t. (11)(12)(13).

$$H^{\mathrm{R}}(\mathbf{w}_{t}^{\mathrm{dl}}, \mathbf{w}_{t}^{\mathrm{ul}}, \mathbf{p}_{t}) = \frac{LD}{2} \left(V_{t} \left[\sum_{k=1}^{K} \left| \mathfrak{Re} \left\{ \frac{\sqrt{p_{k,t}}(\mathbf{w}_{t}^{\mathrm{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}}{\sum_{j=1}^{K} \sqrt{p_{j,t}}(\mathbf{w}_{t}^{\mathrm{ul}})^{\mathsf{H}} \mathbf{g}_{j,t}} \right\} \right| \right] \sum_{k=1}^{K} \frac{\sigma_{d}^{2} \left| \mathfrak{Re} \left\{ \frac{\sqrt{p_{k,t}}(\mathbf{w}_{t}^{\mathrm{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}}{\sum_{j=1}^{K} \sqrt{p_{j,t}}(\mathbf{w}_{t}^{\mathrm{ul}})^{\mathsf{H}} \mathbf{g}_{j,t}} \right\} \right| + \frac{\sigma_{d}^{2} \sum_{k=1}^{K} \frac{p_{k,t} |(\mathbf{w}_{t}^{\mathrm{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}|^{2}}{|(\mathbf{w}_{t}^{\mathrm{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}|^{2}} + \frac{\sigma_{d}^{2}}{2} \left| \sum_{k=1}^{K} \frac{p_{k,t} |(\mathbf{w}_{t}^{\mathrm{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}|^{2}}{|(\mathbf{w}_{t}^{\mathrm{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}|^{2}} + \frac{\sigma_{d}^{2}}{2} \right| \left| \sum_{k=1}^{K} \frac{p_{k,t} |(\mathbf{w}_{t}^{\mathrm{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}|^{2}}{|(\mathbf{w}_{t}^{\mathrm{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}|^{2}} \right| + \frac{\sigma_{d}^{2} \sum_{k=1}^{K} \frac{p_{k,t} |(\mathbf{w}_{t}^{\mathrm{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}|^{2}}{|(\mathbf{w}_{t}^{\mathrm{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}|^{2}}}{\left| \sum_{k=1}^{K} \sqrt{p_{k,t}}(\mathbf{w}_{t}^{\mathrm{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}|^{2}} \right| \left| \sum_{k=1}^{K} \frac{p_{k,t} |(\mathbf{w}_{t}^{\mathrm{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}|^{2}}{|(\mathbf{w}_{t}^{\mathrm{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}|^{2}}} \right| \left| \sum_{k=1}^{K} \frac{p_{k,t} |(\mathbf{w}_{t}^{\mathrm{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}|^{2}}{|(\mathbf{w}_{t}^{\mathrm{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}|^{2}}} \right| \left| \sum_{k=1}^{K} \frac{p_{k,t} |(\mathbf{w}_{t}^{\mathrm{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}|^{2}}{|(\mathbf{w}_{t}^{\mathrm{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}|^{2}}} \right| \left| \sum_{k=1}^{K} \frac{p_{k,t} |(\mathbf{w}_{t}^{\mathrm{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}|^{2}}}{|(\mathbf{w}_{t}^{\mathrm{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}|^{2}}} \right| \left| \sum_{k=1}^{K} \frac{p_{k,t} |(\mathbf{w}_{t}^{\mathrm{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}|^{2}}}{|(\mathbf{w}_{t}^{\mathrm{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}|^{2}}} \right| \left| \sum_{k=1}^{K} \frac{p_{k,t} |(\mathbf{w}_$$

(27)

$$H^{\text{TR}}(\mathbf{w}_{t}^{\text{dl}}, \mathbf{w}_{t}^{\text{ul}}, \mathbf{p}_{t}) = \frac{LD}{2} \left(V_{t} \frac{\sigma_{d}^{2} \sum_{k=1}^{K} \frac{\sqrt{p_{k,t}} |(\mathbf{w}_{t}^{\text{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}|}{|(\mathbf{w}_{t}^{\text{ul}})^{\mathsf{H}} \mathbf{h}_{k,t}|^{2}}}{\sum_{k=1}^{K} \sqrt{p_{k,t}} |(\mathbf{w}_{t}^{\text{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}|} + \frac{\sigma_{d}^{2} \sum_{k=1}^{K} \frac{p_{k,t} |(\mathbf{w}_{t}^{\text{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}|^{2}}{|(\mathbf{w}_{t}^{\text{ul}})^{\mathsf{H}} \mathbf{h}_{k,t}|^{2}} + \frac{\sigma_{u}^{2}}{2}}{\left(\sum_{k=1}^{K} \sqrt{p_{k,t}} |(\mathbf{w}_{t}^{\text{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}|\right)^{2}} \right)$$

$$(28)$$

A. The JDUBF Algorithm

Problem \mathcal{P}_1 is a T-horizon joint optimization problem involving T rounds of the model updates. In particular, $(\mathbf{w}_t^{\rm dl}, \mathbf{w}_t^{\rm ul}, \mathbf{p}_t)$ correlates with $(\mathbf{w}_{t-1}^{\rm dl}, \mathbf{w}_{t-1}^{\rm ul}, \mathbf{p}_{t-1})$ through the model updates $\|\boldsymbol{\theta}_t\|^2$ and $\|\boldsymbol{\theta}_{k,t}^J\|^2$ in constraints (11) and (12), which is challenging to solve. Examining the objective function in (29), we note that for both uplink beamforming cases above, if the condition on $\eta_t J$ in Proposition 1 or Corollary 1 is satisfied, we have $G_t > 0$, $\forall t \in \mathcal{T}$, and thus, $\prod_{s=t+1}^{T-1} G_s > 0$. Following this, we separate $\Psi(\mathbf{w}^{\text{dl}}, \mathbf{w}^{\text{ul}}, \mathbf{p})$ into T per-round objective functions and propose a greedy approach to minimize the perround objective in each communication round t, given by

$$\mathcal{P}_2^t: \min_{\mathbf{w}_t^{\text{ul}}, \mathbf{w}_t^{\text{ul}}, \mathbf{p}_t} H(\mathbf{w}_t^{\text{dl}}, \mathbf{w}_t^{\text{ul}}, \mathbf{p}_t)$$

s.t.
$$\|\mathbf{w}_t^{\text{dl}}\|^2 \|\boldsymbol{\theta}_t\|^2 \le DP^{\text{dl}},$$
 (30)

$$p_{k,t} \|\boldsymbol{\theta}_{k,t}^J\|^2 \le DP_k^{\text{ul}}, \ k \in \mathcal{K}, \tag{31}$$

$$\|\mathbf{w}_t^{\mathrm{ul}}\|^2 = 1 \tag{32}$$

where $H(\mathbf{w}_t^{\text{dl}}, \mathbf{w}_t^{\text{ul}}, \mathbf{p}_t)$ is given by either (27) or (28), depending on the uplink beamforming case considered.

