

Fast Federated Learning by Balancing Communication Trade-Offs

Milad Khademi Nori, Sangseok Yun, and Il-Min Kim, *Senior Member, IEEE*

Abstract—Federated Learning (FL) has recently received a lot of attention for large-scale privacy-preserving machine learning. However, high communication overheads due to frequent gradient transmissions decelerate FL. To mitigate the communication overheads, two main techniques have been studied: (i) local update of weights characterizing the trade-off between communication and computation and (ii) gradient compression characterizing the trade-off between communication and precision. To the best of our knowledge, studying and balancing those two trade-offs *jointly and dynamically* while considering their impacts on convergence has remained unresolved even though it promises significantly faster FL. In this paper, we first formulate our problem to minimize learning error with respect to two variables: *local update coefficients* and *sparsity budgets* of gradient compression who characterize trade-offs between communication and computation/precision, respectively. We then derive an upper bound of the learning error in a given wall-clock time considering the interdependency between the two variables. Based on this theoretical analysis, we propose an enhanced FL scheme, namely Fast FL (FFL), that jointly and dynamically adjusts the two variables to minimize the learning error. We demonstrate that FFL consistently achieves higher accuracies faster than similar schemes existing in the literature.

Index Terms—Communication overhead, communication trade-off, federated learning, gradient compression, local update.

I. INTRODUCTION

FEDERATED Learning (FL) is a training paradigm where a large number of workers collectively train a model using Stochastic Gradient Descent (SGD) [1]. Each worker holds a local (often private) training dataset, with which it contributes to training a global model. Specifically, in each FL round, each worker computes the error of the global model by forward-propagating the samples of its own dataset through the model, followed by comparing the model's outputs with samples' labels and applying an appropriate loss function. By backpropagating this error, each worker calculates the gradients. All workers then send their gradients to the server, who updates the global model and sends the updated model back to all the workers. At this point, one FL round is completed. After a number of (FL) rounds, the global model converges [2].

FL has recently become popular for the following reasons: (i) FL enables privacy-preserving training as workers only share their gradients with the server. That is, FL makes it possible to train a model from the entire datasets of workers without any individuals having to reveal their local datasets to other workers or the server. This characteristic of FL technology

renders it apposite to many applications, the most popular of which is Internet of Things (IoT) where the data from organizations, hospitals, homes (Smart Home), and vehicles (Internet of Vehicles) presumably needs to be collected for processing, inference, and decision-making [3], [4]. FL removes the need to collect the raw data in using machine learning for IoT, thereby it ameliorates the privacy of IoT. (ii) FL is an exemplar of Edge Intelligence, which is the offspring from the union of (mobile) Edge Computing and Artificial Intelligence (AI). Edge Intelligence advocates for pushing the intelligence down from clouds to the proximity of end-users by implementation of AI applications on edge-devices¹, thereby promising better privacy and reliability with lower latency and cost. FL, correspondingly, can be used to facilitate training of large-scale models by parallelizing the computation of gradients on edge-devices. Moreover, FL allows inference to be performed on edge-devices, and thus, fulfills the goals of Edge Intelligence [5]. (iii) FL addresses the intelligence democratization issue: in the AI market, big companies possessing a large amount of data plus computational and storage facilities can get the monopoly of AI businesses while small firms have no chances in this competition. FL enables small businesses to play a role in the market as they do not require a massive amount of data at a datacenter, nor computational and storage facilities.

Unfortunately, the empirical speedup (for training) offered by FL often fails to meet the optimal scaling (in the number of workers) that is ideally desired. It is now widely acknowledged that this speedup saturation is mainly due to the communication overheads, which are largely attributed to frequent gradient transmissions from workers to the server. As the number of parameters in the state-of-the-art models scales to a huge number (e.g., hundreds of millions), the size of the gradients scales proportionally. The communication bottleneck becomes even more pronounced particularly when the workers performing FL are wireless devices (e.g., smartphones and sensors) which communicate through wireless channels and suffer from low-bandwidth, intermittent connections, and expensive mobile data plans [6]–[9]. To mitigate the communication overhead problem and hence speedup learning, in the literature, two main techniques have been studied. The first is local update, characterizing the trade-off between communication and computation, and the second is gradient compression, characterizing the trade-off between communication and precision. Each of these techniques will be discussed in the following.

In the local update technique, within each round, many local updates are performed—instead of only one update. Local update

Milad Khademi Nori and Il-Min Kim are with the Department of Electrical and Computer Engineering, Queen's University, Kingston, ON K7L 3N6, CA (e-mail: 19mkn1@queensu.ca; ilmin.kim@queensu.ca). Sangseok Yun is with the Department of Information and Communications Engineering, Pukyong National University, Busan 48513, South Korea (e-mail: ssyun@pknu.ac.kr).

¹While some papers in literature of FL use the term “edge-device” as to refer to the nodes at which the gradients are calculated, in this paper, we prefer to use the term “worker.”

presents a communication-computation trade-off as *local update coefficient* determines the ratio of computation to communication. The authors of [7]–[9] showed that local updates can significantly alleviate frequent gradient transmissions. Formulating the error-convergence bound in terms of local updates, Adaptive Communication (ADACOMM) [8] dynamically optimized the local update coefficients to accelerate convergence while retaining accuracy. ADACOMM also showed that the communication-computation trade-off introduced by local update coefficients should be dynamically balanced over the course of learning. In [10], the same issue was addressed, but assuming that the loss function was convex. Performing binary compression and benefiting from local updates, Sparse Binary Compression (SBC) [11] could significantly decrease communication overheads.

Gradient compression presents the communication-precision trade-off, i.e., how much compression is productive in different rounds? Gradient compression methods fall into two categories: gradient quantization and sparsification. Gradient quantization approaches the problem of communication overhead by sending a compressed (quantized) version of the gradient. TernGrad [12] used 2 bits for each element of the gradient and authors mathematically proved the convergence of TernGrad under the assumption of a bound on gradients. Quantized SGD (QSGD) [13] investigated the effects of 2 bits, 4 bits, and 8 bits quantizations on different layers of neural networks and showed that the networks converge. Authors of [14] even went further and quantized the gradients to one bit. In [15] and [16], similar quantization techniques were exploited.

Sparsification technique aims to sparsify gradients in order to send only the significant elements of the gradients instead of all. Sparsification techniques fall into two categories depending on in what domain the sparsity is sought for. Some researchers assume that gradients are sparse in their original domain while others find some transformed domains more appropriate. Following the first approach, the scheme proposed in [17] sparsified gradients by setting insignificant entries to zero, where insignificant entries were those entries whose values were between top 0.05% and bottom 0.05%. The authors then used run-length codes to encode the gradients before sending them to the server. In [11], a similar scheme was studied.

The second approach for sparsification aims to find another domain for gradients in which they are more sparse than their original domain [4]. In [18], gradients were transformed to the Singular Value Decomposition (SVD) domain to discover/exploit more sparsity helping to mitigate the communication overheads. Also, the authors of [18] stated that TernGrad [12] and QSGD [13] were special versions of their proposed scheme under certain circumstances. The authors in [19] proved why in general, big matrices are often low-rank.

In the literature, each of the aforementioned techniques, i.e., local updates and gradient compression, has been *individually* demonstrated to mitigate the communication overheads in FL, leading to considerable speedups. Specifically, dynamically determining local update coefficients to balance the communication-computation trade-off has been studied in [8]–[10], but without considering any compression of gradients. Meanwhile, *static* gradient compression (communication-precision trade-off) has been studied in [4], [18], but without dynamically adjusting local update coefficients. Based on these works, it is expected that

conjoining the two techniques would be even more effective in reducing the communication overheads, thereby accelerating FL. To the best of our knowledge, however, studying and balancing those two trade-offs *jointly and dynamically* while considering their impacts on convergence (from both theoretical and empirical perspectives) have not been done in the literature. This unexplored, yet important problem motivated our work. In this paper, we propose such jointly adjusted/balanced scheme, namely Fast FL (FFL), which jointly and dynamically determines the *local update coefficients* and *sparsity budgets* of gradient compression. Our main contributions in this paper are as follows:

- As the first work in the literature, we formulate our FL problem to minimize the error of the global model in a given wall-clock time with respect to *both* local update coefficients and sparsity budgets, which are respectively demonstrated to characterize the trade-off between communication and computation and the trade-off between communication and precision.
- We derive an upper bound of the error of the global model in a given wall-clock time considering the interdependency between the two variables: the local update coefficients and sparsity budgets.
- Using the derived error upper bound, we propose an enhanced FL scheme, FFL, which jointly and dynamically determines the two variables to accelerate learning.
- Our analytical results include almost all of the unbiased compression techniques. This is because our formulation for compression is based on atomic decomposition for sparse representation which is a popular technique for compression from compressed sensing.
- We demonstrate that FFL consistently achieves higher accuracies faster than similar schemes existing in the literature.

The remainder of this paper is organized as follows: in Section II, the fundamental mechanism of FL is presented and our proposed scheme is described at a high level. In Section III, our problem is mathematically formulated, followed by derivation of an error upper bound of the learning error. The FFL scheme is proposed in Section IV. Experiment results are provided in Section V and this paper is concluded in Section VI.

Notations: All vectors are column vectors and denoted by bold font small letters (e.g., \mathbf{x}), while scalars are denoted by normal font small letters (e.g., y). Matrices are denoted by bold font capital letters (e.g., \mathbf{A}). Also, \mathbf{x}^T denotes the transpose of \mathbf{x} . We use “:=” to denote “is defined to be equal to.” We also use $\|\cdot\|_1$ and $\|\cdot\|$ to denote the L_1 and L_2 norms, respectively. For a set \mathcal{S} , $|\mathcal{S}|$ denotes its cardinality. We use $\nabla F(\mathbf{x})$ to represent the gradient of $F(\mathbf{x})$. Expectation with respect to random variable X is denoted by $\mathbb{E}_X[\cdot]$.

II. FL MECHANISM AND THE PROPOSED SCHEME

In this section, we first present the fundamental learning mechanism of FL and then high level descriptions of our proposed scheme.

A. FL Mechanism

Consider an N -sample-size dataset defined by $\mathcal{S} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$, where the pair of (\mathbf{x}_q, y_q)

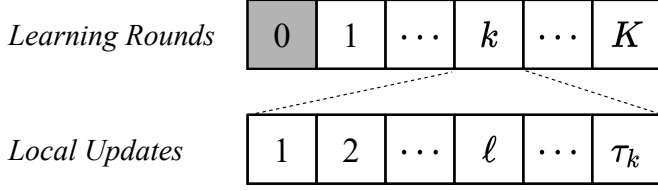


Fig. 1: Two time scales of FL. Each round comprises τ_k number of local updates. The FL system is initialized at $k = 0$.

includes the q th data sample \mathbf{x}_q and its corresponding label y_q . In the FL setting, the dataset \mathcal{S} is distributed among M distinct workers, each of which holds the local dataset $\mathcal{S}_j \subset \mathcal{S}$ for $j = 1, 2, \dots, M$, where $\bigcup_{j=1}^M \mathcal{S}_j = \mathcal{S}$ and $\mathcal{S}_{i'} \cap \mathcal{S}_{j'} = \emptyset$, for $i' \neq j'$. In the k th round, the global loss function $F(\mathbf{w}_k; \mathcal{S})$ on all of the distributed datasets is given by

$$F(\mathbf{w}_k; \mathcal{S}) := \frac{1}{|\mathcal{S}|} \sum_{(\mathbf{x}_q, y_q) \in \mathcal{S}} f(\mathbf{w}_k; \mathbf{x}_q, y_q) \quad (1)$$

for $k = 0, 1, 2, \dots, K$, where $\mathbf{w}_k \in \mathbb{R}^d$ is the global weight vector in the k th round for the global model, d denotes the dimension of the weight vector, and K is the last round. Also, $f(\mathbf{w}_k; \mathbf{x}_q, y_q) \in \mathbb{R}$ stands for the value of loss function for the q th data sample. Because \mathcal{S}_j 's are distributed over multiple distinct workers, it is not possible for the server to directly minimize the global loss function in (1). Instead, in FL, each worker minimizes its own local loss function $F(\mathbf{w}_k^j; \mathcal{S}_j)$, which is again defined by (1) but with local dataset \mathcal{S}_j and local weight vector $\mathbf{w}_k^j \in \mathbb{R}^d$ for its local model.

