

Personalized Federated Learning with Communication Compression

El Houcine Bergou*

*College of Computing
Mohammed VI Polytechnic University
Ben Guerir, Morocco*

Elhoucine.bergou@um6p.ma

Konstantin Burlachenko

King Abdullah University of Science and Technology, KSA

konstantin.burlachenko@kaust.edu.sa

Aritra Dutta†

*Artificial Intelligence Initiative
University of Central Florida
Orlando, Florida-32816*

aritra.dutta@ucf.edu

Peter Richtárik

King Abdullah University of Science and Technology, KSA

peter.richtarik@kaust.edu.sa

Reviewed on OpenReview: <https://openreview.net/forum?id=dZugyhbnFY>

Abstract

In contrast to training traditional machine learning (ML) models in data centers, federated learning (FL) trains ML models over local datasets on resource-constrained heterogeneous edge devices. Existing FL algorithms aim to learn a single global model for all participating devices, which may not be helpful to all devices participating in the training due to the heterogeneity of the data across the devices. Recently, Hanzely and Richtárik (2020) proposed a new formulation for training personalized FL models aimed at balancing the trade-off between the traditional global model and the local models that could be trained by individual devices using their private data only. They derived a new algorithm, called *loopless gradient descent* (L2GD), to solve it and showed that this algorithm leads to improved communication complexity guarantees in regimes when more personalization is required. In this paper, we equip their L2GD algorithm with a *bidirectional* compression mechanism to further reduce the communication bottleneck between the local devices and the server. Unlike other compression-based algorithms used in the FL setting, our compressed L2GD algorithm operates on a probabilistic communication protocol, where communication does not happen on a fixed schedule. Moreover, our compressed L2GD algorithm maintains a similar convergence rate as vanilla SGD without compression. To empirically validate the efficiency of our algorithm, we perform diverse numerical experiments on both convex and non-convex problems and use various compression techniques.¹

1 Introduction

We live in the era of big data, and edge devices have become a part of our daily lives. While the training of ML models using the diverse data stored on these devices is becoming increasingly popular, the traditional data center-based approach to training them faces serious *privacy issues* and has to deal with *high communication*

*Equal contribution

†Equal contribution

¹Our repository is available online: <https://github.com/burlachenkok/compressed-fl-l2gd-code>.

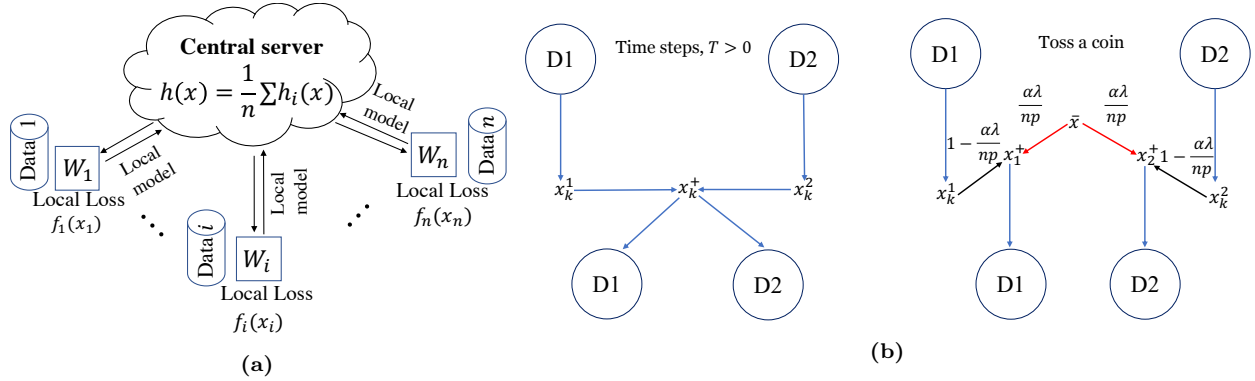


Figure 1: (a) Training n local devices, $\{W_i\}$ on the loss, f_i of their local model, x_i with a central server/master node, where h_i penalizes for dissimilarity between the local model, x_i and the average of all local models, \bar{x} . (b) FedAvg (McMahan et al., 2017) and L2GD (Hanzely & Richtárik, 2020) algorithm on 2 devices (D1 and D2). Unlike FedAvg, L2GD does not communicate after a fixed T local steps, it communicates based on a probabilistic protocol.

and energy cost associated with the transfer of data from users to the data center (Dean et al., 2012). *Federated learning* (FL) provides an attractive alternative to the traditional approach as it aims to train the models directly on *resource constrained* heterogeneous devices without any need for the data to leave them (Konečný et al., 2016; Kairouz et al., 2019).