Note that the solution $(\mathbf{w}_t^{\text{dl}}, \mathbf{w}_t^{\text{ul}}, \mathbf{p}_t)$ to problem \mathcal{P}_2^t is computed at the BS at the beginning of communication round t for the global model update in round t. However, constraint (31) contains the local model update $\|\boldsymbol{\theta}_{k,t}^{J}\|^2$, which is unavailable at the beginning of round t. To address this issue, we further propose to solve \mathcal{P}_2^t in an online optimization manner by replacing $\|\boldsymbol{\theta}_{k,t}^J\|^2$ in (31) with $\|\boldsymbol{\theta}_t\|^2$, *i.e.*, the current global model available at the BS at the beginning of round t, arriving at the following joint optimization problem

$$\mathcal{P}_{3}^{t} : \min_{\mathbf{w}_{t}^{\text{dl}}, \mathbf{w}_{t}^{\text{ul}}, \mathbf{p}_{t}} H(\mathbf{w}_{t}^{\text{dl}}, \mathbf{w}_{t}^{\text{ul}}, \mathbf{p}_{t})$$
s.t. (30), (32),
$$p_{k,t} \|\boldsymbol{\theta}_{t}\|^{2} \leq DP_{k}^{\text{ul}}, \ k \in \mathcal{K}. \tag{33}$$

The expression of $H(\mathbf{w}_t^{\text{dl}}, \mathbf{w}_t^{\text{ul}}, \mathbf{p}_t)$ given in either (28) or (27) is a complicated nonconvex function of $(\mathbf{w}_t^{\text{dl}}, \mathbf{w}_t^{\text{ul}}, \mathbf{p}_t)$. Thus, finding the optimal solution to \mathcal{P}_3^t is challenging. To compute a solution to \mathcal{P}_3^t efficiently, we adopt the BCD method [33] to update the downlink beamforming \mathbf{w}_t^{dl} and uplink beamforming \mathbf{w}_t^{ul} and \mathbf{p}_t alternatingly. Let $\mathcal{W}_t^{\text{ul}} \triangleq \{\mathbf{w}_t^{\text{ul}}: \|\mathbf{w}_t^{\text{ul}}\|^2 \|\boldsymbol{\theta}_t\|^2 \leq DP^{\text{ul}}\}$, $\mathcal{W}_t^{\text{ul}} \triangleq \{\mathbf{w}_t^{\text{ul}}: \|\mathbf{w}_t^{\text{ul}}\|^2 = 1\}$, and $\mathcal{Y}_t \triangleq \{\mathbf{p}_t: p_{k,t} \|\boldsymbol{\theta}_t\|^2 \leq DP_k^{\text{ul}}, \ k \in \mathcal{K}\}$. The BCD updating procedure to solve \mathcal{P}_3^t is given as follows: At iteration i,

$$\mathbf{w}_{t}^{\text{dl}(i+1)} = \underset{\mathbf{w}_{t}^{\text{dl}} \in \mathcal{W}_{t}^{\text{dl}}}{\operatorname{arg\,min}} H(\mathbf{w}_{t}^{\text{dl}}, \mathbf{w}_{t}^{\text{ul}(i)}, \mathbf{p}_{t}^{(i)}), \tag{34}$$

$$\mathbf{w}_{t}^{\text{ul}(i+1)} = \underset{\mathbf{w}_{t}^{\text{ul}} \in \mathcal{W}_{t}^{\text{ul}}}{\min} H(\mathbf{w}_{t}^{\text{dl}(i+1)}, \mathbf{w}_{t}^{\text{ul}}, \mathbf{p}_{t}^{(i)}),$$
(35)
$$\mathbf{p}_{t}^{(i+1)} = \underset{\mathbf{p}_{t} \in \mathcal{Y}_{t}}{\min} H(\mathbf{w}_{t}^{\text{dl}(i+1)}, \mathbf{w}_{t}^{\text{ul}(i+1)}, \mathbf{p}_{t}).$$
(36)

$$\mathbf{p}_{t}^{(i+1)} = \underset{\mathbf{p}_{t} \in \mathcal{Y}_{t}}{\operatorname{arg\,min}} H(\mathbf{w}_{t}^{\operatorname{dl}(i+1)}, \mathbf{w}_{t}^{\operatorname{ul}(i+1)}, \mathbf{p}_{t}). \tag{36}$$

Subproblems (34)–(36) are respectively a downlink beamforming problem, an uplink receive beamforming problem, and an uplink transmit power design problem. Since the objective function of each subproblem is nonconvex with a complicated expression, we propose to apply PGD to solve each subproblem. PGD [34] is a first-order iterative algorithm that uses gradient updates to solve a constrained minimization problem given by $\min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$, where \mathcal{X} is the convex feasible set for \mathbf{x} . Its updating procedure at iteration j is given by

$$\mathbf{x}_{j+1} = \Pi_{\mathcal{X}} \left(\mathbf{x}_j - \beta \nabla_{\mathbf{x}} f(\mathbf{x}_j) \right) \tag{37}$$

where $\beta > 0$ is the step size and $\Pi_{\mathcal{X}}(\mathbf{x})$ denotes the projection of x onto \mathcal{X} . PGD is guaranteed to find an approximate stationary point for each subproblem in polynomial time [36].

PGD can be applied to our problem, since the objective functions in (34)–(36) are differentiable and the feasible set of each subproblem is convex. Moreover, we adopt PGD because it is particularly suitable for subproblems (34)–(36), where both the gradients and projections can be obtained in closed-form for fast update. Specifically, the gradient of the objective function in each of subproblems (34)-(36) can be obtained in closed-form based on the expression in (28) or (27). Also, the projection operation for each subproblem is given by

• For subproblem (34):

$$\Pi_{\mathcal{W}_t^{\mathrm{dl}}}(\mathbf{w}_t^{\mathrm{dl}}) = \begin{cases} \mathbf{w}_t^{\mathrm{dl}} & \text{if } \mathbf{w}_t^{\mathrm{dl}} \in \mathcal{W}_t^{\mathrm{dl}}, \\ \sqrt{\frac{DP^{\mathrm{dl}}}{\|\mathbf{w}_t^{\mathrm{dl}}\|^2 \|\boldsymbol{\theta}_t\|^2}} \mathbf{w}_t^{\mathrm{dl}} & \text{otherwise}. \end{cases}$$

- For subproblem (35): $\Pi_{\mathcal{W}_t^{\text{ul}}}(\mathbf{w}_t^{\text{ul}}) = \frac{\mathbf{w}_t^{\text{ul}}}{\|\mathbf{w}_t^{\text{ul}}\|}$
- For subproblem (36):

$$[\Pi_{\mathcal{Y}_t}(\mathbf{p}_t)]_k = \min\left\{p_{k,t}, \frac{DP_k^{\mathrm{ul}}}{\|\boldsymbol{\theta}_t\|^2}\right\}, k \in \mathcal{K}$$

where $[\Pi_{\mathcal{Y}_t}(\mathbf{p}_t)]_k$ denotes the k-th element in $\Pi_{\mathcal{Y}_t}(\mathbf{p}_t)$.

In summary, our proposed JDUBF algorithm for \mathcal{P}_1 uses a greedy approach and solves the online version \mathcal{P}_3^t per-round. It uses BCD [33] to solve subproblems (34)–(36) alternatingly, and each is then solved via PGD [34]. The initial point of PGD for each subproblem is set to be the computed solution from the previous BCD iteration. For example, for solving (34) at BCD iteration i, $\mathbf{w}_t^{\mathrm{dl}(i)}$ is used as the initial point for PGD to compute $\mathbf{w}_t^{\mathrm{dl}(i+1)}$.

We denote JDUBF under the uplink receive beamforming only case as JDUBF-R, and that under the uplink joint transmitreceive beamforming case as JDUBF-TR. After obtaining the solution $(\mathbf{w}_t^{\sf dl}, \mathbf{w}_t^{\sf ul}, \mathbf{p}_t)$ to \mathcal{P}_3^t in round t, JDUBF-R and JDUBF-TR have different post-processing procedures at each device:

- JDUBF-R: The benefit of receiver only beamforming is that the BS does not need to send any additional information to devices, and thus the communication overhead is low. In this case, each device directly sets $a_{k,t}$ to its meet its power budget
- JDUBF-TR: For device k to apply the transmit beamforming weight $a_{k,t}$ in (9), the BS computes $a_{k,t}$ based on $(\mathbf{w}_t^{\text{ul}}, p_{k,t})$ and sends it via the downlink signaling channel to each device k. Device k then recovers its power scaling solution $p_{k,t} =$ $|a_{k,t}|$. However, since $p_{k,t}$ is computed at the BS using the global model update θ_t in (33), applying it to the local model update $\theta_{k,t}^J$ may not satisfy the transmit power constraint (11). Each device k needs to further adjust the power scaling factor to ensure the transmit power constraint (11) is met. Thus, the actual power scaling factor used at device k is determined by $\tilde{p}_{k,t}=\min\big\{p_{k,t},\frac{DP_k^{\text{ul}}}{\|\boldsymbol{\theta}_{k,t}^{J}\|^2}\big\}$. Following this, the transmit beamforming weight is updated as $a_{k,t} \leftarrow \sqrt{\tilde{p}_{k,t}} \frac{a_{k,t}}{|a_{k,t}|}$, which yields $a_{k,t} = \sqrt{\tilde{p}_{k,t}} \frac{\mathbf{g}_{k,t}^{\mathsf{H}} \mathbf{w}_t^{\mathsf{ul}}}{|(\mathbf{w}_t^{\mathsf{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}|}$.