The FL learning process proceeds iteratively as follows: at the initialization stage ($k = 0$), all local weight vectors \mathbf{w}_0^j 's at different workers are initialized to the same value. At the beginning of each round (for $k > 0$), new values of local weight vectors are computed by performing τ_k number of consecutive local updates via SGD at the workers. Specifically, at the ℓ th local update, the local weight vector is updated from $\mathbf{w}_k^{j, \ell-1}$ to $\mathbf{w}_k^{j, \ell}$ as follows:

$$\mathbf{w}_k^{j, \ell} = \mathbf{w}_k^{j, \ell-1} - \eta \mathbf{g}(\mathbf{w}_k^{j, \ell-1}; \xi_j) \quad (2)$$

for $j = 1, 2, \dots, M$, $k = 0, 1, 2, \dots, K$, and $\ell = 1, 2, \dots, \tau_k$, where η is the learning rate and $\xi_j \subset \mathcal{S}_j$ is a randomly selected mini-batch with replacement at the j th worker. Also, $\mathbf{g}(\mathbf{w}_k^{j, \ell}; \xi_j) = \nabla F(\mathbf{w}_k^{j, \ell}; \xi_j)$ is the gradient of $F(\mathbf{w}_k^{j, \ell}; \xi_j)$ at the ℓ th local update calculated with weight vector $\mathbf{w}_k^{j, \ell}$, on the mini-batch ξ_j . The initial local weight vector $\mathbf{w}_k^{j, \ell=0}$ is the global weight vector transmitted from the server in the end of the previous round ($k - 1$). As soon as the j th worker performs τ_k number of local updates according to (2), the aggregated gradient $\mathbf{g}(\mathbf{w}_k^j)$ in the k th round is determined by

$$\mathbf{g}(\mathbf{w}_k^j) := \sum_{\ell=1}^{\tau_k} \mathbf{g}(\mathbf{w}_k^{j, \ell}; \xi_j) \quad (3)$$

for $j = 1, 2, \dots, M$, and $k = 0, 1, 2, \dots, K$. The local update coefficient τ_k determines the computation (due to local update) to communication ratio in the k th round and characterizes the communication-computation trade-off that will be discussed later.

After calculating the aggregated gradient by (3), all M workers

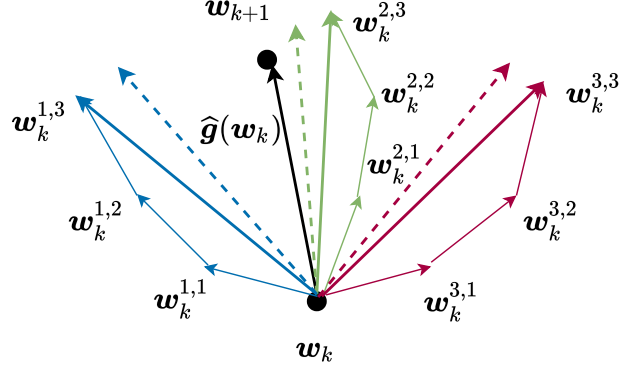


Fig. 2: In each round, workers locally update their weights multiple times based on their private datasets. Then they compress the *local aggregated gradient* and send them to the server who averages the received gradients and applies them to the global weight. It is clear that $\tau_k = 3$ and we have 3 workers.

TABLE I: Main Notations.

Notation	Description
τ_k	Local update coefficient of the k th round
s_k	Sparsity budget of the k th round
\mathbf{w}_k	Global weights of the k th round
\mathbf{w}_k^j	Weights of the j th worker (k th round)
$\mathbf{w}_k^{j, \ell}$	The ℓ th local weights of the j th worker
$\mathbf{g}(\mathbf{w}_k)$	The global gradient of the k th round
$\mathbf{g}(\mathbf{w}_k^j)$	Local aggregated gradient of the j th worker
$\mathbf{g}(\mathbf{w}_k^{j, \ell}; \xi_j)$	The ℓ th local update of the j th worker
$\hat{\mathbf{g}}(\mathbf{w}_k)$	The global compressed gradient (k th round)
$\hat{\mathbf{g}}(\mathbf{w}_k^j)$	The compressed local gradient (j th worker)

compress their locally aggregated gradients $\mathbf{g}(\mathbf{w}_k^j)$ to $\hat{\mathbf{g}}(\mathbf{w}_k^j)$ with a given *sparsity budget* of gradient compression, s_k (will be defined later), and send them to the server²³. Note that here the second trade-off, communication-precision trade-off emerges by introducing gradient compression. Then, at the server, the compressed global gradient vector, denoted by $\hat{\mathbf{g}}(\mathbf{w}_k)$, is obtained by averaging all received gradient vectors from the workers as follows:

$$\hat{\mathbf{g}}(\mathbf{w}_k) = \frac{1}{M} \sum_{j=1}^M \hat{\mathbf{g}}(\mathbf{w}_k^j), \quad k = 0, 1, 2, \dots, K. \quad (4)$$

Using the average of gradient vectors, the global weight vector

²³In the literature, two different approaches have been used for the global aggregation at the server: gradient-averaging and weight-averaging. The former requires each worker to send the computed gradients to the server, whereas in the latter weights themselves are transmitted. In this paper, we use gradient-averaging because gradients are far more sparse than weights in almost any domain. This helps to reduce the communication overhead, which is the main challenge of FL.

³If weight-averaging were adopted, the workers *would* transmit the latest weights. In our paper, however, the workers transmit the aggregate of local gradients as in (3).

\mathbf{w}_k is updated by SGD (or its variants) as follows:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \eta \hat{\mathbf{g}}(\mathbf{w}_k), \quad k = 0, 1, 2, \dots, K. \quad (5)$$

The server then broadcasts the updated global weight vector \mathbf{w}_{k+1} to synchronize all workers. This is the process of a single round involving M workers, each of which individually performing τ_k local updates with gradient compression. Fig. 1 shows the time scales of the rounds and local updates while Fig. 2 portrays local updates, aggregated gradients, compressed aggregated gradients, and the global update. Table I lists/describes our main notations. In this section, we have introduced two inherent trade-offs to be balanced for achieving fast FL, and in the next section, we will discuss the trade-offs.

B. Proposed Scheme for FL: Adjusting Two Key Variables $\{\tau_k\}$ and $\{s_k\}$

Achieving “fast FL” is the main concern of our proposed scheme. However, achieving “fast FL” and “communication-efficient FL” should not be conflated: although these two objectives overlap/correlate in practice, they are still different at a conceptual/theoretical level. Achieving fast FL requires communication-efficiency; but, it is not limited to it. For achieving fast FL, every time-consumer should be taken into account: (i) downlink communication time from the server to the workers, (ii) time for local training of neural networks (i.e., computation time) at the workers, and (iii) uplink communication time from workers to the server. In a special case, when the communication data rate is low, the communication time becomes the dominant time-consumer—and the bottleneck. In such a case, “achieving fast FL” *reduces* to “achieving communication-efficient FL.”

The problem of achieving “fast FL” can be framed as to jointly and dynamically find balances considering the two inherent trade-offs: (i) communication-computation and (ii) communication-precision. The first trade-off is characterized by the local update coefficient τ_k , for which striking a balance matters because one extreme (high communication and low computation) slows the convergence whereas the other extreme compromises the accuracy [8]; either way, the learning is decelerated. The second trade-off is characterized by the sparsity budget of gradient compression where seeking a balance is imperative since high compression/imprecision (one extreme) stalls the learning process while no compression incurs heavy/unnecessary communication overheads. Having that framing of the problem in mind, (in Section III) we will mathematically formulate our problem to minimize the learning error in a given wall-clock time with respect to local update coefficient and sparsity budget of gradient compression.

Before delving into formulating our problem, we define two key variables of the proposed scheme: the first one is the local update coefficient τ_k which was used in (3). As we have mentioned, in the proposed scheme, instead of fixing the local update coefficient τ_k over different rounds, we dynamically adjust it in the course of the entire learning process to facilitate the convergence while minimizing the dispensable gradient transmissions.

The second key variable is called the sparsity budget. Assume $\mathbf{A} \in \mathbb{R}^{m \times n}$ is a big matrix, which needs to be transmitted. To minimize the communication overheads, instead of sending

the exact matrix \mathbf{A} as is, we send its approximation $\hat{\mathbf{A}}$ that is constructed only by s_k number of selected *basis components*, which are obtained by the *compressor*. In the real-world applications, for most domains, due to sparsity, many (or even most) of the basis components of big matrices are often small and negligible [19], and in this case, $\hat{\mathbf{A}}$ can precisely approximate \mathbf{A} by judiciously selecting s_k basis components. In the context of the proposed FL, sending an approximation of matrix \mathbf{A} means sending an approximation of vector $\mathbf{g}(\mathbf{w}_k^j)$ of all users.

After all, in the proposed scheme, both variables $\{\tau_k\}$ and $\{s_k\}$ will be jointly and dynamically optimized over different rounds, while considering their interdependency. Note that adjusting the local update coefficients $\{\tau_k\}$ was studied in [8], but without considering compression of $\mathbf{g}(\mathbf{w}_k^j)$, not to mention the *dynamic* compression (characterized by $\{s_k\}$). Meanwhile, it was studied that compression could reduce the communication overheads in FL [12]–[14], [18]. However, gradient compression was not dynamically controlled; that is, $\{s_k\}$ were assumed to be fixed to a pre-determined s . Furthermore, in [12]–[14], [18], the local updates were not adopted; that is, $\{\tau_k\}$ were assumed to be all fixed to one. To the best of our knowledge, jointly adjusting $\{\tau_k\}$ and $\{s_k\}$ has not been studied in the literature, which motivated our work.

III. PROBLEM FORMULATION AND THE ERROR UPPER BOUND ANALYSIS

In this section, we first mathematically formulate our fast FL problem and then derive a mathematical expression for the error upper bound of the FL loss function in a given wall-clock time, which will be needed to develop the proposed scheme in the next section.

A. Problem Formulation

In the proposed scheme, both $\{\tau_k\}$ and $\{s_k\}$ are jointly and dynamically optimized. Mathematically, the problem is to determine the optimal values of both $\{\tau_k\}$ and $\{s_k\}$ in different rounds so that the error in a given wall-clock time is minimized. This is formulated as follows:

$$\begin{aligned} \min_{\{\tau_k\}, \{s_k\}} \mathbb{E}_{\{\xi_j\}} \left[\min_{k \in \{1, 2, \dots, K\}} F(\mathbf{w}_k; \mathcal{S}) \right] \\ \text{s.t.} \quad \sum_{k=1}^K (D_k + Y_k) = T \end{aligned} \quad (6)$$

where $\tau_k \in \{1, 2, \dots, \tau_{ub}\}$ and $s_k \in [1, s_{ub}]$ represent the local update coefficient and the sparsity budget of the k th round, respectively. Also, $F(\mathbf{w}_k; \mathcal{S})$ denotes the global loss function defined in (1) with the global weight vector \mathbf{w}_k in the k th round and dataset \mathcal{S} . Communication and computation times in the k th round are denoted by D_k and Y_k . The given wall-clock time within which the global loss function is to be minimized is denoted by T . This time constraint helps to avoid trivial solutions: (i) a solution with too large values of local update coefficients that would be time-consuming and/or (ii) a solution with its sparsity budgets set to the maximum value that extremely/unnecessarily aggravates communication overheads and consumes time. For solving the optimization problem defined in (6), we need the expression of learning error in terms of both

$\{\tau_k\}_{k=1}^K$ and $\{s_k\}_{k=1}^K$ after K rounds in a given wall-clock time T . It is generally impossible to find such an exact analytical expression [8]–[10]. It is even difficult to find an upper bound of the learning error in terms of both $\{\tau_k\}_{k=1}^K$ and $\{s_k\}_{k=1}^K$ after K rounds. In this paper, to make the analysis tractable, we will derive an upper bound of the learning error (i.e., the cost function of (6)), in terms of both τ_k and s_k after each round in a given wall-clock time. This provides the mathematical expression that will be needed to optimize the values of τ_k and s_k .

B. Error Upper Bound Analysis

In this subsection, we first derive an upper bound for the learning error after a given wall-clock time as a function of τ_k *without* considering the gradient compression that is characterized by s_k . After then, the effects of gradient compression will be reflected on the upper bound. Conjoining local update and gradient compression presents four theoretical challenges, which are addressed in our paper: (i) gradient compression must be unbiased to ensure convergence of the learning, which is the same condition that we have for SGD to be an unbiased estimation of Full-batch Gradient Descent (FGD). (ii) The variance between compressed gradient and uncompressed gradient must be made as small as possible because this way it has been shown that the learning is accelerated [18]. (iii) The communication time has to be written in terms of sparsity budget as it becomes dependent on the sparsity budget of gradient compression. (iv) Gradient compression induces a variance⁴ from uncompressed gradient (pure SGD) that needs to be incorporated into the variance of SGD. Note that SGD itself has a variance from FGD. We need to merge these two variances.