The prevalent paradigm for training FL models is empirical risk minimization—to train a *single global model* using the aggregate of all the training data stored across all participating devices. Among the popular algorithms for training FL models for this formulation belong FedAvg (McMahan et al., 2017), Local GD (Khaled et al., 2019; 2020), local SGD (Stich, 2019; Khaled et al., 2020; Gorbunov et al., 2021) and Shifted Local SVRG (Gorbunov et al., 2021). All these methods require the participating devices to perform a local training procedure (e.g., by taking multiple steps of some optimization algorithm) and subsequently communicate the resulting model to an orchestrating server for aggregation; see Figure 1a. This process is repeated until a model of suitable qualities is found. For more variants of local methods and further pointers to the literature, we refer the reader to Gorbunov et al. (2021).

1.1 Personalized FL

In contrast, Hanzely & Richtárik (2020) recently introduced a new formulation of FL as an alternative to the existing “single-model-suits-all” approach embodied by empirical risk minimization. Their formulation explicitly aims to find a *personalized* model for every device; see Figure 1a. In particular, Hanzely & Richtárik (2020) considered the formulation²

$$\min_{x \in \mathbb{R}^{nd}} [F(x) := f(x) + h(x)] \quad (1)$$

for simultaneous training of n personalized FL models $x_1, \dots, x_n \in \mathbb{R}^d$ for n participating devices. They chose

$$f(x) := \frac{1}{n} \sum_{i=1}^n f_i(x_i), \quad \text{and} \quad h(x) := \frac{1}{n} \sum_{i=1}^n h_i(x),$$

where f_i represents the loss of model x_i over the local data stored on device i . Function h_i penalizes for dissimilarity between the local model x_i and the average of all local models $\bar{x} := \frac{1}{n} \sum_{i=1}^n x_i$, and is defined to be

$$h_i(x) = \frac{\lambda}{2} \|x_i - \bar{x}\|_2^2,$$

where $\lambda > 0$ controls for the strength of this penalization. At one extreme, $\lambda \rightarrow \infty$ forces the local models to be equal to their average, and hence, mutually identical. Therefore, equation 1 reduces to the classical

²Zhang et al. (2015) considered a similar model in a different context and with different motivations.

empirical risk minimization formulation of FL

$$\min_{z \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n f_i(z).$$

On the other hand, for $\lambda = 0$ problem equation 1 is equivalent to each client (node) training independently using their data only. In particular, the i^{th} client solves

$$\min_{x_i \in \mathbb{R}^d} f_i(x_i).$$

By choosing λ to a value in between these two extremes, i.e., $0 < \lambda < \infty$, we control for the level of similarity we want the personalized models $\{x_i\}_{i=1}^n$ to possess.

We remark that local methods such as FedAvg by McMahan et al. (2017) (also see similar methods in (Haddadpour et al., 2019b; Stich, 2019; Wang & Joshi, 2019; Zhou & Cong, 2018; Lin et al., 2020)), are popular for training FL models. Nevertheless, their main drawback in the heterogeneous setting with data and device heterogeneity is inefficient communication. Hanzely & Richtárik (2020) proposed this new personalization to tackle heterogeneous data, and we are using their model to build our compressed, personalized FL.

To solve equation 1, Hanzely & Richtárik (2020) proposed a *probabilistic* gradient descent algorithm for which they coined the name loopless local gradient descent (L2GD). Hanzely & Richtárik (2020) shows how L2GD can be interpreted as a simple variant of FedAvg, typically presented as a method for solving the standard empirical risk minimization (ERM) formulation of FL. However, alongside Hanzely & Richtárik (2020) argue, L2GD is better seen as an algorithm for solving the personalized FL formulation equation 1. By doing so, they interpret the nature of local steps in classical FL: the role of local steps in classical FL methods is to provide personalization and not communication efficiency as was widely believed—FedAvg can diverge on highly non-identical data partitions (McMahan et al., 2017). Instead, communication efficiency in local methods comes from their tendency to gear towards personalization, and personalized models are provably easier to train.