 Computational Complexity: For JDUBF, the main com-

putational complexity lies in computing the gradients of $H(\mathbf{w}_t^{\text{dl}}, \mathbf{w}_t^{\text{ul}}, \mathbf{p}_t)$ in beamforming subproblems (34)(35) in each PGD iteration. In particular, the required matrix-vector computation for each beamforming gradient requires O(NK) flops. From our experiments, for $N=16\sim64$ antennas and $K=10\sim 20$ devices, it typically takes $5\sim 1500$ iterations for PGD to converge. We point out that since the updates are all in closed-form, the computation is fast despite it may take 1500 iterations.

B. The Initialization Method

JDUBF requires an initial point $(\mathbf{w}_t^{\text{dl}(0)}, \mathbf{w}_t^{\text{ul}(0)}, \mathbf{p}_t^{(0)})$ for the BCD iterations. A good initial point is desirable as it can facilitate fast convergence. We propose an efficient initialization method that is based on separate downlink and uplink beamforming solutions.

1) Downlink: The downlink beamforming subproblem (34) is a multicast beamforming problem for sending the global model to all K devices. The part of the objective function $H(\mathbf{w}_t^{\mathrm{dl}}, \mathbf{w}_t^{\mathrm{ul}}, \mathbf{p}_t)$ in (28) or (27) w.r.t. $\mathbf{w}_t^{\mathrm{dl}}$ is a summation term in the form of $\frac{\sigma_t^2}{|(\mathbf{w}_t^{\mathrm{dl}})^{\mathrm{H}}\mathbf{h}_{k,t}|^2}$, *i.e.*, the inverse of received SNR at device k. Heuristically, the summation would be dominated by the term with the minimum SNR. Thus, to minimize the summation, it is effective to maximize the minimum SNR, which is the MMF multicast beamforming problem [30]-[32]. Thus, we propose to use the beamforming solution to the MMF problem given by

$$\max_{\mathbf{w}_{t}^{\text{dl}} \in \mathcal{W}_{t}^{\text{dl}}} \min_{k \in \mathcal{K}} |(\mathbf{w}_{t}^{\text{dl}})^{\mathsf{H}} \mathbf{h}_{k,t}|^{2}.$$
 (38)

Note that the asymptotic multicast beamforming solution to the above problem, as the number of antennas $N \to \infty$, is obtained in closed-form in (49) of [31]. Thus, for fast initialization, we propose to use this asymptotic MMF multicast beamforming solution as the initial $\mathbf{w}_t^{\mathrm{dl}(\hat{0})}$. It is in a form of MMSE beamformer with all parameters obtained in closedform. We refer readers to [31] for the expression detail.

2) Uplink: For uplink subproblems (35)(36), we notice in (28) or (27) that the term in $H(\mathbf{w}_t^{\text{dl}}, \mathbf{w}_t^{\text{ul}}, \mathbf{p}_t)$ related to uplink aggregation is an inverse to received SNR as to uplink aggregation is an inverse $\frac{\sigma_{\mathrm{u}}^2/2}{(\sum_{k=1}^K \sqrt{p_{k,t}}|(\mathbf{w}_t^{\mathrm{ul}})^{\mathrm{H}}\mathbf{g}_{k,t}|)^2}$. We can further lower bound the SNR as $\left(\sum_{k=1}^K \sqrt{p_{k,t}}|(\mathbf{w}_t^{\mathrm{ul}})^{\mathrm{H}}\mathbf{g}_{k,t}|\right)^2 \geq \sum_{k=1}^K p_{k,t}|(\mathbf{w}_t^{\mathrm{ul}})^{\mathrm{H}}\mathbf{g}_{k,t}|^2$. Thus, we set $(\mathbf{w}_t^{\mathrm{ul}}), \mathbf{p}_t^{(0)}$ as the solution of the following uplink received SNR maximization problem:

$$\max_{\mathbf{w}_t^{\text{ul}} \in \mathcal{W}_t^{\text{ul}}, \mathbf{p}_t \in \mathcal{Y}_t} \frac{1}{\sigma_u^2} \sum_{k=1}^K p_{k,t} |(\mathbf{w}_t^{\text{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}|^2.$$
(39)

The power solution to the above problem is the maximum power satisfying constraint (33): $p_{k,t} = \frac{DP_k^{\text{ul}}}{\|\boldsymbol{\theta}_t\|^2}, \forall k$. The optimal solution \mathbf{w}_t^{ul} to problem (39) is the eigenvector corresponding to the largest eigenvalue of $\sum_{k=1}^K p_{k,t} \mathbf{g}_{k,t} \mathbf{g}_{k,t}^{\mathsf{H}}$. From the above proposed method, we obtain the initial point $(\mathbf{w}_t^{\text{ul}(0)}, \mathbf{w}_t^{\text{ul}(0)}, \mathbf{p}_t^{(0)})$ all in closed-form.

C. Separate Downlink and Uplink Beamforming Design

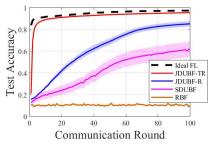
To compare our joint beamforming approach JDUBF, we also consider the conventional method where downlink and uplink transmissions are designed separately for the communication system, which we refer to as SDUBF. This method closely resembles the downlink and uplink beamforming problems we consider in our initialization method above:

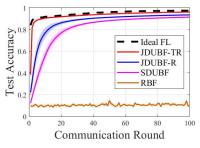
- 1) Downlink: The downlink transmission of the common global model to all K devices is a multicast beamforming problem, which can be formulated as the MMF problem exactly as in (38) considered in the initialization method for our proposed method. Instead of the asymptotic solution used there, we can solve the optimization problem more accurately by the projected subgradient algorithm (PSA) [37].
- 2) Uplink: For the uplink over-the-air aggregation, each device uses transmit beamforming weight $a_{k,t} = \sqrt{p_{k,t}} \frac{\mathbf{g}_{k,t}^{\mathsf{H}} \mathbf{w}_{t}^{\mathsf{ul}}}{|(\mathbf{w}_{t}^{\mathsf{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}|}$ in (9) to phase-align all $(\mathbf{w}_{t}^{\mathsf{ul}})^{\mathsf{H}} \mathbf{g}_{k,t}$'s, where $p_{k,t} = \frac{DP_{k}^{\mathsf{ul}}}{\|\boldsymbol{\theta}_{k,t}^{\mathsf{ul}}\|^{2}}$ is the maximum power scaling factor that the maximum power scaling factor that meets the power budget. Assume all $p_{k,t}$'s are known perfectly at the BS. The uplink receive beamforming problem is to maximize the received SNR of the aggregated signal, which is exactly the problem (39).

VI. SIMULATION RESULTS

A. Simulation Setup

We consider the real-world dataset for image classification via a wireless FL system using typical 5G wireless specifications [38]. We consider system bandwidth 10 MHz and carrier frequency 2 GHz. The maximum transmit power at the BS is 47 dBm, and that at each device is 23 dBm. We assume the devices use 1 MHz bandwidth for their uplink transmission. The channels of all devices are generated i.i.d. as $\mathbf{h}_{k,t} = \sqrt{G_k} \mathbf{\bar{h}}_{k,t}$ with $\mathbf{\bar{h}}_{k,t} \sim \mathcal{CN}(\mathbf{0}, \mathbf{I})$ for downlink, and $\mathbf{g}_{k,t} = \sqrt{G_k} \bar{\mathbf{g}}_{k,t}$ with $\bar{\mathbf{g}}_{k,t} \sim \mathcal{CN}(\mathbf{0},\mathbf{I})$ for uplink, where G_k is the channel variance. We model G_k based on the pathloss model: $G_k[dB] = -161.3 - 35 \log_{10} d_k - \psi_k$, where d_k is the BS-device distance in kilometers, and ψ_k is the shadowing random variable with standard deviation 8 dB. Receiver noise power spectral density is $N_0 = -174$ dBm/Hz, and noise figure is set to $N_F = 8$ dB and 2 dB at the device and BS receivers, respectively.