From [8], [20], we present Theorem 1 (with some adjustment of notation fitting our scheme), which holds if the following three assumptions hold: (i) the global loss function $F(\mathbf{x}; \mathcal{S})$ is differentiable, and Lipschitz smooth, which means $\|\nabla F(\mathbf{x}; \mathcal{S}) - \nabla F(\mathbf{y}; \mathcal{S})\| \leq L\|\mathbf{x} - \mathbf{y}\|$, with a lower bound F_{\inf} ; (ii) the SGD is an unbiased estimator of the FGD $\mathbb{E}_{\{\xi_j\}}[\mathbf{g}(\mathbf{w}_k)] = \nabla F(\mathbf{w}_k; \mathcal{S})$ for all k ; and (iii) the variance of the calculated global mini-batch gradient is bounded as $\mathbb{E}_{\{\xi_j\}}[\|\mathbf{g}(\mathbf{w}_k) - \nabla F(\mathbf{w}_k; \mathcal{S})\|^2] \leq \beta\|\nabla F(\mathbf{w}_k; \mathcal{S})\|^2 + \sigma$, for all k , where β and σ are non-negative constants and inversely proportional to the mini-batch size. In this inequality, σ represents the variance between the SGD, $\mathbf{g}(\mathbf{w}_k)$, and the FGD, $\nabla F(\mathbf{w}_k; \mathcal{S})$.

Theorem 1 (Error Upper Bound *without* Gradient Compression [8]): *Let Y_k and D_k denote the computation and communication times at the k th round, respectively. If the learning rate satisfies $\eta L + \eta^2 L^2 \tau_k (\tau_k - 1) \leq 1$, Y_k and D_k are constant in the k th round, and the weight vectors of all workers are initialized at the same point \mathbf{w}_k , then after T_k wall-clock time, the expression in (6) over the k th round will be bounded by:*

$$\frac{2[F(\mathbf{w}_k) - F_{\inf}]}{\eta T_k} \left(Y_k + \frac{D_k}{\tau_k} \right) + \frac{\eta L \sigma}{M} + \eta^2 L^2 \sigma (\tau_k - 1) \quad (7)$$

⁴The gradient compressor we are using in this paper is a probabilistic compressor whose mean is the same as the uncompressed gradient (unbiased). However, the gradient compressor induces a variance from the uncompressed gradient. If we did not use gradient compression before transmission, there would be no such a variance.

where L is the Lipschitz constant of the loss function.

Proof: See Appendix of [8]. \square

Theorem 1 specifies the dynamics of our first trade-off, communication-computation trade-off. As we mentioned, this trade-off is characterized by τ_k which determines the computation to communication ratio. In the error upper bound of (7), the variable τ_k is both in numerator and denominator, which indicates that both too small and large values of τ_k can contribute in increasing the error upper bound. Therefore, seeking a balance is necessary. Also, due to presence of the loss value $F(\mathbf{w}_k)$ in the error upper bound and noting that the value of loss is time-varying, the trade-off must be balanced dynamically over the course of learning.

Theorem 1 would be enough if the notion of local update were solely used *without* considering the gradient compression—communication-precision trade-off. In the proposed scheme, however, we also adopt gradient compression, which complicates the error upper bound analysis. Specifically, incorporating compression to the *aggregated local gradient* of the local updates before each transmission complicates the error upper bound given by Theorem 1 in two ways: (i) the communication time D_k becomes dependent on sparsity budget s_k . (ii) Compression introduces extra imprecisions (extra terms) to the variance σ in (7): the variance becomes dependent on the sparsity budget s_k .

As discussed in Section II-B, for compression, the matrix to be sent is written as a weighted sum of multiple *atom matrices* (basis components), where a *compressor*⁵ is used to extract the atom matrices. Due to sparsity [18], some atom matrices have more contribution than others in precisely approximating the original matrix. Therefore, the problem is to perform an unbiased selection of atom matrices so that the variance⁶ is minimized. Accordingly, we write the j th worker's gradient $\mathbf{g}(\mathbf{w}_k^j)$ as follows:

$$\mathbf{g}(\mathbf{w}_k^j) = \sum_{i=1}^B \lambda^i(\mathbf{w}_k^j) \mathbf{a}^i(\mathbf{w}_k^j) \quad (8)$$

where $\mathbf{a}^i(\mathbf{w}_k^j) \in \mathbb{R}^d$ is the i th atom, $\lambda^i(\mathbf{w}_k^j)$ is its corresponding coefficient, and B is the number of atom matrices that are summed. Note that, for notational simplicity, we assume that atoms matrices are flattened; that is why $\mathbf{a}^i(\mathbf{w}_k^j) \in \mathbb{R}^d$. Note that our formulation for compression relies on writing a matrix as a combination of atom matrices. In *compressed sensing*, this is called *atomic decomposition for sparse representation*. The formulation of compression via atomic decomposition makes our work widely inclusive of almost all of the unbiased compression techniques: specifically, TernGrad [12] and QSGD [13], two important quantization schemes, are special cases of this formulation⁷. Meanwhile, sparsification techniques such as spectral-ATOMO [18] and element-wise sparsification (e.g., top- k in DGC [17]) also comply with this formulation. Now, the problem is how and which coefficients, $\lambda^i(\mathbf{w}_k^j)$'s, should be selected. As the first requirement, we are interested in an

⁵The compressor can be any compressor that is based on the atomic decomposition for sparse representation in compressed sensing.

⁶Note that this variance is distinct from the variance (denoted by σ) mentioned earlier. This variance characterizes the variance of compression process by the estimator whereas the variance denoted by σ characterized the variance of mini-batch SGD. Later, in Theorem 3, we will merge these two variances.

⁷This has been shown in [18].

“estimator” that is unbiased. The following estimator qualifies this requirement of unbiasedness (adopted from [18] with some adaptation):

$$\hat{\mathbf{g}}(\mathbf{w}_k^j) = \sum_{i=1}^B \frac{\lambda^i(\mathbf{w}_k^j) e^i(\mathbf{w}_k^j)}{p^i(\mathbf{w}_k^j)} \mathbf{a}^i(\mathbf{w}_k^j) \quad (9)$$

where $e^i(\mathbf{w}_k^j) \sim \text{Bernoulli}(p^i(\mathbf{w}_k^j))$ and $e^i(\mathbf{w}_k^j)$'s are independent, for $0 < p^i(\mathbf{w}_k^j) \leq 1$. Also, $p^i(\mathbf{w}_k^j)$ denotes the probability characterizing the Bernoulli distribution, which will be optimized later. We derive two key properties for the estimator in (9): (i) the estimator of (9) is unbiased as it is presented in Lemma 1. Unbiasedness is necessary for guaranteeing theoretical convergence as it was a requirement of Theorem 1. (ii) The variance of the estimator in (9) is derived in Lemma 2.

Lemma 1 (Unbiased Estimator): *The estimator given in (9) is unbiased: $\mathbb{E}_{\{e^i(\mathbf{w}_k^j)\}}[\hat{\mathbf{g}}(\mathbf{w}_k^j)] = \mathbf{g}(\mathbf{w}_k^j)$.*

Proof: The proof is straightforward via the definition of expectation. \square

Lemma 2 (Variance of Estimator): *The variance of estimator in (9) is given by $\mathbb{E}_{\{e^i(\mathbf{w}_k^j)\}}[\|\hat{\mathbf{g}}(\mathbf{w}_k^j) - \mathbf{g}(\mathbf{w}_k^j)\|^2] = \sum_{i=1}^B \lambda^i(\mathbf{w}_k^j)^2 \left(\frac{1}{p^i(\mathbf{w}_k^j)} - 1 \right)$.*

Proof: See Appendix A. \square

Having derived the variance of the estimator, we can formulate an optimization problem to minimize the variance. Note that making the variance of the estimator as small as possible is important since it is known that the smaller the variance, the closer the estimated/compressed gradient $\hat{\mathbf{g}}(\mathbf{w}_k^j)$ is to the uncompressed one $\mathbf{g}(\mathbf{w}_k^j)$, which leads to faster convergence of training [18]. In the optimization problem to be formulated, the variables to be determined are $p^i(\mathbf{w}_k^j)$'s. Thus, the optimization problem is given by

$$\min \sum_{i=1}^B \frac{\lambda^i(\mathbf{w}_k^j)^2}{p^i(\mathbf{w}_k^j)} \text{ s.t. } 0 < p^i(\mathbf{w}_k^j) \leq 1 \text{ and } \sum_{i=1}^B p^i(\mathbf{w}_k^j) = s_k. \quad (10)$$

Under the assumption of s_k -balancedness in the following definition, the solution to the optimization problem in (10) is derived in Theorem 2.

Definition 1 (s_k -balancedness): *An atomic decomposition $\mathbf{g}(\mathbf{w}_k^j) = \sum_{i=1}^B \lambda^i(\mathbf{w}_k^j) \mathbf{a}^i(\mathbf{w}_k^j)$ is s_k -unbalanced at the i th entry if $\lambda^i(\mathbf{w}_k^j) s_k > \|\lambda(\mathbf{w}_k^j)\|_1$. If at no entry $\mathbf{g}(\mathbf{w}_k^j)$ is s_k -unbalanced, then we call it s_k -balanced.*

Theorem 2 (Solution to the Optimization Problem in (10)): *If $\mathbf{g}(\mathbf{w}_k^j)$ is s_k -balanced, the solution to the optimization problem in (10) is given by*

$$p^i(\mathbf{w}_k^j) = \frac{\lambda^i(\mathbf{w}_k^j) s_k}{\|\lambda(\mathbf{w}_k^j)\|_1}. \quad (11)$$

Proof: This can be proven via Lagrangian multiplier. \square

After introducing the estimator in (9) and deriving the optimal probabilities $p^i(\mathbf{w}_k^j)$'s minimizing the estimator's variance in Theorem 2, we can now establish the effects of compression on (i) communication time D_k and (ii) variance σ . For communication time D_k , assuming that $\mathbf{A} \in \mathbb{R}^{m \times n}$ is a matrix of gradients to be sent, instead of \mathbf{A} , one can send $\hat{\mathbf{A}}$ with sparsity budget s_k , i.e., only s_k number of atoms are sent. As a result, we

have communication time D_k as a function of sparsity budget s_k given by $D_k = \alpha s_k$, where α is the communication time per atom. For variance σ , based on Theorem 2, which gave the solution to the optimization problem in (10), we derive the variance of *compressed aggregated gradient* in the following theorem.

Theorem 3 (Variance of the Compressed Gradient SGD): *The variance of the compressed gradient SGD is bounded as follows:*

$$\mathbb{E}_{\{\xi_j\}, \{e^i(\mathbf{w}_k^j)\}} [\|\hat{\mathbf{g}}(\mathbf{w}_k) - \nabla F(\mathbf{w}_k)\|^2] \leq \beta \|\nabla F(\mathbf{w}_k)\|^2 + \frac{\sigma_1}{s_k} + \sigma_2 \quad (12)$$

where β , σ_1 , and σ_2 are non-negative constants and inversely proportional to the mini-batch size.

Proof: See Appendix B. \square

Now we can re-write the third assumption of Theorem 1 for the case where the gradient compression is performed, as follows: the upper bound on the variance of compressed SGD evaluated on a mini-batch from \mathcal{S}_j is given by (12). Note that, in the new assumption, σ is replaced by $\sigma_1/s_k + \sigma_2$, which does not affect the proof of Theorem 1 given in [8]. Therefore, with a change of variable $\sigma = \sigma_1/s_k + \sigma_2$, the same proof can be exploited. As a result of considering the effects of compression on communication ($D_k = \alpha s_k$) and variance ($\sigma = \sigma_1/s_k + \sigma_2$), the following theorem which is an extension of Theorem 1 can be proved.

Theorem 4 (Error Upper Bound with Gradient Compression): *If the learning rate satisfies $\eta L + \eta^2 L^2 \tau_k (\tau_k - 1) \leq 1$, Y_k and $D_k = \alpha s_k$ are constants in the k th round, and the weight vectors of all workers are initialized at the same point \mathbf{w}_k , then after T_k wall-clock time, the expression in (6) over the k th round will be bounded by:*

$$\psi_k(\tau_k, s_k) = \frac{2[F(\mathbf{w}_k) - F_{\text{inf}}]}{\eta T_k} \left(Y_k + \frac{\alpha s_k}{\tau_k} \right) + \frac{\eta L(\frac{\sigma_1}{s_k} + \sigma_2)}{M} + \eta^2 L^2 (\frac{\sigma_1}{s_k} + \sigma_2) (\tau_k - 1) \quad (13)$$

where $\psi_k(\tau_k, s_k)$ is the upper bound of the error in the k th round.