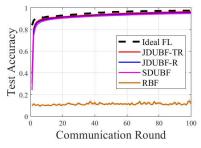


Fig. 1. Test accuracy vs. communication round T with equal-distance devices for (N,K)=(64,10). Left: $G_k=-145$ dB, $\forall k$, $({\rm SNR_{dl}},{\rm SNR_{ul}})=(-2,-10)$ dB. Middle: $G_k=-140$ dB, $\forall k$, $({\rm SNR_{dl}},{\rm SNR_{ul}})=(3,-5)$ dB. Right: $G_k=-133$ dB, $\forall k$, $({\rm SNR_{dl}},{\rm SNR_{ul}})=(10,2)$ dB.

We adopt the MNIST dataset [39] for model training and testing. MNIST consists of 6×10^4 training samples and 1×10^4 test samples from 10 different classes. Each sample is a labeled image of size 28×28 pixels with $\mathbf{s}_{k,i} \in \mathbb{R}^{784 \times 1}$ and $v_{k,i} \in \{0,\ldots,9\}$ indicating the class. We consider training a convolutional neural network with an $8 \times 3 \times 3$ ReLU convolutional layer, a 2×2 max pooling layer, a ReLU fully-connected layer, and a softmax output layer, resulting in $D = 1.361 \times 10^4$ model parameters in total. We use the 10^4 test samples to measure the test accuracy of the global model update θ_t at each round t. The training samples are randomly and evenly distributed over K devices, with the local dataset at device k having $S_k=\frac{6\times 10^4}{K}$ samples. For the local training via the SGD at each device, we set J=30, the mini-batch size $|\mathcal{B}_{k,t}^{\tau}| = \frac{2 \times 10^3}{K}, \forall k, \tau, t$, and the learning rate $\eta_t = 0.1, \ \forall t$. For both JDUBF-TR and JDUBF-R, we set the step size of PGD to $\beta = 0.01$.

B. Performance Comparison

We evaluate the performance of our joint downlink-uplink beamforming design using the proposed JDUBF-TR and JDUBF-R method. We also consider the following three approaches for comparison:

- Ideal FL [2]: FL with error-free downlink and uplink and perfect recovery of model parameters at the BS and devices. Specifically, it uses the global model update in (7), with receiver noise $\tilde{\mathbf{n}}_{k,t}^{\text{dl}} = \tilde{\mathbf{n}}_t^{\text{ul}} = \mathbf{0}$ and receiver post-processing weight $\rho_{k,t} = \frac{1}{K}, \forall k,t$. This benchmark provides the performance upper bound for all schemes.
- SDUBF: The conventional separate downlink and uplink beamforming design based on SNR maximization, which is described in Section V-C. The step size of PSA used in SDUBF is set to 0.01.
- RBF: Randomly generated downlink and uplink beamforming vectors \mathbf{w}_t^{dl} and \mathbf{w}_t^{ul} . The devices use the maximum transmit power and do not perform any transmit beamforming phase alignment. This scheme is used to demonstrate the performance gain of a properly designed beamforming solution.
- 1) Performance with Equal-Distance Devices: We first consider all devices have the same average channel quality and study the test accuracy performance of different methods over communication round T. In particular, all devices are at the same distance d_k from the BS. In this case, all device channels have the same pathloss G_k , and thus, the same downlink/uplink average nominal SNR when the BS/device transmits with full power using a single antenna. The average nominal SNRs at downlink and uplink are denoted by SNR_{dl} and SNR_{ul}, respectively. We use them to indicate the level of channel quality in our experiments. We set (N,K)=(64,10). All results are obtained by averaging over 20 channel realizations.

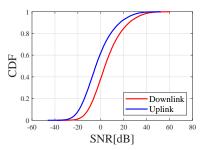
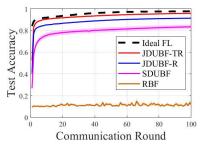
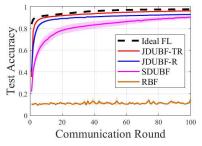


Fig. 2. The CDF of downlink and uplink SNRs with randomly located devices.

Fig. 1-Left shows the test accuracy performance for $(SNR_{dl}, SNR_{ul}) = (-2, -10) dB$, a very low SNR environment on both links. We also plot the 90% confidence interval of each curve as the shadowed area. We see that JDUBF-TR substantially outperforms the other schemes. In particular, it nearly attains the upper bound under the Ideal FL for T > 20and achieves an accuracy above 95% for $T \approx 100$. JDUBF-R is much worse than JDUBF-TR in this case. This is expected since JDUBF-R only uses uplink receive beamforming without transmit beamforming at devices, and this results in much more noisy local model aggregation at a low SNR environment. Nonetheless, it still noticeably outperforms SDUBF, which employs uplink transmit-receive beamforming. This demonstrates the effectiveness of our joint downlink-uplink beamforming design for FL that considers both downlink and uplink communication effect for the global model update. In comparison, SDUBF has a much slower training convergence rate, only reaching 60% test accuracy after T = 100. For RBF, no training convergence is observed. This is because that RBF provides no beamforming gain, leading to highly suboptimal communication performance that affects the learning performance.

Fig. 1-Middle shows the test accuracy for $(SNR_{dl}, SNR_{ul}) =$ (3 dB, -5 dB), a moderate downlink SNR environment and a low uplink SNR environment. We see that as the channel quality improves, the learning performance improves, leading to a higher test accuracy under all beamforming design methods. JDUBF-TR has the fastest convergence rate, nearly attaining the optimal performance after 10 rounds, while JDUBF-R and SDUBF approach the upper bound with slower convergence rates and are worse than JDUBF-TR after 100 rounds. JDUBF-R still converges faster than SDUBF with a slightly better accuracy after 100 rounds. Fig. 1-Right shows the case for $(SNR_{dl}, SNR_{ul}) = (10 \text{ dB}, 2 \text{ dB}), \text{ a high downlink SNR}$ environment and a moderate uplink SNR environment. We see that the performance of all three methods JDUBF-TR, JDUBF-R and SDUBF are close to the Ideal FL, in terms of both convergence rate and the achieved test accuracy. This is expected as when the device channel quality improves, the





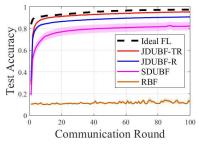


Fig. 3. Test accuracy vs. communication round T with randomly located devices. Left: N = 16, K = 20. Middle: N = 64, K = 20. Right: N = 64, K = 10.

impact of noisy communication on the performance reduces, and the different beamforming design methods can all lead to good learning performance close to the Ideal FL.

2) Performance with Randomly Located Devices: We now consider randomly located devices within the cell radius, where the BS-device distance $d_k \in (0.02 \text{ km}, 0.5 \text{ km})$ for each device. We study the test accuracy performance of the proposed methods over communication round T under different settings of (N, K). All results are obtained by averaging over 5 drops of device locations, each with 20 channel realizations. We first show the cumulative distribution functions (CDFs) of SNR_{dl} and SNR_{ul} among devices at donwlink and uplink, respectively, in Fig. 2. It shows that SNR_{dl} is in the range of (-20 dB, 40 dB). SNR_{ul} at uplink is in the range of (-30 dB, 35 dB), about $5 \sim 7 \text{ dB}$ worse than downlink. This is expected due to different maximum transmit power at the BS and devices.

Fig. 3-Left shows the test accuracy performance for (N,K)=(16,20). We see that JDUBF-TR converges fast and nearly attains the upper bound under the Ideal FL after 20 communication rounds. It achieves an accuracy above 95% at 100 rounds. JDUBF-R is worse than JDUBF-TR without uplink transmit beamforming, but it still achieves above 90% accuracy at 100 rounds. SDUBF has a much slower training convergence rate. After 100 rounds, it reaches test accuracy between 80%-85%. RBF again performs poorly without beamforming gain. No training convergence is observed for RBF, which leads to an accuracy $\sim 10\%$ for all rounds.