Proof: Similar to the proof of Theorem 1 [8] except that σ is replaced by $\sigma_1/s_k + \sigma_2$, and $D_k = \alpha s_k$. \square

This upper bound $\psi_k(\tau_k, s_k)$ in (13) specifies the dynamics of our *both* trade-offs: trade-offs between communication and computation/precision characterized by τ_k and s_k , respectively. Also, $\psi_k(\tau_k, s_k)$ gives us insight into how to balance these trade-offs dynamically over the course of learning. One extreme is when τ_k is set to its minimum value $\tau_k = 1$, where after a single local update (computation), the gradients are communicated. The other extreme is when τ_k is set to a very large value $\tau_k \gg 1$ (which corresponds to too many computations/local updates). The first term in (13) shows the positive effects of employing local updates on the error upper bound $\psi_k(\tau_k, s_k)$. Specifically, the first term implies that the larger the value of τ_k , the smaller the $\psi_k(\tau_k, s_k)$ (as τ_k is in the denominator). But, larger τ_k comes with a repercussion that manifests itself in the third term which shows that employing local updates causes an error as a result of a growing discrepancy among local models due to less often communication (synchronization). Therefore, the aim

is to strike a balance between these two extremes: too many communications when $\tau_k = 1$ and too many computations for $\tau_k \gg 1$.

For s_k , one extreme is to set s_k to its minimum value $s_k = 1$, where the communication overhead is at its minimum (due to high compression) and consequently the precision of the communicated gradients is the minimum. The other extreme is when $s_k \gg 1$, i.e., communicating precise gradients—high communication overheads. The first term in (13) entails reducing s_k (which is in the numerator), thereby mitigating the communication overheads, whereas the second and third terms require the value of s_k to be large (which is in the denominator) so that the introduced imprecision as a result of compression is minimized. Neither $s_k = 1$ nor $s_k \gg 1$ is an optimal choice: the former communicates too imprecise gradients resulting in a prolonged convergence, while the latter incurs too much dispensable communication overheads. The goal is to find the optimal balance between these two extremes.

IV. PROPOSED FL SCHEME: FFL

We now propose our enhanced FL scheme, namely, FFL, which jointly and dynamically adjusts the values of τ_k and s_k over different rounds in order to reduce the convergence time. Mathematically speaking, in the k th round, FFL minimizes the upper bound of the error in (13) as follows:

$$\tau_k^*, s_k^* = \arg \min_{\tau_k, s_k} \psi_k(\tau_k, s_k), \quad k = 1, 2, \dots, K \quad (14)$$

where $\tau_k \in \{1, 2, \dots, \tau_{ub}\}$ and $s_k \in [1, s_{ub}]$. We present Theorem 5 that proves the convexity of $\psi_k(\tau_k, s_k)$, and Theorem 6 that gives the optimal and approximate solutions to the problem in (14).

Theorem 5 (Convexity of $\psi_k(\tau_k, s_k)$): *If assumptions (i) $\tau_k \geq 2$, (ii) $\eta^5 \approx 0$, (iii) $(L^4 T_k \sigma_1 / 2\alpha (F(\mathbf{w}_k) - F_{\inf}) s_k^4) < \infty$, and (iv) $2\eta^2 L T_k \sigma_1 \tau_k \geq \alpha M s_k^2 (F(\mathbf{w}_k) - F_{\inf})$ hold, then $\psi_k(\tau_k, s_k)$ is convex.*

Proof: See Appendix C. \square

In the above theorem, assumption (i) excludes only one value of τ_k which is $\tau_k = 1$. Nonetheless, this does not matter because when $\tau_k = 1$, the learning has already converged—we will see this in Section V. Assumption (ii) is reasonable because $0 < \eta \ll 1$. Assumption (iii) always holds in practice. Finally, assumption (iv) excludes only a half-space of the involved parameters.

Theorem 6 (Optimal and Approximate Values of τ_k and s_k): *The value of upper error bound $\psi_k(\tau_k, s_k)$ is minimized after T_k wall-clock time when the τ_k and s_k are as follows:*

$$\tau_k = \sqrt{\frac{2\alpha [F(\mathbf{w}_k) - F_{\inf}] s_k^2}{\eta^3 L^2 (\sigma_1 + \sigma_2 s_k)}} \quad (15)$$

$$s_k = \sqrt{\frac{\sigma_1 \eta^2 L T_k (1 - \eta L (\tau_k - 1)) \tau_k}{2\alpha [F(\mathbf{w}_k) - F_{\inf}]}}. \quad (16)$$

Under assumptions (a) $\sigma_1 \ll \sigma_2 s_k$ and (b) $\eta L (\tau_k - 1) \ll 1$, the approximate values of τ_k and s_k which involve no hyperparameters are given as follows:

$$\frac{\tau_{k+1}}{\tau_k} = \sqrt{\frac{F(\mathbf{w}_{k+1}) - F_{\inf}}{F(\mathbf{w}_k) - F_{\inf}}} \sqrt{\frac{\sigma_1 + \sigma_2 s_k}{\sigma_1 + \sigma_2 s_{k+1}}} \frac{s_{k+1}}{s_k}$$

$$\approx \sqrt{\frac{F(\mathbf{w}_{k+1})}{F(\mathbf{w}_k)}} \sqrt{\frac{s_{k+1}}{s_k}} \quad (17)$$

$$\begin{aligned} \frac{s_{k+1}}{s_k} &= \sqrt{\frac{F(\mathbf{w}_k) - F_{\inf}}{F(\mathbf{w}_{k+1}) - F_{\inf}}} \sqrt{\frac{1 - \eta L (\tau_{k+1} - 1)}{1 - \eta L (\tau_k - 1)}} \sqrt{\frac{\tau_{k+1}}{\tau_k}} \\ &\approx \sqrt{\frac{F(\mathbf{w}_k)}{F(\mathbf{w}_{k+1})}} \sqrt{\frac{\tau_{k+1}}{\tau_k}}. \end{aligned} \quad (18)$$

Proof: This can be proven by setting partial derivatives of $\psi_k(\tau_k, s_k)$ to zero and applying the assumptions. \square

In Theorem 6, assumption (a) is justifiable because σ_2 is often considerably larger than σ_1 —besides, we know that $s_k \geq 1$. To explain the reason why $\sigma_2 \gg \sigma_1$, we first answer the question that where do σ_2 and σ_1 come from? When using gradient compression in addition to local update SGD for FL, there are two sources of variance with respect to the FGD: (i) calculating the gradients based on mini-batch SGD (instead of FGD) accounts for the first variance, σ_2 , and (ii) performing compression on top of SGD causes the second variance from the FGD, σ_1 . For sparse matrices (e.g., gradients)⁸, compression does not distort the matrix considerably because with few *atom* matrices, one can almost precisely represent/reconstruct the original matrix. Therefore, the compression does not have to resort to an aggressively lossy compression. Expectedly, the variance σ_1 , which is introduced as a result of compression is negligible (close to zero) compared to the variance of SGD σ_2 . Meanwhile, the value of σ_2 depends on the mini-batch size. Specifically, we know that the variance of the mini-batch SGD is inversely proportional to the mini-batch size [8], [20]. Because the mini-batch size of SGD is often significantly smaller than the size of the full-batch, the variance of mini-batch SGD σ_2 becomes large and therefore it is safe to infer that σ_2 is larger than the variance of compression σ_1 . Assumption (b) of Theorem 6 is reasonable as in practice the learning rate is usually small (near 0.01), $L < 1$, and $(\tau_k - 1) < 50$. Also, F_{\inf} is usually considered to be zero [8], [18]. Recall that in Theorem 6, assumptions (a) and (b) help to remove the dependency of our optimal solution to hyperparameters. Consequently, Equations (17) and (18) help to jointly and dynamically balance trade-offs between communication and computation/precision over the course of learning. The values of τ_k and s_k in (17) and (18) are mutually dependent whose dependence can be decoupled by substituting (17) in (18). The results are as follows:

$$\frac{\tau_{k+1}}{\tau_k} = \sqrt[3]{\frac{F(\mathbf{w}_{k+1})}{F(\mathbf{w}_k)}}, \quad \frac{s_{k+1}}{s_k} = \sqrt[3]{\frac{F(\mathbf{w}_k)}{F(\mathbf{w}_{k+1})}}. \quad (19)$$

Also, we can write them in terms of initial values $F(\mathbf{w}_0)$, τ_0 , and s_0 as follows:

$$\tau_k = \sqrt[3]{\frac{F(\mathbf{w}_k)}{F(\mathbf{w}_0)}} \tau_0, \quad s_k = \sqrt[3]{\frac{F(\mathbf{w}_0)}{F(\mathbf{w}_k)}} s_0. \quad (20)$$

Equation (20) yields conclusive results about the profiles of τ_k and s_k over the course of learning. It implies that as the learning proceeds and the loss value $F(\mathbf{w}_k)$ gets smaller, the value of τ_k should also decrease, while the value of s_k needs to increase.

⁸The sparsity of gradients (our supposition) has been observed and exploited in many studies [17]–[19].

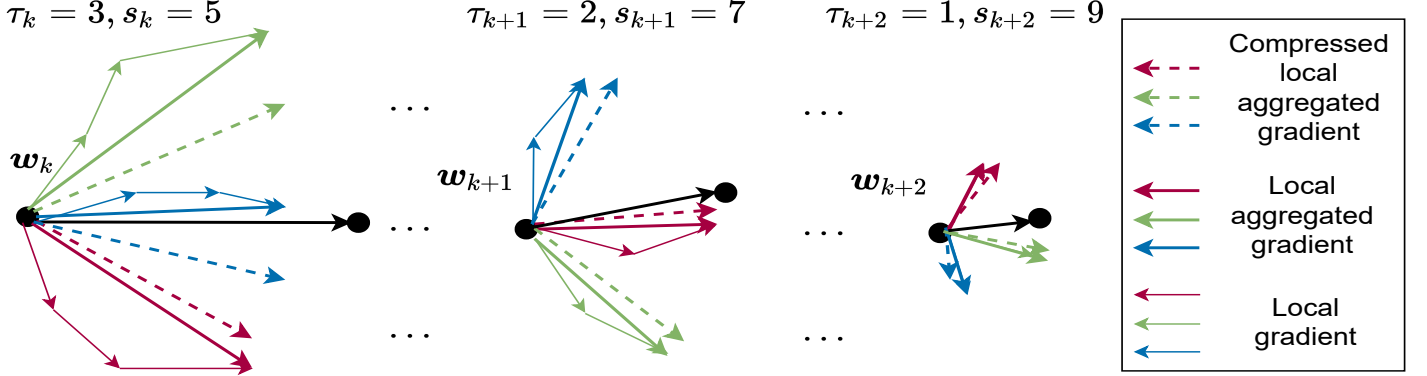


Fig. 3: Over the course of learning, the value of τ_k decreases whereas the value of s_k increases.

We view both the communication-computation trade-off and the communication-precision trade-off as two distinct exploration-exploitation trade-offs, respectively—see Fig. 3. This view/perspective provides the rationale behind the joint and dynamic adjustment of τ_k and s_k ; it also immensely assists to elucidate our proposed scheme FFL. Our two distinct exploration-exploitation trade-offs are explained as follows: at the beginning of learning, when the loss value is large, higher τ_k can accelerate learning (optimization) via aggressive exploration. However, once the loss becomes small, a large τ_k would fail to exploit. For s_k , meanwhile, we start with low values and over time increase it because when the loss is large, even imprecise gradients can contribute to decreasing the loss effectively—facilitate exploration. At lower losses, however, more precise gradients should be used—for exploitation. Fig. 3 portrays the exploration-exploitation trade-offs: as time passes, for k th, $(k+1)$ th, and $(k+2)$ th rounds, the values of τ_k , τ_{k+1} , and τ_{k+2} decrease from 3 and 2 to 1. Meantime, for s_k , over time, the values of s_k , s_{k+1} , and s_{k+2} increase from 5 and 7 to 9. It can be seen that over time the vector of *compressed local aggregated gradient* gets closer to *local aggregated gradient* because the compression is becoming more precise—more exploitation.

Algorithm 1 presents FFL at the server and workers. Interestingly, FFL is similar to Model-Agnostic Meta-Learning (MAML) [21] in three aspects: (i) in FFL, we have *workers* who possess their own datasets (with different distributions) while in MAML, there are *tasks* with their corresponding datasets. (ii) In FFL, we have *local updates* that improve performances (loss/accuracy) of individual workers whereas in MAML, there are *inner-loop updates* that do so for individual tasks. (iii) In FFL, in each round, we apply *global gradient* to weights, which improves the overall performance of all workers while in MAML, the *outer-loop update* does so for all of tasks. Meanwhile, FFL is dissimilar to MAML in three facets: (i) FFL compresses *local aggregated gradients* while MAML does not. (ii) FFL dynamically adjusts the local updates whereas MAML keeps the number of inner-loop updates fixed. (iii) In FFL, we only have *one* update with global gradient while MAML can have multiple outer-loop updates. This coincidental similarity between FFL and MAML is promising as MAML has shown considerable success in *accelerating* adaptation/generalization to different environments/tasks and the objective of FFL is achieving *fast* FL by balancing trade-offs between communication and computation/precision.

Algorithm 1: FFL at the server and workers.

```

1 The server broadcasts  $w_0$ ;
2 Workers receive and initialize  $w_0$ ;
3 for  $k = 1, \dots, K$  do
4   The server calculates  $\tau_k$  and  $s_k$  using (20);
5   The server broadcasts  $\tau_k$  and  $s_k$ ;
6   for  $j = 1, \dots, M$  in parallel do
7     Workers receive  $\tau_k$  and  $s_k$ ;
8     for  $\ell = 1, 2, \dots, \tau_k$  do
9       Workers compute  $g(w_k^{j,\ell}; \xi_j)$ ;
10      Workers update  $w_k^{j,\ell}$  as in (2);
11    end
12    Workers compute  $g(w_k^j)$  as in (3);
13    Workers compress  $g(w_k^j)$  as in (9) and (10);
14    Workers transmit  $\hat{g}(w_k^j)$  to the server;
15  end
16  The server receives  $\hat{g}(w_k^j)$ 's from workers;
17  The server averages  $\hat{g}(w_k^j)$ 's as in (4);
18  The server updates the global model as in (5);
19  The server broadcast  $w_{k+1}$  to workers;
20 end

```

In practical applications, our proposed scheme (Algorithm 1) is implemented in the following steps: (i) each worker sends a participation signal to the server that serves as a participation request/permission. (ii) When the server, via referring to its database, approves (the reliability/trustworthiness of) the worker, it sends the latest weights as well as the local update coefficient and sparsity budget to the worker. The local update coefficient determines the number of gradient computations before communication, while the sparsity budget defines the rate of compression at which the worker is supposed to compress the gradients prior to communication. This way, the worker is synchronized with the FL system. (iii) The workers then, in each round, continuously compute gradients as many times as the local update coefficient allows and those workers transmit the compressed gradients at the compression rate that is specified by the sparsity budget. (iv) The server, meanwhile, receives the compressed gradients from the workers and computes the average gradient with which it updates the global weights. Later, the server jointly optimizes the local update coefficient and sparsity budget (via Theorem

6) for the upcoming round based on the latest loss value. The updated weights as well as the optimal local update coefficient and sparsity budget are then broadcast to workers.

V. EXPERIMENTS

We used the TensorFlow 2.x deep learning library to conduct the experiments for performance evaluation of the proposed scheme. Specifically, we compare FFL with the following two state-of-the-art schemes: (i) ADACOMM [8], which exploits the idea of dynamically adjusting τ_k , yet the workers do not compress the gradients for transmission and (ii) ATOMO [18], which makes use of gradient compression, but without dynamically adjusting s_k over different rounds nor benefiting from dynamically adjusting τ_k . We examine aforementioned schemes on two learning tasks: a Fully connected Neural Network (FNN) with the architecture of [784, 400, 400, 10] on MNIST, and a Convolutional Neural Network (CNN) that is VGG16 on CIFAR10. In all simulation scenarios, learning rate is set to 0.01 and SGD optimizer is used—a simple optimizer lest simulation results get obscured [22]. For compression, as in [18], we adopt Singular Value Decomposition (SVD). However, for implementation efficiency, we use thin-SVD [23] that calculates only s number of singular values of $\mathbf{A} \in \mathbb{R}^{m \times n}$ with time-complexity of $O(mns)$ for $s \leq \sqrt{\min(m, n)}$ as opposed to full-SVD with $O(mn \times \min(m, n))$. Empirical observations suggest that most of the singular values of gradient matrices are often (very) small, i.e., (very) close to zero [17]–[19]; hence, calculating a small number of singular values for gradient matrices, as it is done in thin-SVD, often suffices. Note that neither our analysis nor the proposed scheme is not limited to SVD, which is used only as an example of compression in the experiments.

The additional computation cost introduced by the SVD procedure reflects itself in the time axis of Figs. 4, 5, 6, 7, 8, and 9. Specifically, our scheme FFL uses the SVD compression while ADACOMM, which is one of our baselines, does not—ATOMO which is the other baseline does. This additional computation time of FFL that is spent for SVD compression causes computation part of each round of FFL to take slightly longer time compared to ADACOMM. Nonetheless, the SVD compression substantially reduces communication overheads, thereby saving considerable amount of time in compensation. Therefore, the overall time taken for a round is shorter in our scheme. This trade-off, which is spending time for computation in order to compress the gradients, is significantly time-saving because usually it is the communication time that is the bottleneck. Also, the amount of time spent for computation required for SVD compression is negligible compared to the gradient computation that includes performing τ number of local updates.

We plot variations of desired parameters (e.g., accuracy) against wall-clock time instead of FL rounds, epochs, or iterations. This is because unlike traditional centralized learning where different iterations take almost the same time and therefore the number of iterations reflects the speed of convergence, for FL that does not hold. The convergence speed of FL is determined by two factors: (i) the number of rounds and (ii) the time spent for each round that can *significantly* differ over time

because each round comprises (a) a downlink broadcast of global weights from the server to the workers with different link speeds (different delays), (b) local gradient computations that are again dependent on many factors such as the number of local updates and hardware heterogeneity, and finally, (c) transmissions of the compressed gradients to the server by workers whose delays depend on compression ratios and uplink speeds. Hence, simply reporting the number of rounds is neither informative nor conclusive for the speed of convergence.

In experiments, for each worker, we first *measure* the time it consumes for performing local computation and communication of gradients in each round. Then we select the maximum time consumed corresponding to the slowest worker as the time consumption of the round, because it is the slowest worker (the straggler) in each round that determines the time consumption of that round. We used the *time()* function of *time* package available in Python programming language for measuring the time spent for training process (local computation) and communication of various workers.

We explore the performance results of the three schemes along *seven* distinct dimensions which characterize our learning environment: (i) Neural Network (NN) model, (ii) dataset, (iii) number of workers, (iv) uplink/downlink data rates, (v) momentum at workers, (vi) non-IID-ness, and (vii) noisy channels with packet failure. Throughout simulations, the data distribution among workers is balanced and non-overlapping. We adopt reasonable mini-batch sizes of 64 and 128 for MNIST and CIFAR10, respectively [22]. In FL, it is crucial to adopt reasonable mini-batch sizes: because on one extreme, a too small mini-batch size cripples/prolongs the learning process due to high-variance SGD gradients—demanding even more communication overheads/rounds which in turn further slow the convergence [18], [22]. Recall that the variance of SGD gradients (denoted by σ_2 in this paper) is inversely proportional to the mini-batch size, and therefore the smaller the mini-batch size, the more high-variance/imprecise the gradients are [8], [20]⁹. However, we should not overshoot: the size of the mini-batch must be just as large as to ensure that the mini-batch is fairly representative. As soon as representative-ness of the mini-batch is achieved, increasing the mini-batch size further only consumes dispensable resources—the other extreme. Because a larger mini-batch size proportionately demands larger memory and more computation which might be unaffordable particularly for resource-constrained IoT nodes and edge-devices [22]. Interestingly, determining the mini-batch size itself involves balancing an inherent communication-computation trade-off: a too small mini-batch size (which translates to too little computation per communication) prolongs/cripples the learning, whereas an unnecessarily large mini-batch size (too much computation per communication) slows the learning.

In our experiments, the supported instantaneous communication data rate (more precisely, the instantaneous channel capacity)

⁹The main cause of high-variance gradients is having a non-representative mini-batch: a too small mini-batch size or a poorly shuffled mini-batch. Having a representative mini-batch to reduce the variance of gradients is necessary for acceleration of convergence [18] both in general, and particularly for FL because in FL the learning process is distributed; this makes it more prone to gradient conflict/discrepancy and poor convergence. To secure representative-ness, a large (well-shuffled) mini-batch is vital.

is time-invariant. Specifically, in the t -th time slot, the supported instantaneous communication data rate R_t (bps) is given by

$$R_t = W \log_2 \left(1 + \frac{|h_t|^2 P_t}{\sigma^2} \right) \quad (21)$$

where W is the assigned bandwidth (Hz), P_t is the transmission power (watts) in the t -th time slot, h_t is the channel coefficient in the t -th time slot, and σ^2 is the noise power (watts). In our experiments, we assume that the data rate R_t is fixed throughout the course of learning by employing power control. For all data rates R_t adopted in our experiments such as 10Kbps, 100Kbps, and 10Mbps, the bandwidth W is constant, equal to 1MHz in all cases, channel coefficient h_t varies among workers, and noise power is also kept constant at $1mW$. The transmission power P_t is controlled so that the received signal to noise ratio, $|h_t|^2 P_t / \sigma^2$ yields the required data rate R_t .

Experiment results of the FNN on MNIST are shown in Fig. 4 for 32 workers. Figs. 4a, 4b, and 4c show the values of accuracy, τ_k , and s_k , respectively. Fig. 4a shows that, compared to the other two schemes, FFL achieves higher accuracies faster in terms of wall-clock time. In Fig. 4b, since ATOMO does not use local updates, it is set to one. The values of s_k are plotted in Fig. 4c except for ADACOMM because it does not use gradient compression, whereas ATOMO does, but at a fixed s_k . Unlike these two schemes, FFL dynamically adjusts the value of τ_k and s_k over time: for τ_k , FFL starts with higher values to enable workers to perform an aggressive exploration. However, once the accuracy rises, the exploitation begins with lower τ_k . For s_k , FFL starts with low values and over time increases because when the accuracy is small, even coarse gradients (i.e., highly compressed gradients due to low s_k) can still contribute to improving the accuracy of the model effectively. At higher accuracies, however, finer gradients are used for exploitation. Although all schemes in the figure perform 1560 rounds, ADACOMM takes twice as much time as others to converge because FFL and ATOMO are compressing their gradients in the uplink at a compression ratio of about 1/50 associated with $s_k \in [5, 9]$ that results in making the uplink communication time negligible compared to the downlink. Results of Fig. 4 carry over to scenarios with different number of workers such as 4, 8, 16, and 64; we, therefore, avoid discussing those results again. In the rest of scenarios, hyperparameters are set the same as in Fig. 4; otherwise, it is mentioned.

In Fig. 5, we investigate three learning scenarios whose uplink and downlink data rates are different. First, in Fig. 5a, we consider high uplink/downlink communication data rates that equal to 10 Mbps—a relatively high data rate so as to dwarf the importance of compression. In such a case, not only compression is not that important in reducing the overall delay but also it can be detrimental due to its introduction of a compression time/delay. Indeed, here computation has become more of a bottleneck than communication. Expectedly, in Fig. 5a, ADACOMM is outperforming ATOMO despite its employment of compression. In this case, FFL approaches the performance of ADACOMM, even though FFL benefits from compression. Since employing compression is no longer the primary contributor in reducing the convergence time; what here plays the role instead is whether a scheme is using local updates or not. However, FFL is still outperforming ADACOMM

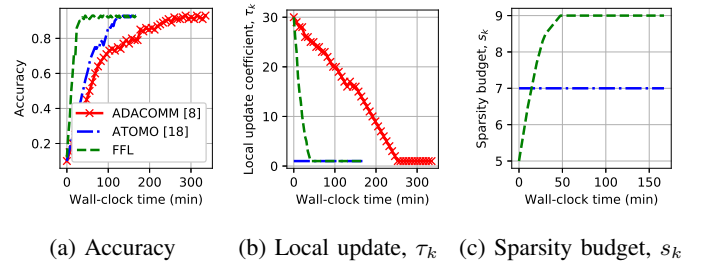


Fig. 4: Accuracy, τ_k , and s_k values over wall-clock time for ADACOMM [8], ATOMO [18], and our proposed FFL with 32 workers for the FNN trained on MNIST. The dataset distribution is IID and all workers participate without momentum, while the server uses a momentum of 0.9. The uplink and downlink data rates are set equal to 100 Kbps. Note that $\tau_k \in \{1, \dots, 30\}$ and $s_k \in [5, 9]$.