Fig. 3-Middle shows the case for (N,K)=(64,20). Compared with those in Fig. 3-Left, both JDUBF-R and SDUBF have noticeable performance improvement due to more BS antennas leading to higher beamforming again at downlink and uplink for improved wireless FL communication. This results in improved overall learning performance. The further performance improvement of JDUBF-TR is less noticeable, as it is already close to the Ideal FL for N=16 in Fig. 3-Left.

Fig. 3-Right shows the test accuracy for (N,K)=(64,10). Comparing Figs. 3-Middle and Right, we see that as K reduces from 20 to 10, the change of test accuracy of JDUBF-TR and JDUBF-R is hardly noticeable. However, that of SDUBF has a more significant reduction. The performance depredation is due to fewer devices resulting in less (distributed) transmit beamforming gain for uplink over-the-air aggregation. The learning performance of SDUBF under separate downlink and uplink design is more sensitive to the number of participating devices for aggregation. In summary, from Fig. 3 with different N and K values, we see that JDUBF-TR and JDUBF-R based on joint downlink-uplink beamforming design are effective methods to combat noisy and imperfect communication links to facilitate wireless FL under various system configurations. They achieve fast training convergence and high test accuracy.

VII. CONCLUSION

This paper considers the transceiver beamforming design for wireless FL experiencing fluctuated wireless links and noisy reception that degrade the learning performance. We propose a joint downlink-uplink beamforming design approach to maximize the FL training convergence over such a wireless environment. We first obtain the round-trip global model updating equation, which captures the impact of transmitter and receiver processing, noisy reception, and local model training. These factors are then reflected in the upper bound we derive on the expected global training loss after T rounds, and we formulate the joint downlink-uplink beamforming optimization problem to minimize this upper bound after T rounds. Depending on whether the knowledge of CSI is available at the devices, we propose two joint downlink-uplink beamforming methods, JDUBF-TR and JDUBF-R. Each method uses a greedy approach to solve a per-round joint online optimization problem, which is further decomposed into three subproblems that are solved alternatingly. PGD is used for each subproblem, which yields fast closed-from updates. We also propose a closedform initialization method to accelerate the fast convergence of JDUBF-TR and JDUBF-R. Simulation results show that our proposed methods substantially outperform the conventional separate-link-based beamforming design for various number of antennas and devices. In particular, JDUBF-TR nearly attains the learning performance of ideal FL with error-free communication links.

APPENDIX A PROOF OF LEMMA 1

Proof: For bounding $A_{1,t}$, we separately bound the two terms in (17). Specifically, let $A_{2,t}$ and $A_{3,t}$ denote the first and second terms in (17), respectively. For $A_{2,t}$, we have

$$\begin{split} &A_{2,t} \stackrel{(a)}{=} -\eta_t \sum_{\tau=0}^{J-1} \mathbb{E} \bigg[\bigg(\sum_{k=1}^K \mathfrak{Re}\{\rho_{k,t}\} \nabla F_k(\boldsymbol{\theta}_{k,t}^{\tau}) \bigg)^\mathsf{T} \nabla F(\boldsymbol{\theta}_t) \bigg] \\ &\stackrel{(b)}{\leq} -\frac{\eta_t}{2} \sum_{\tau=0}^{J-1} \mathbb{E} \left[\|\nabla F(\boldsymbol{\theta}_t)\|^2 \right] \\ &+ \frac{\eta_t}{2} \sum_{\tau=0}^{J-1} \mathbb{E} \bigg[\|\nabla F(\boldsymbol{\theta}_t) - \sum_{k=1}^K \mathfrak{Re}\{\rho_{k,t}\} \nabla F_k(\boldsymbol{\theta}_{k,t}^{\tau}) \bigg\|^2 \bigg] \\ &\stackrel{(c)}{\leq} -\frac{\eta_t J}{2} \mathbb{E} \left[\|\nabla F(\boldsymbol{\theta}_t)\|^2 \right] \\ &+ \eta_t \sum_{\tau=0}^{J-1} \mathbb{E} \bigg[\left\| \nabla F(\boldsymbol{\theta}_t) - \sum_{k=1}^K \mathfrak{Re}\{\rho_{k,t}\} \nabla F_k(\boldsymbol{\theta}_t) \bigg\|^2 \bigg] \\ &+ \eta_t \sum_{\tau=0}^{J-1} \mathbb{E} \bigg[\left\| \sum_{k=1}^K \mathfrak{Re}\{\rho_{k,t}\} (\nabla F_k(\boldsymbol{\theta}_t) - \nabla F_k(\boldsymbol{\theta}_{k,t}^{\tau})) \right\|^2 \bigg] \end{split}$$

$$\stackrel{(d)}{\leq} - \frac{\eta_{t}J}{2} \mathbb{E} \left[\|\nabla F(\boldsymbol{\theta}_{t})\|^{2} \right] + \eta_{t}J\delta
+ \eta_{t}(\beta_{t}^{\text{re}})^{2} \sum_{\tau=0}^{J-1} \mathbb{E} \left[\left\| \sum_{k=1}^{K} \frac{|\Re \mathfrak{e}\{\rho_{k,t}\}|}{\beta_{t}^{\text{re}}} \|\nabla F_{k}(\boldsymbol{\theta}_{t}) - \nabla F_{k}(\boldsymbol{\theta}_{k,t}^{\tau})\| \right\|^{2} \right]
\stackrel{(e)}{\leq} - \frac{\eta_{t}J}{2} \mathbb{E} \left[\|\nabla F(\boldsymbol{\theta}_{t})\|^{2} \right] + \eta_{t}J\delta
+ \eta_{t}\beta_{t}^{\text{re}} \sum_{k=1}^{K} |\Re \mathfrak{e}\{\rho_{k,t}\}| \sum_{\tau=0}^{J-1} \mathbb{E} \left[\left\|\nabla F_{k}(\boldsymbol{\theta}_{t}) - \nabla F_{k}(\boldsymbol{\theta}_{k,t}^{\tau})\right\|^{2} \right]
\stackrel{(f)}{\leq} - \frac{\eta_{t}J}{2} \mathbb{E} \left[\|\nabla F(\boldsymbol{\theta}_{t})\|^{2} \right] + \eta_{t}J\delta
+ \eta_{t}L^{2}\beta_{t}^{\text{re}} \sum_{k=1}^{K} |\Re \mathfrak{e}\{\rho_{k,t}\}| \sum_{\tau=0}^{J-1} \mathbb{E} \left[\|\boldsymbol{\theta}_{t} - \boldsymbol{\theta}_{k,t}^{\tau}\|^{2} \right]$$

$$(40)$$

where (a) follows $\Delta \boldsymbol{\theta}_{k,t} = -\eta_t \sum_{\tau=0}^{J-1} \nabla F_k(\boldsymbol{\theta}_{k,t}^{\tau}; \mathcal{B}_{k,t}^{\tau})$ and the unbiasedness of the mini-batch SGD from Assumption 2, (b) and (c) are based on $\|\mathbf{x} + \mathbf{y}\|^2 \leq 2\|\mathbf{x}\|^2 + 2\|\mathbf{y}\|^2, \forall \mathbf{x}, \mathbf{y}, (d)$ follows Assumption 3, (e) is based on the Jensen's inequality, and (f) follows the L-smoothness of $F_k(\cdot)$ from Assumption 1. Let $A_{4,t}$ denote $\sum_{\tau=0}^{J-1} \mathbb{E}[\|\boldsymbol{\theta}_t - \boldsymbol{\theta}_{k,t}^{\tau}\|^2]$ in (40). Then,