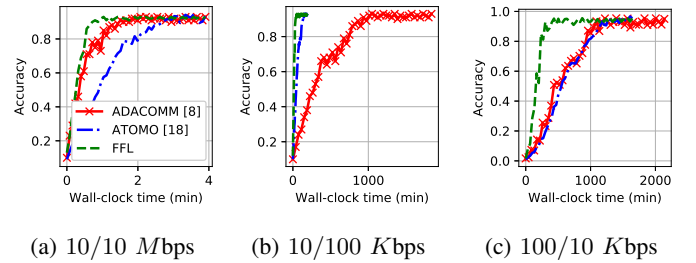


Fig. 5: Accuracy for different uplink/downlink data rates.

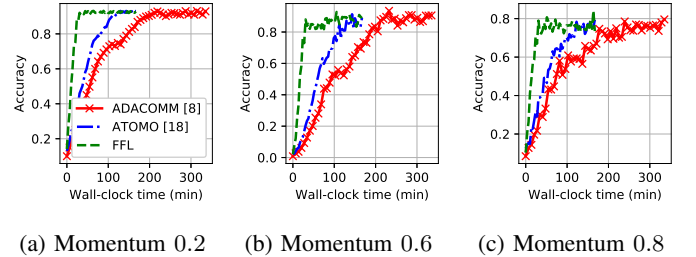
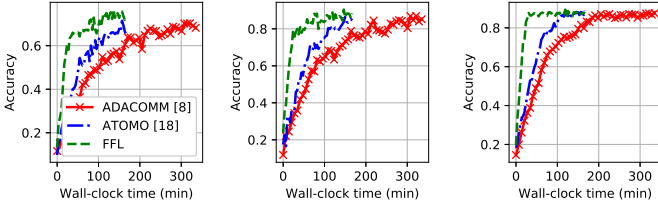


Fig. 6: Accuracy for different momentums at workers.

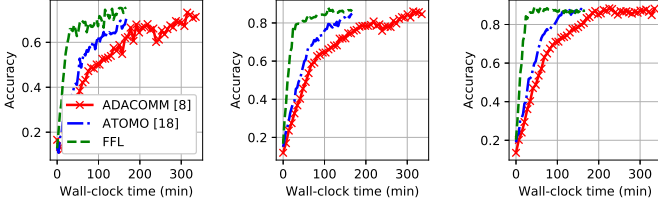
due to a slight benefit it is taking from compression. Albeit, when uplink/downlink data rates approach infinity, ADACOMM would grow closer to FFL because the communication time of each round would approach zero and then it would not matter whether we compress gradients or not. In Fig. 5b, realistic values of uplink/downlink data rates (10/100 Kbps) are chosen; here uplink communication is the bottleneck. In this case, the schemes which compress the gradients in the uplink win the race. Also, it can be seen that FFL has a slight advantage over ATOMO thanks to local updates. Conversely, in Fig. 5c, the downlink communication is the bottleneck which renders uplink compression futile. Consequently, ATOMO, despite its compression in the uplink, exhibits the same performance as ADACOMM; nonetheless, FFL outperforms both.

In Fig. 6, the impacts of using momentum (block momentum) at the workers are studied: we assess three values for momentum. In Fig. 6a, we notice that a small value for momentum at workers accelerates the learning compared to Fig. 4a where no momentum was used. Nevertheless, for higher values of



(a) 3 classes (b) 6 classes (c) 9 classes

Fig. 7: Accuracy for different levels of non-IID-ness.



(a) Probability of 0.7 (b) Probability of 0.4 (c) Probability of 0.1

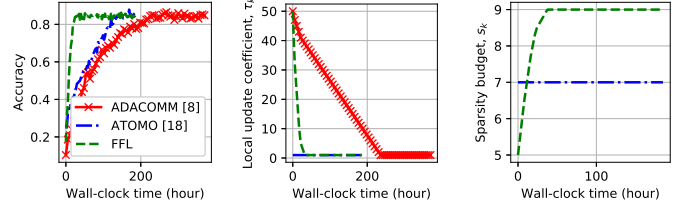
Fig. 8: Accuracy for different probabilities of packet failure.

momentum, the learning process starts to destabilize and yields lower accuracies, near 89% and 73% in Figs. 6b and 6c. This is because higher momentums at workers contribute in increasing the discrepancy among local gradients particularly because the dataset distribution is non-overlapping, which aggravates *gradient conflict* among workers. The same behavior was observed in [22] where high momentum crippled the learning. Yet, FFL is still more tolerant of higher momentums at workers than others.

Investigating non-IID data distribution over workers is demonstrated in Fig. 7 where to each worker is assigned only 3, 6, and 9 classes in Figs. 7a, 7b, and 7c, respectively. Assigning 9 classes to workers (close to the IID scenario) does not inflict any noticeable harm, while for 6 classes, the baselines exhibit minor deterioration in performance; nonetheless, FFL is still on top. Finally, for 3 classes the accuracies undergo a considerable drop to 60%, which is unsurprising since non-IID learning is still an open problem even in centralized machine learning, let alone FL [24].

Performances for noisy channels with packet failure are examined in Fig. 8 where at each global update only a percentage of the entire gradients reach the server; however, it is assumed that workers are still synchronized with the latest weights. The challenge of noisy channels with packet failure is that it can stall learning when only a set of *non-representative gradients* reach the server and steer the learning away from the minimum and might even cause *catastrophic forgetting* by directing the learning in the opposite direction. Results for noisy channels with packet failure probability of 0.1, 0.4, and 0.7 are shown in Figs. 8a, 8b, and 8c.

Experiment results of the VGG16 on CIFAR10 are shown in Fig. 9 for 16 workers. In this figure, compared to the previous results, FFL even more significantly outperforms the other two schemes. Numerically, whereas in the MNIST experiment in Fig. 4, FFL was 4 \times faster than ATOMO, here in CIFAR10 experiment, FFL is 11 \times faster, because either of two techniques,



(a) Accuracy (b) Local update, τ_k (c) Sparsity budget, s_k

Fig. 9: Accuracy, τ_k , and s_k for VGG16 on CIFAR10.

local updates and gradient compression, becomes more necessary for larger NN models. Specifically, the larger the NN model, the more significant the role of dynamic local updates, which can accelerate the learning/optimization by dynamic adjustment of local update coefficients corresponding to the communication-computation trade-off. Meanwhile, when the NN model is small, the gradient constituents are also small matrices, which usually are not as low-rank as large matrices and therefore are not as suitable for compression, whereas large matrices tend to be highly low-rank [17]–[19]. Results of Fig. 9 carry over to different number of workers such as 4 and 8. The other results observed for the FNN on MNIST such as effects of uplink/downlink data rates, momentum, non-IID-ness, and noisy channels with packet failure carry over to VGG16 on CIFAR10; hence, we avoid presenting them.

A summary of main results from the experiments are as follows: (i) as NN models scale, either of local update and gradient compression becomes more effective in accelerating learning, (ii) both dynamic adjustments of local update and gradient compression can be thought of as exploration-exploitation trade-offs: we start with high exploration (large τ_k and small s_k) and end with pure exploitation ($\tau_k = 1$ and large s_k), (iii) FFL provides formulations to balance these trade-offs, and (iv) high momentums at *workers*, non-IID-ness, and noisy channels with packet failure, which are the main causes of non-representative gradients inflict the least harm on FFL.

VI. CONCLUSION

In this paper, our goal was to accelerate FL via minimizing learning error in a given wall-clock time with respect to local updates and gradient compression that correspond to trade-offs between communication and computation/precision, respectively. To this end, we first derived an upper bound of the learning error in a given wall-clock time considering the interdependency between the two variables: local update coefficient and sparsity budget. Based on this theoretical analysis, we then proposed an enhanced FL scheme, namely Fast FL (FFL), which jointly and dynamically adjusted the two variables to achieve *fast* FL. In experiments, we demonstrated that FFL consistently achieved higher accuracies faster than similar schemes in the literature.

REFERENCES

- [1] E. P. Xing, Q. Ho, W. Dai, J. K. Kim, J. Wei, S. Lee, X. Zheng, P. Xie, A. Kumar, and Y. Yu, “Petuum: A new platform for distributed machine learning on big data,” *IEEE Trans. Big Data*, vol. 1, no. 2, pp. 49–67, Jun. 2015.

APPENDIX A
PROOF OF LEMMA 2

- [2] P. Goyal, P. Dollar, R. Girshick, P. Noordhuis, L. Wesolowski, A. Kyrola, A. Tulloch, Y. Jia, and K. He, "Accurate, large minibatch SGD: Training ImageNet in 1 hour," *arXiv:1706.02677*, 2017.
- [3] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Proc. Artificial Intelligence and Statistics*, 2017, pp. 1273–1282.
- [4] J. Konecny, H. B. McMahan, F. X. Yu, P. Richtarik, A. T. Suresh, and D. Bacon, "Federated learning: Strategies for improving communication efficiency," *arXiv:1610.05492*, 2016.
- [5] X. Wang, Y. Han, V. C. Leung, D. Niyato, X. Yan, and X. Chen, "Convergence of edge computing and deep learning: A comprehensive survey," *IEEE Communications Surveys & Tutorials*, vol. 22, no. 2, pp. 869–904, Jan. 2020.
- [6] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith, "Federated learning: Challenges, methods, and future directions," *arXiv:1908.07873*, 2019.
- [7] T. Lin, S. U. Stich, K. K. Patel, and M. Jaggi, "Don't use large mini-batches, use local SGD," *arXiv:1808.07217*, 2018.
- [8] J. Wang and G. Joshi, "Adaptive communication strategies to achieve the best error-runtime trade-off in local-update SGD," in *Proc. Conference on Machine Learning and Systems*, vol. 1, 2019, pp. 126–145.
- [9] F. Zhou and G. Cong, "On the convergence properties of a k -step averaging stochastic gradient descent algorithm for nonconvex optimization," *arXiv:1708.01012*, 2017.
- [10] S. Wang, T. Tuor, T. Salonidis, K. K. Leung, C. Makaya, T. He, and K. Chan, "Adaptive federated learning in resource constrained edge computing systems," *IEEE J. Sel. Areas Comm.*, vol. 37, no. 6, pp. 1205–1221, Mar. 2019.
- [11] F. Sattler, S. Wiedemann, K.-R. Müller, and W. Samek, "Sparse binary compression: Towards distributed deep learning with minimal communication," in *Proc. International Joint Conference on Neural Networks (IJCNN)*, 2019, pp. 1–8.
- [12] W. Wen, C. Xu, F. Yan, C. Wu, Y. Wang, Y. Chen, and H. Li, "TernGrad: Ternary gradients to reduce communication in distributed deep learning," in *Proc. Advances in Neural Information Processing Systems*, 2017, pp. 1509–1519.
- [13] D. Alistarh, D. Grubic, J. Li, R. Tomioka, and M. Vojnovic, "QSGD: Communication-efficient SGD via gradient quantization and encoding," in *Proc. Advances in Neural Information Processing Systems*, 2017, pp. 1709–1720.
- [14] F. Seide, H. Fu, J. Droppo, G. Li, and D. Yu, "1-bit stochastic gradient descent and its application to data-parallel distributed training of speech DNNs," in *Proc. International Speech Communication Association*, 2014, pp. 1–8.
- [15] J. Wu, W. Huang, J. Huang, and T. Zhang, "Error compensated quantized SGD and its applications to large-scale distributed optimization," *arXiv:1806.08054*, 2018.
- [16] P. Jiang and G. Agrawal, "A linear speedup analysis of distributed deep learning with sparse and quantized communication," in *Proc. Advances in Neural Information Processing Systems*, 2018, pp. 2525–2536.
- [17] Y. Lin, S. Han, H. Mao, Y. Wang, and B. Dally, "Deep gradient compression: Reducing the communication bandwidth for distributed training," in *Proc. International Conference on Learning Representations*, 2018, pp. 9850–9861.
- [18] H. Wang, S. Sievert, S. Liu, Z. Charles, D. Papailiopoulos, and S. Wright, "ATOMO: Communication-efficient learning via atomic sparsification," in *Proc. Advances in Neural Information Processing Systems*, 2018, pp. 9850–9861.
- [19] M. Udell and A. Townsend, "Why are big data matrices approximately low rank?" *SIAM Journal on Mathematics of Data Science*, vol. 1, no. 1, pp. 144–160, Feb. 2019.
- [20] J. Wang and G. Joshi, "Cooperative SGD: A unified framework for the design and analysis of communication-efficient SGD algorithms," *arXiv:1808.07576*, 2018.
- [21] C. Finn, P. Abbeel, and S. Levine, "Model-agnostic meta-learning for fast adaptation of deep networks," in *Proc. International Conference on Machine Learning*, 2017, pp. 1126–1135.
- [22] F. Sattler, S. Wiedemann, K. R. Müller, and W. Samek, "Robust and communication-efficient federated learning from non-i.i.d. data," *IEEE Trans. Neural Networks and Learning Systems*, vol. 31, no. 9, pp. 3400–3413, Sep. 2020.
- [23] M. Brand, "Fast low-rank modifications of the thin singular value decomposition," *Linear Algebra and Its Applications*, vol. 415, no. 1, pp. 20–30, May 2006.
- [24] S. P. Karimireddy, S. Kale, M. Mohri, S. Reddi, S. Stich, and A. T. Suresh, "Scaffold: Stochastic controlled averaging for federated learning," in *Proc. International Conference on Machine Learning*, 2020, pp. 5132–5143.

We determine the variance of the estimator in (9) starting from $\mathbb{E}_{\{e^i(\mathbf{w}_k^j)\}} \left[\|\hat{\mathbf{g}}(\mathbf{w}_k^j) - \mathbf{g}(\mathbf{w}_k^j)\|^2 \right]$ as follows:

$$\begin{aligned} & \mathbb{E}_{\{e^i(\mathbf{w}_k^j)\}} \left[\|\hat{\mathbf{g}}(\mathbf{w}_k^j) - \mathbf{g}(\mathbf{w}_k^j)\|^2 \right] \\ &= \mathbb{E}_{\{e^i(\mathbf{w}_k^j)\}} \left[\left\| \sum_{i=1}^B \left(\frac{\lambda^i(\mathbf{w}_k^j) e^i(\mathbf{w}_k^j)}{p^i(\mathbf{w}_k^j)} \right) \mathbf{a}^i(\mathbf{w}_k^j) - \lambda^i(\mathbf{w}_k^j) \mathbf{a}^i(\mathbf{w}_k^j) \right\|^2 \right] \end{aligned} \quad (\text{A.1})$$

$$\begin{aligned} &= \mathbb{E}_{\{e^i(\mathbf{w}_k^j)\}} \left[\left(\sum_{i=1}^B \lambda^i(\mathbf{w}_k^j) \mathbf{a}^i(\mathbf{w}_k^j) \left(\frac{e^i(\mathbf{w}_k^j) - p^i(\mathbf{w}_k^j)}{p^i(\mathbf{w}_k^j)} \right) \right)^T \times \left(\sum_{i=1}^B \lambda^i(\mathbf{w}_k^j) \mathbf{a}^i(\mathbf{w}_k^j) \left(\frac{e^i(\mathbf{w}_k^j) - p^i(\mathbf{w}_k^j)}{p^i(\mathbf{w}_k^j)} \right) \right) \right] \end{aligned} \quad (\text{A.2})$$

$$\begin{aligned} &= \sum_{i=1}^B \lambda^i(\mathbf{w}_k^j)^2 \|\mathbf{a}^i(\mathbf{w}_k^j)\|^2 \\ &\quad \times \mathbb{E}_{\{e^i(\mathbf{w}_k^j)\}} \left[\left(\frac{e^i(\mathbf{w}_k^j) - p^i(\mathbf{w}_k^j)}{p^i(\mathbf{w}_k^j)} \right)^2 \right] \\ &\quad + \sum_{r,t;r \neq t}^B \lambda^r(\mathbf{w}_k^j) \lambda^t(\mathbf{w}_k^j) \langle \mathbf{a}^r(\mathbf{w}_k^j), \mathbf{a}^t(\mathbf{w}_k^j) \rangle \\ &\quad \times \mathbb{E}_{\{e^i(\mathbf{w}_k^j)\}} \left[\left(\frac{e^r(\mathbf{w}_k^j) - p^r(\mathbf{w}_k^j)}{p^r(\mathbf{w}_k^j)} \right) \left(\frac{e^t(\mathbf{w}_k^j) - p^t(\mathbf{w}_k^j)}{p^t(\mathbf{w}_k^j)} \right) \right]. \end{aligned} \quad (\text{A.3})$$

Now, we need to determine $\mathbb{E}_{\{e^i(\mathbf{w}_k^j)\}} \left[\left(\frac{e^i(\mathbf{w}_k^j) - p^i(\mathbf{w}_k^j)}{p^i(\mathbf{w}_k^j)} \right)^2 \right]$ and $\mathbb{E}_{\{e^i(\mathbf{w}_k^j)\}} \left[\left(\frac{e^r(\mathbf{w}_k^j) - p^r(\mathbf{w}_k^j)}{p^r(\mathbf{w}_k^j)} \right) \left(\frac{e^t(\mathbf{w}_k^j) - p^t(\mathbf{w}_k^j)}{p^t(\mathbf{w}_k^j)} \right) \right]$. Starting with the first one, we have:

$$\begin{aligned} & \mathbb{E}_{\{e^i(\mathbf{w}_k^j)\}} \left[\left(\frac{e^i(\mathbf{w}_k^j) - p^i(\mathbf{w}_k^j)}{p^i(\mathbf{w}_k^j)} \right)^2 \right] \\ &= (1 - p^i(\mathbf{w}_k^j)) \times \left(\frac{0 - p^i(\mathbf{w}_k^j)}{p^i(\mathbf{w}_k^j)} \right)^2 \\ &\quad + p^i(\mathbf{w}_k^j) \times \left(\frac{1 - p^i(\mathbf{w}_k^j)}{p^i(\mathbf{w}_k^j)} \right)^2 \end{aligned} \quad (\text{A.4})$$

$$= \left(\frac{1}{p^i(\mathbf{w}_k^j)} - 1 \right). \quad (\text{A.5})$$

The second one can be written as follows:

$$\mathbb{E}_{\{e^i(\mathbf{w}_k^j)\}} \left[\left(\frac{e^r(\mathbf{w}_k^j) - p^r(\mathbf{w}_k^j)}{p^r(\mathbf{w}_k^j)} \right) \times \left(\frac{e^t(\mathbf{w}_k^j) - p^t(\mathbf{w}_k^j)}{p^t(\mathbf{w}_k^j)} \right) \right] \quad (\text{A.6})$$

$$\begin{aligned}
&= \mathbb{E}_{\{e^i(\mathbf{w}_k^j)\}} \left[\left(\frac{e^r(\mathbf{w}_k^j) - p^r(\mathbf{w}_k^j)}{p^r(\mathbf{w}_k^j)} \right) \right] \\
&\quad \times \mathbb{E}_{\{e^i(\mathbf{w}_k^j)\}} \left[\left(\frac{e^t(\mathbf{w}_k^j) - p^t(\mathbf{w}_k^j)}{p^t(\mathbf{w}_k^j)} \right) \right]. \quad (\text{A.7})
\end{aligned}$$

We show that $\mathbb{E}_{\{e^i(\mathbf{w}_k^j)\}} \left[\left(\frac{e^r(\mathbf{w}_k^j) - p^r(\mathbf{w}_k^j)}{p^r(\mathbf{w}_k^j)} \right) \right]$ and $\mathbb{E}_{\{e^i(\mathbf{w}_k^j)\}} \left[\left(\frac{e^t(\mathbf{w}_k^j) - p^t(\mathbf{w}_k^j)}{p^t(\mathbf{w}_k^j)} \right) \right]$ are zero.

$$\begin{aligned}
&\mathbb{E}_{\{e^i(\mathbf{w}_k^j)\}} \left[\left(\frac{e^r(\mathbf{w}_k^j) - p^r(\mathbf{w}_k^j)}{p^r(\mathbf{w}_k^j)} \right) \right] \\
&= (1 - p^r(\mathbf{w}_k^j)) \left(\frac{0 - p^r(\mathbf{w}_k^j)}{p^r(\mathbf{w}_k^j)} \right) + p^r(\mathbf{w}_k^j) \left(\frac{1 - p^r(\mathbf{w}_k^j)}{p^r(\mathbf{w}_k^j)} \right) \\
&= 0. \quad (\text{A.8})
\end{aligned}$$

In the same way, the second term is zero. Resuming with (A.3), we have:

$$\begin{aligned}
&\mathbb{E}_{\{e^i(\mathbf{w}_k^j)\}} \left[\|\widehat{\mathbf{g}}(\mathbf{w}_k^j) - \mathbf{g}(\mathbf{w}_k^j)\|^2 \right] = \sum_{i=1}^B \lambda^i(\mathbf{w}_k^j)^2 \|\mathbf{a}^i(\mathbf{w}_k^j)\|^2 \\
&\quad \times \mathbb{E}_{\{e^i(\mathbf{w}_k^j)\}} \left[\left(\frac{e^i(\mathbf{w}_k^j) - p^i(\mathbf{w}_k^j)}{p^i(\mathbf{w}_k^j)} \right)^2 \right] \\
&\quad + \sum_{r,t;r \neq t}^B \lambda^r(\mathbf{w}_k^j) \lambda^t(\mathbf{w}_k^j) \langle \mathbf{a}^r(\mathbf{w}_k^j), \mathbf{a}^t(\mathbf{w}_k^j) \rangle \\
&\quad \times \mathbb{E}_{\{e^i(\mathbf{w}_k^j)\}} \left[\left(\frac{e^r(\mathbf{w}_k^j) - p^r(\mathbf{w}_k^j)}{p^r(\mathbf{w}_k^j)} \right) \right. \\
&\quad \left. \times \left(\frac{e^t(\mathbf{w}_k^j) - p^t(\mathbf{w}_k^j)}{p^t(\mathbf{w}_k^j)} \right) \right]. \quad (\text{A.10})
\end{aligned}$$

The second term in (A.10) was shown to be zero in (A.9); thus, we have:

$$\begin{aligned}
&\mathbb{E}_{\{e^i(\mathbf{w}_k^j)\}} \left[\|\widehat{\mathbf{g}}(\mathbf{w}_k^j) - \mathbf{g}(\mathbf{w}_k^j)\|^2 \right] = \sum_{i=1}^B \lambda^i(\mathbf{w}_k^j)^2 \|\mathbf{a}^i(\mathbf{w}_k^j)\|^2 \\
&\quad \times \mathbb{E}_{\{e^i(\mathbf{w}_k^j)\}} \left[\left(\frac{e^i(\mathbf{w}_k^j) - p^i(\mathbf{w}_k^j)}{p^i(\mathbf{w}_k^j)} \right)^2 \right]. \quad (\text{A.11})
\end{aligned}$$

From (A.5) we know $\mathbb{E}_{\{e^i(\mathbf{w}_k^j)\}} \left[\left(\frac{e^i(\mathbf{w}_k^j) - p^i(\mathbf{w}_k^j)}{p^i(\mathbf{w}_k^j)} \right)^2 \right] = \left(\frac{1}{p^i(\mathbf{w}_k^j)} - 1 \right)$. Also, $\|\mathbf{a}^i(\mathbf{w}_k^j)\|^2 = 1$. Therefore, we can write:

$$\begin{aligned}
&\mathbb{E}_{\{e^i(\mathbf{w}_k^j)\}} \left[\|\widehat{\mathbf{g}}(\mathbf{w}_k^j) - \mathbf{g}(\mathbf{w}_k^j)\|^2 \right] \\
&= \sum_{i=1}^B \lambda^i(\mathbf{w}_k^j)^2 \left(\frac{1}{p^i(\mathbf{w}_k^j)} - 1 \right). \quad (\text{A.12})
\end{aligned}$$

This completes the proof. \square

APPENDIX B PROOF OF THEOREM 3

We want to find an upper bound for $\mathbb{E}_{\{\xi_j\}, \{e^i(\mathbf{w}_k^j)\}} \left[\|\widehat{\mathbf{g}}(\mathbf{w}_k) - \nabla F(\mathbf{w}_k)\|^2 \right]$; thus, we start by writing:

$$\begin{aligned}
&\mathbb{E}_{\{\xi_j\}, \{e^i(\mathbf{w}_k^j)\}} \left[\|\widehat{\mathbf{g}}(\mathbf{w}_k) - \nabla F(\mathbf{w}_k)\|^2 \right] \\
&= \mathbb{E}_{\{\xi_j\}, \{e^i(\mathbf{w}_k^j)\}} \left[\|\widehat{\mathbf{g}}(\mathbf{w}_k) - \mathbf{g}(\mathbf{w}_k)\| \right. \\
&\quad \left. + (\mathbf{g}(\mathbf{w}_k) - \nabla F(\mathbf{w}_k))\|^2 \right] \quad (\text{B.1})
\end{aligned}$$

$$\begin{aligned}
&\leq \mathbb{E}_{\{e^i(\mathbf{w}_k^j)\}} \left[\|\widehat{\mathbf{g}}(\mathbf{w}_k) - \mathbf{g}(\mathbf{w}_k)\|^2 \right] \\
&\quad + \mathbb{E}_{\{\xi_j\}} \left[\|\mathbf{g}(\mathbf{w}_k) - \nabla F(\mathbf{w}_k)\|^2 \right]. \quad (\text{B.2})
\end{aligned}$$