$$\begin{split} &A_{4,t} = \mathbb{E}\left[\|\boldsymbol{\theta}_{t} - \boldsymbol{\theta}_{k,t}^{0}\|^{2}\right] + \sum_{\tau=1}^{J-1} \mathbb{E}\left[\|\boldsymbol{\theta}_{t} - \boldsymbol{\theta}_{k,t}^{\tau}\|^{2}\right] \\ &\leq 2\sum_{\tau=0}^{J-1} \mathbb{E}\left[\|\hat{\mathbf{n}}_{k,t}^{\text{dl}}\|^{2}\right] + 2\sum_{\tau=1}^{J-1} \mathbb{E}\left[\left\|\eta_{t}\sum_{\tau'=0}^{\tau-1} \nabla F_{k}(\boldsymbol{\theta}_{k,t}^{\tau'}; \boldsymbol{\mathcal{B}}_{k,t}^{\tau'})\right\|^{2}\right] \\ &\leq \frac{JD\sigma_{d}^{2}}{|(\mathbf{w}_{t}^{\text{dl}})^{\mathsf{H}}\mathbf{h}_{k,t}|^{2}} + 2\eta_{t}^{2}\sum_{\tau=1}^{J-1} \tau\sum_{\tau'=0}^{\tau-1} \mathbb{E}\left[\left\|\nabla F_{k}(\boldsymbol{\theta}_{k,t}^{\tau'}; \boldsymbol{\mathcal{B}}_{k,t}^{\tau'})\right\|^{2}\right] \\ &\stackrel{(a)}{=} \frac{JD\sigma_{d}^{2}}{|(\mathbf{w}_{t}^{\text{dl}})^{\mathsf{H}}\mathbf{h}_{k,t}|^{2}} + 2\eta_{t}^{2}\sum_{\tau=1}^{J-1} \tau\sum_{\tau'=0}^{\tau-1} \mathbb{E}\left[\left\|\nabla F_{k}(\boldsymbol{\theta}_{k,t}^{\tau'}; \boldsymbol{\mathcal{B}}_{k,t}^{\tau'}) - \nabla F_{k}(\boldsymbol{\theta}_{k,t}^{\tau'})\right\|^{2}\right] \\ &+ 2\eta_{t}^{2}\sum_{\tau=1}^{J-1} \tau\sum_{\tau'=0}^{\tau-1} \mathbb{E}\left[\left\|\nabla F_{k}(\boldsymbol{\theta}_{k,t}^{\tau'}; \boldsymbol{\mathcal{B}}_{k,t}^{\tau'}) - \nabla F_{k}(\boldsymbol{\theta}_{k,t}^{\tau'})\right\|^{2}\right] \\ &\stackrel{(b)}{\leq} \frac{JD\sigma_{d}^{2}}{|(\mathbf{w}_{t}^{\text{dl}})^{\mathsf{H}}\mathbf{h}_{k,t}|^{2}} + 2\eta_{t}^{2}\sum_{\tau=1}^{J-1} \tau\sum_{\tau'=0}^{\tau-1} \mathbb{E}\left[\left\|\nabla F_{k}(\boldsymbol{\theta}_{k,t}^{\tau'}; \boldsymbol{\mathcal{B}}_{k,t}^{\tau'}) - \nabla F_{k}(\boldsymbol{\theta}_{k,t}^{\tau'})\right\|^{2}\right] + 2\eta_{t}^{2}J^{3}\mu \end{split}$$

where (a) and (b) follow Assumption 2. Let $A_{5,t}$ denote $\sum_{\tau'=0}^{\tau-1} \mathbb{E}[\|\nabla F_k(\boldsymbol{\theta}_{k,t}^{\tau'})\|^2]$, which is bounded by

$$A_{5,t} \leq 2 \sum_{\tau'=0}^{\tau-1} \mathbb{E}\left[\left\|\nabla F_k(\boldsymbol{\theta}_{k,t}^{\tau'}) - \nabla F_k(\boldsymbol{\theta}_t)\right\|^2\right] + 4J\mathbb{E}\left[\left\|\nabla F(\boldsymbol{\theta}_t)\right\|^2\right] + 4J\mathbb{E}\left[\left\|\nabla F_k(\boldsymbol{\theta}_t) - \nabla F(\boldsymbol{\theta}_t)\right\|^2\right]$$

$$\leq 2L^2 \sum_{\tau'=0}^{\tau-1} \mathbb{E}\left[\left\|\boldsymbol{\theta}_{k,t}^{\tau'} - \boldsymbol{\theta}_t\right\|^2\right] + 4J\mathbb{E}\left[\left\|\nabla F(\boldsymbol{\theta}_t)\right\|^2\right] + 4J\delta. \tag{41}$$

Next, for $A_{3,t}$, based on Assumption 4, we obtain $A_{3,t} \leq \sum_{k=1}^K |\mathfrak{Im}\{\rho_{k,t}\}| \, \mathbb{E}\left[\left\|\Delta \bar{\boldsymbol{\theta}}_{k,t}\right\| \, \|\nabla F(\boldsymbol{\theta}_t)\|\right] \leq \nu \zeta \beta_t^{\mathrm{im}}$. Combining the above bounds with (17), we have (18).

APPENDIX B PROOF OF LEMMA 2

Proof: For bounding $B_{1,t}$, we bound the first term in (19), denoted by $B_{2,t}$. Specifically, we have

$$B_{2,t} \leq 2\mathbb{E}\left[\left\|\sum_{k=1}^{K} \rho_{k,t} \Delta \tilde{\boldsymbol{\theta}}_{k,t}\right\|^{2}\right] + 2\mathbb{E}\left[\left\|\sum_{k=1}^{K} \rho_{k,t} \tilde{\mathbf{n}}_{k,t}^{\text{dl}}\right\|^{2}\right]$$

$$= 2\mathbb{E}\left[\left\|\sum_{k=1}^{K} \mathfrak{Re}\{\rho_{k,t}\}\Delta\tilde{\boldsymbol{\theta}}_{k,t}\right\|^{2}\right] + 2\mathbb{E}\left[\left\|\sum_{k=1}^{K} \mathfrak{Im}\{\rho_{k,t}\}\Delta\tilde{\boldsymbol{\theta}}_{k,t}\right\|^{2}\right] + D\sigma_{d}^{2} \sum_{k=1}^{K} \frac{|\rho_{k,t}|^{2}}{|(\mathbf{w}_{t}^{dl})^{\mathsf{H}}\mathbf{h}_{k,t}|^{2}}.$$

$$(42)$$

Let $B_{3,t}$ and $B_{4,t}$ respectively denote the first and second terms in (42). Let $B_{3,t}$ and $B_{4,t}$ respectively denote the first and second terms in (42). For $B_{3,t}$, we have

$$B_{3,t} \leq 2(\beta_t^{\text{re}})^2 \mathbb{E} \left[\left\| \sum_{k=1}^K \frac{|\mathfrak{Re}\{\rho_{k,t}\}|}{\beta_t^{\text{re}}} \right\| \Delta \tilde{\boldsymbol{\theta}}_{k,t} \right\| \right\|^2 \right]$$

$$\stackrel{(a)}{\leq} 2\beta_t^{\text{re}} \sum_{k=1}^K |\mathfrak{Re}\{\rho_{k,t}\}| \mathbb{E}[\|\Delta \boldsymbol{\theta}_{k,t}\|^2]$$

$$(43)$$

where (a) applies the Jensen's inequality. We now bound the term $\mathbb{E}[\left\|\Delta\theta_{k,t}\right\|^2]$ in (43), given by

$$\mathbb{E}\left[\left\|\Delta\boldsymbol{\theta}_{k,t}\right\|^{2}\right] = \eta_{t}^{2}\mathbb{E}\left[\left\|\sum_{\tau=0}^{J-1}\nabla F_{k}(\boldsymbol{\theta}_{k,t}^{\tau}; \mathcal{B}_{k,t}^{\tau})\right\|^{2}\right]$$

$$\stackrel{(a)}{=} \eta_{t}^{2}\sum_{\tau=0}^{J-1}\mathbb{E}\left[\left\|\nabla F_{k}(\boldsymbol{\theta}_{k,t}^{\tau}; \mathcal{B}_{k,t}^{\tau}) - \nabla F_{k}(\boldsymbol{\theta}_{k,t}^{\tau})\right\|^{2}\right]$$