An upper bound for the second term in (B.2), following [8], is obtained as $\beta \|\nabla F(\mathbf{w}_k)\|^2 + \sigma$. For the first term, the upper bound is determined as follows:

$$\begin{aligned}
&\mathbb{E}_{\{e^i(\mathbf{w}_k^j)\}} \left[\|\widehat{\mathbf{g}}(\mathbf{w}_k) - \mathbf{g}(\mathbf{w}_k)\|^2 \right] \\
&= \mathbb{E}_{\{e^i(\mathbf{w}_k^j)\}} \left[\left\| \frac{1}{M} \sum_{j=1}^M (\widehat{\mathbf{g}}(\mathbf{w}_k^j) - \mathbf{g}(\mathbf{w}_k^j)) \right\|^2 \right]. \quad (\text{B.3})
\end{aligned}$$

From Lemma 2, we know that $\mathbb{E}_{\{e^i(\mathbf{w}_k^j)\}} \left[\|\widehat{\mathbf{g}}(\mathbf{w}_k^j) - \mathbf{g}(\mathbf{w}_k^j)\|^2 \right] = \sum_{i=1}^B \lambda^i(\mathbf{w}_k^j)^2 \left(\frac{1}{p^i(\mathbf{w}_k^j)} - 1 \right)$. Also, from Lemma 3, we know that $\sum_{i=1}^B \lambda^i(\mathbf{w}_k^j)^2 \left(\frac{1}{p^i(\mathbf{w}_k^j)} - 1 \right) = \frac{\sigma_{1,j}^k}{s_k} + \sigma_{2,j}^k$. Thus, we continue:

$$\begin{aligned}
&\mathbb{E}_{\{e^i(\mathbf{w}_k^j)\}} \left[\|\widehat{\mathbf{g}}(\mathbf{w}_k) - \mathbf{g}(\mathbf{w}_k)\|^2 \right] \\
&\leq \frac{1}{M} \sum_{j=1}^M \sum_{i=1}^B \lambda^i(\mathbf{w}_k^j)^2 \left(\frac{1}{p^i(\mathbf{w}_k^j)} - 1 \right) \quad (\text{B.4})
\end{aligned}$$

$$= \frac{1}{M} \left(\sum_{j=1}^M \frac{\sigma_{1,j}^k}{s_k} + \sigma_{2,j}^k \right). \quad (\text{B.5})$$

Therefore, from (B.2) and (B.5), we have:

$$\begin{aligned}
&\mathbb{E}_{\{\xi_j\}, \{e^i(\mathbf{w}_k^j)\}} \left[\|\widehat{\mathbf{g}}(\mathbf{w}_k) - \nabla F(\mathbf{w}_k)\|^2 \right] \\
&\leq \beta \|\nabla F(\mathbf{w}_k)\|^2 + \frac{\sigma_1}{s_k} + \sigma_2 \quad (\text{B.6})
\end{aligned}$$

where $\sigma_1 = \max_k \left(\frac{1}{M} \sum_{j=1}^M \sigma_{1,j}^k \right)$ and $\sigma_2 = \max_k \left(\frac{1}{M} \sum_{j=1}^M \sigma_{2,j}^k \right) + \sigma$. \square

APPENDIX C PROOF OF THEOREM 5

For $\psi_k(\tau_k, s_k)$ to be convex, its Hessian matrix must be positive semidefinite. We derive the Hessian matrix of $\psi_k(\tau_k, s_k)$ as follows:

$$\mathbf{H}(\psi_k(\tau_k, s_k)) = \begin{bmatrix} 2A \frac{\alpha s_k}{\tau_k^3} & -A \frac{\alpha}{\tau_k^2} - C \frac{\sigma_1}{s_k^2} \\ -A \frac{\alpha}{\tau_k^2} - C \frac{\sigma_1}{s_k^2} & 2B \frac{\sigma_1}{s_k^3} + 2C(\tau_k - 1) \frac{\sigma_1}{s_k^3} \end{bmatrix} \quad (\text{C.1})$$

where $A = \frac{2[F(\mathbf{w}_0) - F_{\text{inf}}]}{\eta T_k}$, $B = \frac{\eta L}{M}$, and $C = \eta^2 L^2$. If we ensure that both the diagonal elements and the determinant

are positive, then we have proven that the matrix is positive semidefinite which would in turn guarantee the convexity of $\psi_k(\tau_k, s_k)$. Clearly, the elements on the diagonal are positive; meanwhile for the determinant we have:

$$4\alpha\sigma_1 \left(\frac{AB}{\tau_k^3 s_k^2} + \frac{AC(\tau_k - 1)}{\tau_k^3 s_k^2} - \left(\frac{A^2\alpha}{\sigma_1 \tau_k^4} + \frac{C^2\sigma_1}{\alpha s_k^4} + \frac{2AC}{\tau_k^2 s_k^2} \right) \right) \quad (\text{C.2})$$

$$= 4\alpha\sigma_1 \left(\frac{4AB}{\tau_k^3 s_k^2} + \frac{2AC}{\tau_k^2 s_k^2} - \left(\frac{A^2\alpha}{\sigma_1 \tau_k^4} + \frac{C^2\sigma_1}{\alpha s_k^4} + \frac{4AC}{\tau_k^3 s_k^2} \right) \right) \quad (\text{C.3})$$

$$= \frac{2AC}{\tau_k^2 s_k^2} \left(1 - \frac{2}{\tau_k} \right) + \frac{4AB}{\tau_k^3 s_k^2} - \left(\frac{A^2\alpha}{\sigma_1 \tau_k^4} + \frac{C^2\sigma_1}{\alpha s_k^4} \right) \quad (\text{C.4})$$

$$\geq \frac{4B}{\tau_k^3 s_k^2} - \left(\frac{A\alpha}{\sigma_1 \tau_k^4} + \frac{C^2\sigma_1}{A\alpha s_k^4} \right) \quad (\text{C.5})$$

$$= \frac{4B}{\tau_k^3 s_k^2} - \left(\frac{A\alpha}{\sigma_1 \tau_k^4} + \frac{\eta^5 L^4 T_k \sigma_1}{2\alpha [F(\mathbf{w}_0) - F_{\inf}] s_k^4} \right) \quad (\text{C.6})$$

$$\approx \frac{4B}{\tau_k^3 s_k^2} - \frac{A\alpha}{\sigma_1 \tau_k^4} \quad (\text{C.7})$$

$$= \frac{1}{M\eta T_k \sigma_1 s_k^2 \tau_k^4} (2\eta L T_k \sigma_1 \tau_k - \alpha M s_k^2 (F(\mathbf{w}_0) - F_{\inf})) \quad (\text{C.8})$$

$$\geq 0. \quad (\text{C.9})$$

According to assumption (i), $\tau_k \geq 2$, that justifies the move to (C.5). Assumptions (ii) and (iii) help to ignore the third term in (C.6). The move to (C.8) was due to assumption (iv). \square

APPENDIX D

LEMMA 3

Lemma 3: The difference between compressed gradient $\hat{\mathbf{g}}(\mathbf{w}_k^j)$ and uncompressed gradient $\mathbf{g}(\mathbf{w}_k^j)$ for the j th worker is given by

$$\mathbb{E}_{\{e^i(\mathbf{w}_k^j)\}} \left[\|\hat{\mathbf{g}}(\mathbf{w}_k^j) - \mathbf{g}(\mathbf{w}_k^j)\|^2 \right] = \frac{\sigma_{1,j}^k}{s_k} + \sigma_{2,j}^k \quad (\text{D.1})$$

where $\sigma_{1,j}^k = \sum_{i=1}^B \left(\lambda^i(\mathbf{w}_k^j) \|\lambda(\mathbf{w}_k^j)\|_1 \right)$ and $\sigma_{2,j}^k = -\sum_{i=1}^B \lambda^i(\mathbf{w}_k^j)^2$.

Proof: To derive the difference between the compressed gradient $\hat{\mathbf{g}}(\mathbf{w}_k^j)$ and uncompressed gradient $\mathbf{g}(\mathbf{w}_k^j)$ for the j th worker, using Lemma 2 and Theorem 2, we have:

$$\mathbb{E}_{\{e^i(\mathbf{w}_k^j)\}} \left[\|\hat{\mathbf{g}}(\mathbf{w}_k^j) - \mathbf{g}(\mathbf{w}_k^j)\|^2 \right] = \sum_{i=1}^B \lambda^i(\mathbf{w}_k^j)^2 \left(\frac{1}{p^i(\mathbf{w}_k^j)} - 1 \right) \quad (\text{D.2})$$

$$= \sum_{i=1}^B \lambda^i(\mathbf{w}_k^j)^2 \left(\frac{\|\lambda(\mathbf{w}_k^j)\|_1}{\lambda^i(\mathbf{w}_k^j) s_k} - 1 \right) \quad (\text{D.3})$$

$$= \frac{1}{s_k} \sum_{i=1}^B \left(\lambda^i(\mathbf{w}_k^j) \|\lambda(\mathbf{w}_k^j)\|_1 \right) - \sum_{i=1}^B \lambda^i(\mathbf{w}_k^j)^2 \quad (\text{D.4})$$

$$= \frac{\sigma_{1,j}^k}{s_k} + \sigma_{2,j}^k \quad (\text{D.5})$$

where $\sigma_{1,j}^k = \sum_{i=1}^B \left(\lambda^i(\mathbf{w}_k^j) \|\lambda(\mathbf{w}_k^j)\|_1 \right)$ and $\sigma_{2,j}^k = -\sum_{i=1}^B \lambda^i(\mathbf{w}_k^j)^2$. \square



Milad Khademi Nori received his B.Sc. degree in electronics engineering as the 1st rank student from Semnan University, Semnan, Iran, in September 2016, and later his M.Sc. degree in digital electronics engineering from Amirkabir University of Technology, Tehran, Iran, in February 2019. Next, in August 2019, he joined the Electrical and Computer Engineering Department of Queen's University, Kingston, ON, Canada, where he is currently a Ph.D. student working on the junction of artificial intelligence and (wireless) communication in the Wireless Artificial Intelligence

Laboratory (WAI lab). In February 2016, he ranked 1st in the first phase of Iran's national universities olympiad in the field of electrical engineering, and in November 2018, he was awarded in Iran's national Information and Computer Technology (ICT) competition held by Huawei. His research interests include (Wireless) Communications, IoT, Machine Learning, and Federated Learning.



Sangseok Yun received the Ph.D. degree in electrical engineering from the Korea Advanced Institute of Science and Technology (KAIST), Daejeon, South Korea, in 2018. In 2018, he was a Post-Doctoral Fellow at Korea Advanced Institute of Science and Technology (KAIST), Daejeon, South Korea, and visited Queen's University, Kingston, Canada to conduct collaborative research. In 2019 and 2020, he was a Post-Doctoral Fellow at Queen's University, Kingston, Canada. Since 2021, he has been with Pukyong National University, Busan, South Korea, where he is currently an Assistant

Professor with the Department of Information and Communications Engineering. His research interests include theories and applications of deep learning, machine learning, wireless artificial intelligence, wireless communications, and physical layer security.



Il-Min Kim (SM'06) received the B.Sc. degree in electronics engineering from Yonsei University, Seoul, South Korea, in 1996, and the M.S. and Ph.D. degrees in electrical engineering from Korea Advanced Institute of Science and Technology (KAIST), Daejeon, South Korea, in 1998 and 2001, respectively. From October 2001 to August 2002, he was with the Department of Electrical Engineering and Computer Sciences, MIT, Cambridge, MA, USA, and from September 2002 to June 2003, he was with the Department of Electrical Engineering, Harvard University, Cambridge,

MA, USA, as a Post-Doctoral Research Fellow. In July 2003, he joined the Department of Electrical and Computer Engineering, Queen's University, Kingston, Canada, where he is currently a Professor. His current research interests include deep learning, reinforcement learning, continual learning, federated learning, Geo-AI, edge AI, deep learning for IoT and 6G. He has been an Editor for the IEEE TRANSACTIONS ON WIRELESS COMMUNICATIONS, the IEEE WIRELESS COMMUNICATIONS LETTERS, and the Journal of Communications and Networks (JCN).