$$+ \eta_{t}^{2}\mathbb{E}\left[\left\|\sum_{\tau=0}^{J-1}\nabla F_{k}(\boldsymbol{\theta}_{k,t}^{\tau})\right\|^{2}\right]$$

$$\stackrel{(b)}{\leq} \eta_{t}^{2}J\mu + 2\eta_{t}^{2}\mathbb{E}\left[\left\|\sum_{\tau=0}^{J-1}(\nabla F_{k}(\boldsymbol{\theta}_{k,t}^{\tau}) - \nabla F_{k}(\boldsymbol{\theta}_{t}))\right\|^{2}\right]$$

$$+ 2\eta_{t}^{2}\mathbb{E}\left[\left\|\sum_{\tau=0}^{J-1}\nabla F_{k}(\boldsymbol{\theta}_{t})\right\|^{2}\right]$$

$$(44)$$

where (a) and (b) are from Assumption 2. Let $B_{5,t}$ and $B_{6,t}$ respectively denote the second and third terms in (44). For $B_{5,t}$, we have

$$B_{5,t} \leq 2\eta_{t}^{2} J \sum_{\tau=0}^{J-1} \mathbb{E}\left[\|\nabla F_{k}(\boldsymbol{\theta}_{k,t}^{\tau}) - \nabla F_{k}(\boldsymbol{\theta}_{t})\|^{2}\right]$$

$$\stackrel{(a)}{\leq} 2\eta_{t}^{2} J L^{2} \sum_{\tau=0}^{J-1} \mathbb{E}\left[\|\boldsymbol{\theta}_{k,t}^{\tau} - \boldsymbol{\theta}_{t}\|^{2}\right]$$

$$\stackrel{(b)}{\leq} \frac{(1 - Q_{t})^{2}}{L^{2} Q_{t}} \mathbb{E}[\|\nabla F(\boldsymbol{\theta}_{t})\|^{2}] + \frac{(1 - Q_{t}) D \sigma_{d}^{2}}{2 Q_{t} |(\mathbf{w}_{t}^{dl})^{\mathsf{H}} \mathbf{h}_{k,t}|^{2}}$$

$$+ \frac{(1 - Q_{t})^{2}}{2 L^{2} Q_{t}} (4\delta + \mu) \tag{45}$$

where (a) follows the L-smoothness of $F_k(\cdot)$ in Assumption 1, and (b) applies the bounds on $A_{4,t}$. For $B_{6,t}$, we have

$$B_{6,t} = 2\eta_t^2 J^2 \mathbb{E} \left[\|\nabla F_k(\boldsymbol{\theta}_t)\|^2 \right]$$

$$\leq 4\eta_t^2 J^2 \mathbb{E} \left[\|\nabla F_k(\boldsymbol{\theta}_t) - \nabla F(\boldsymbol{\theta}_t)\|^2 \right] + 4\eta_t^2 J^2 \mathbb{E} \left[\|\nabla F(\boldsymbol{\theta}_t)\|^2 \right]$$

$$\stackrel{(a)}{\leq} 4\eta_t^2 J^2 \delta + 4\eta_t^2 J^2 \mathbb{E} \left[\|\nabla F(\boldsymbol{\theta}(t))\|^2 \right]$$
(46)

where (a) follows Assumption 3. Next, for $B_{4,t}$, based on Assumption 4, we obtain $B_{4,t} \leq 2\mathbb{E}[\|\sum_{k=1}^K |\mathfrak{Im}\{\rho_{k,t}\}|\|\Delta\tilde{\theta}_{k,t}\|\|^2] \leq 2\nu^2(\beta_t^{\mathrm{im}})^2$. Combining the above bounds with (42)(19), we have (20).

APPENDIX C PROOF OF PROPOSITION 1

Proof: We apply Lemmas 1 and 2 to (21). Let $Y_t \triangleq$ Thosp. We apply Echimas 1 and 2 to (21). Let $T_t = \frac{\eta_t J}{2} \left(4V_t(\beta_t^{\rm re})^2 - 1 \right)$. Using the definition of V_t , we have $Y_t = \frac{\eta_t J}{2Q_t} \left(4(1-Q_t)(\beta_t^{\rm re})^2 + 4\sqrt{1-Q_t}(\beta_t^{\rm re})^2 - Q_t \right)$. Letting $x \triangleq \sqrt{1-Q_t}$, we have

$$Y_t = \frac{\eta_t J}{2Q_t} \left(\left(4(\beta_t^{\text{re}})^2 + 1 \right) x^2 + 4(\beta_t^{\text{re}})^2 x - 1 \right).$$

We first obtain $Q_t>0$ for $\eta_t J<\frac{1}{2L}$. Next, the positive root of equation $\left(4(\beta_t^{\rm re})^2+1\right)x^2+4(\beta_t^{\rm re})^2x-1=0$ is given by

$$x^{o} = \frac{-4(\beta_{t}^{\text{re}})^{2} + \sqrt{16(\beta_{t}^{\text{re}})^{4} + 16(\beta_{t}^{\text{re}})^{2} + 4}}{8(\beta_{t}^{\text{re}})^{2} + 2}$$
$$= \frac{1}{4(\beta_{t}^{\text{re}})^{2} + 1}.$$

Thus, for $x \geq x^o$, we have $(4(\beta_t^{\rm re})^2 + 1)x^2 + 4(\beta_t^{\rm re})^2x - 1 \geq 0$. Re-express $x \geq x^o$ as $\eta_t J \geq \frac{1}{2L(4(\beta_t^{\rm re})^2 + 1)}$. Note that $\beta_t^{\rm re} = \sum_{k=1}^K |\Re \{ \rho_{k,t} \}| \geq |\Re \{ \sum_{k=1}^K \rho_{k,t} \}| = 1$, leading to $\frac{1}{2L(4(\beta_t^{\rm re})^2 + 1)} < \frac{1}{2L}$. Based on these, we have $Y_t \geq 0$ for $\frac{1}{2L(4(\beta_t^{\rm re})^2 + 1)} \leq \eta_t J < \frac{1}{2L}$.

Let $Z_t \triangleq H(\mathbf{w}_t^{\text{dl}}, \boldsymbol{\rho}_t) + C_t$. We have $Z_t > 0$ due to $H(\mathbf{w}_t^{\text{dl}}, \boldsymbol{\rho}_t) > 0$ and $C_t > 0$. Then, after combining (18), (20), and (21), for $\frac{1}{2L(4(\beta_t^{\text{re}})^2 + 1)} \leq \eta_t J < \frac{1}{2L}$, we have

$$\mathbb{E}[F(\boldsymbol{\theta}_{t+1})] - F^{\star} \leq \mathbb{E}[F(\boldsymbol{\theta}_{t})] - F^{\star} + Y_{t}\mathbb{E}[\|\nabla F(\boldsymbol{\theta}_{t})\|^{2}] + Z_{t}$$

$$\stackrel{(a)}{\leq} \mathbb{E}[F(\boldsymbol{\theta}_{t})] - F^{\star} + L^{2}Y_{t}\mathbb{E}[\|\boldsymbol{\theta}_{t} - \boldsymbol{\theta}^{\star}\|^{2}] + Z_{t}$$

$$\stackrel{(b)}{\leq} \mathbb{E}[F(\boldsymbol{\theta}_{t})] - F^{\star} + \frac{2L^{2}Y_{t}}{\lambda} (\mathbb{E}[F(\boldsymbol{\theta}_{t})] - F^{\star}) + Z_{t}$$

$$(47)$$

where (a) and (b) follow from (1) and the L-smoothness and λ -strong-convexity of $F_k(\cdot)$ in Assumption 1, respectively. Summing up both sides of (47) over $t \in \mathcal{T}$ and rearranging the terms, we have (22).

REFERENCES

- [1] C. Zhang, M. Dong, B. Liang, A. Afana, and Y. Ahmed, "Joint downlinkuplink beamforming for wireless multi-antenna federated learning," in Proc. WiOpt, Aug. 2023, pp. 1-8.
- [2] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. Y. Arcas, "Communication-efficient learning of deep networks from decentralized data," in Proc. AISTATS, Apr. 2017, pp. 1273-1282
- G. Zhu, D. Liu, Y. Du, C. You, J. Zhang, and K. Huang, "Toward an intelligent edge: Wireless communication meets machine learning," IEEE
- Commun. Mag., vol. 58, no. 1, pp. 19–25, Jan. 2020.

 Y. Du, S. Yang, and K. Huang, "High-dimensional stochastic gradient quantization for communication-efficient edge learning," *IEEE Trans*. Signal Process., vol. 68, pp. 2128-2142, Mar. 2020.
- Z. Yang, M. Chen, W. Saad, C. S. Hong, and M. Shikh-Bahaei, "Energy efficient federated learning over wireless communication networks," IEEE Trans. Wireless Commun., vol. 20, no. 3, pp. 1935–1949, Mar. 2021.
 [6] J. Xu and H. Wang, "Client selection and bandwidth allocation in
- wireless federated learning networks: A long-term perspective," IEEE Trans. Wireless Commun., vol. 20, no. 2, pp. 1188-1200, Feb. 2021.
- M. M. Amiri, D. Gündüz, S. R. Kulkarni, and H. V. Poor, "Convergence of update aware device scheduling for federated learning at the wireless edge," IEEE Trans. Wireless Commun., vol. 20, no. 6, pp. 3643-3658,
- [8] Y. Wang, Y. Xu, Q. Shi, and T.-H. Chang, "Quantized federated learning under transmission delay and outage constraints," IEEE J. Sel. Areas
- Commun., vol. 40, no. 1, pp. 323–341, Jan. 2022. 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.
- N. Zhang and M. Tao, "Gradient statistics aware power control for overthe-air federated learning," IEEE Trans. Wireless Commun., vol. 20, no. 8, pp. 5115-5128, Aug. 2021.

- [11] J. Wang, B. Liang, M. Dong, G. Boudreau, and H. Abou-Zeid, "Joint online optimization of model training and analog aggregation for wireless edge learning," IEEE/ACM Trans. Netw., vol. 32, no. 2, pp. 1212-1228, Apr. 2024.
- [12] L. Qu, S. Song, C.-Y. Tsui, and Y. Mao, "How robust is federated learning to communication error? A comparison study between uplink and downlink channels," in Proc. IEEE WCNC, Jul. 2024, pp. 1-6.
- M. M. Amiri, D. Gündüz, S. R. Kulkarni, and H. V. Poor, "Convergence of federated learning over a noisy downlink," IEEE Trans. Wireless Commun., vol. 21, no. 3, pp. 1422–1437, Mar. 2022.
 [14] X. Wei and C. Shen, "Federated learning over noisy channels: Conver-
- gence analysis and design examples," IEEE Trans. Cogn. Commun. Netw., vol. 8, no. 2, pp. 1253-1268, Jun. 2022.
- W. Guo, R. Li, C. Huang, X. Qin, K. Shen, and W. Zhang, "Joint device selection and power control for wireless federated learning, Areas Commun., vol. 40, no. 8, pp. 2395-2410, Aug. 2022
- S. M. Shah, L. Su, and V. K. N. L. Lau, "Robust federated learning over noisy fading channels," *IEEE Internet Things J.*, vol. 10, no. 9, pp. 7993–
- 8013, May 2023. [17] Z. Wang, Y. Zhou, Y. Shi, and W. Zhuang, "Interference management for over-the-air federated learning in multi-cell wireless networks," IEEE J. Sel. Areas Commun., vol. 40, no. 8, pp. 2361–2377, Aug. 2022.
 [18] L. Chen, N. Zhao, Y. Chen, F. R. Yu, and G. Wei, "Over-the-air
- computation for IoT networks: Computing multiple functions with antenna arrays," IEEE Internet Things J., vol. 5, no. 6, pp. 5296-5306, Dec. 2018.
- G. Zhu, L. Chen, and K. Huang, "MIMO over-the-air computation: Beamforming optimization on the Grassmann manifold," in *Proc. IEEE*
- GLOBECOM, Dec. 2018, pp. 1–6. K. Yang, T. Jiang, Y. Shi, and Z. Ding, "Federated learning via overthe-air computation," *IEEE Trans. Wireless Commun.*, vol. 19, no. 3, pp. 2022-2035, Mar. 2020.
- H. Liu, X. Yuan, and Y.-J. A. Zhang, "Reconfigurable intelligent surface enabled federated learning: A unified communication-learning design approach," IEEE Trans. Wireless Commun., vol. 20, no. 11, pp. 7595-7609, Nov. 2021.
- M. Kim, A. L. Swindlehurst, and D. Park, "Beamforming vector design and device selection in over-the-air federated learning," *IEEE Trans.* and device selection in over-the-air federated learning, Wireless Commun., vol. 22, no. 11, pp. 7464–7477, Nov. 2023.
 [23] F. M. Kalarde, B. Liang, M. Dong, Y. A. E. Ahmed, and H. T. Cheng,
- "Power minimization in federated learning with over-the-air aggregation and receiver beamforming," in Proc. MSWiM, Oct. 2023, p. 259-267.
- F. M. Kalarde, M. Dong, B. Liang, Y. A. E. Ahmed, and H. T. Cheng, "Beamforming and device selection design in federated learning with overthe-air aggregation," IEEE Open J. Commun. Soc., vol. 5, pp. 1710-1723, Mar. 2024.
- C. Zhang, M. Dong, B. Liang, A. Afana, and Y. Ahmed, "Uplink overthe-air aggregation for multi-model wireless federated learning," in Proc. IEEE SPAWC, Sept. 2024, pp. 36-40.
- -, "Multi-model wireless federated learning with downlink beamforming," in Proc. IEEE ICASSP, Apr. 2024, pp. 9146-9150.
- G. Zhu, Y. Du, D. Gündüz, 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.
- A. Şahin, "Over-the-air computation based on balanced number systems for federated edge learning," *IEEE Trans. Wireless Commun.*, vol. 23, no. 5, pp. 4564–4579, May 2024.
- [29] L. Qiao, Z. Gao, M. B. Mashhadi, and D. Gündüz, "Massive digital overthe-air computation for communication-efficient federated edge learning, IEEE J. Sel. Areas Commun., vol. 42, no. 11, pp. 3078–3094, Nov. 2024.
 [30] N. D. Sidiropoulos, T. N. Davidson, and Z.-Q. Luo, "Transmit beamform
- ing for physical-layer multicasting," IEEE Trans. Signal Process., vol. 54, no. 6, pp. 2239-2251, Jun. 2006.
- [31] M. Dong and Q. Wang, "Multi-group multicast beamforming: Optimal structure and efficient algorithms," *IEEE Trans. Signal Process.*, vol. 68, pp. 3738-3753, May 2020.
- Zhang, M. Dong, and B. Liang, "Ultra-low-complexity algorithms with structurally optimal multi-group multicast beamforming in large-scale systems," IEEE Trans. Signal Process., vol. 71, pp. 1626–1641, Apr. 2023.
- [33] D. Bertsekas, Nonlinear Programming. Athena Scientific, 2016.
- [34] E. S. Levitin and B. T. Polyak, "Constrained minimization methods," USSR Comput. Math. Math. Phys., vol. 6, no. 5, pp. 787-823, 1966.
- [35] S. Bubeck, "Convex optimization: Algorithms and complexity," Found. Trends Mach. Learn., vol. 8, no. 3-4, pp. 231–357, 2015.
 [36] A. Mokhtari, A. Ozdaglar, and A. Jadbabaie, "Escaping saddle points in
- constrained optimization," in Proc. NeurIPS, Dec. 2018, pp. 3629-3639.
- C. Zhang, M. Dong, and B. Liang, "Fast first-order algorithm for largescale max-min fair multi-group multicast beamforming," IEEE Wireless Commun. Lett., vol. 11, no. 8, pp. 1560–1564, Aug. 2022. Y. Yang, J. Xu, G. Shi, and C.-X. Wang, 5G wireless systems. Springer,
- 2018
- Y. LeCun, C. Cortes, and C. Burges. The MNIST Database of Handwritten Digits. 1998. [Online]. Available: http://yann.lecun.com/exdb/mnist/.