Communication compression. We observe that the *L2GD algorithm does not support any compression mechanism* for the master-worker and worker-master communication that needs to happen—This is the starting point of our work. *We believe that equipping personalized FL with fast and theoretically tractable communication compression mechanisms is an important open problem.*

In distributed training of deep neural network (DNN) models, synchronous data-parallelism (Dean et al., 2012) is most widely used and adopted by mainstream deep learning toolkits (such as PyTorch, TensorFlow). However, exchanging the stochastic gradients in the network for aggregation creates a communication bottleneck, and this results in slower training (Xu et al., 2021a). One way to save on communication costs is to use compression operators (Alistarh et al., 2017; Horváth et al., 2019; Xu et al., 2021a). Gradient compression techniques, such as quantization (Alistarh et al., 2017; Bernstein et al., 2018; Horváth et al., 2019; Beznosikov et al., 2020; Safaryan et al., 2020), sparsification (Suresh et al., 2017; Konečný & Richtárik, 2018; Aji & Heafield, 2017; Sahu et al., 2021; Stich et al., 2018; Dutta et al., 2020; Safaryan et al., 2020), hybrid compressors (Strom, 2015; Basu et al., 2019), and low-rank methods (Vogels et al., 2019) have been proposed to overcome this issue.³

Although recent works have introduced compression in traditional FL formulation (Konečný et al., 2016; Reisizadeh et al., 2020; Philippenko & Dieuleveut, 2020; Shlezinger et al., 2020; Amiri et al., 2020; Xu et al., 2021b); except (Horváth et al., 2019; Philippenko & Dieuleveut, 2020; Gorbunov et al., 2020; Amiri et al., 2020), others use compression only for the *throughput limited uplink* channel, that is, to upload the local models from the devices to the central server. But limited bandwidth in the downlink channel may pose a communication latency between the server and the devices and consequently, slow down the training; see detailed discussion in §2. As of now, no study combines *bidirectional* compression techniques with a probabilistic communication protocol in the FL set-up by using a mixture of a local and global model as in equation 1. In this work, we combine these aspects and make subsequent contributions.

³Model compression (Guo, 2018; Chraïbi et al., 2019) is orthogonal to gradient compression and not in the scope of this work.

1.2 Contributions

(i) **L2GD algorithm with bidirectional compression.** Communication compression is prevalent in recent local FL training algorithms, but these algorithms are not robust to data and device heterogeneity. L2GD by Hanzely & Richtárik (2020) remedies this issue by introducing personalization in FL. However, integrating compression with the L2GD algorithm is a nontrivial task—unlike other FL algorithms, L2GD does not communicate after fixed local steps, it communicates based on a probabilistic protocol; see §3 and Figure 1b. Additionally, due to this probabilistic protocol, the communication involves local model updates, as well as gradients; see §3. To reduce the communication bottleneck in L2GD, we use compression techniques on top of its probabilistic communication at both master and the participating local devices; see §4. To the best of our knowledge, we are the first to integrate *bidirectional compression techniques* with a probabilistic communication in the FL set-up, and we call our algorithm *compressed L2GD*; see Algorithm 1.

(ii) **Convergence analysis.** In §5, we prove the convergence of our *compressed L2GD* algorithm based on the most recent theoretical development, such as expected smoothness as in Gower et al. (2019). Admittedly, convergence analysis of first-order optimization algorithms with bidirectional compression exists in the literature, see (Tang et al., 2019; Horváth et al., 2019; Amiri et al., 2020; Dutta et al., 2020), integrating arbitrary unbiased compressors with a probabilistic communication protocol into personalized FL, and showing convergence are nontrivial and the first one in its class. Our compressed L2GD algorithm maintains a similar asymptotic convergence rate as the baseline vanilla SGD without compression in both strongly convex and smooth nonconvex cases; see Theorem 1 and 2 in §5.

(iii) **Optimal rate and communication.** We optimized the complexity bounds of our algorithm as a function of the parameters involved in the algorithm. This leads us to the *optimal* setting of our algorithm. Mainly, we derived the optimal expected number of local steps to get the optimal iteration complexity and communication rounds; see §6. Although our analysis is based on some hard-to-compute constants in real life, e.g., the Lipschitz constant, this may help the practitioners to get an insight into the iteration complexity and communication trade-off; see Theorem 3 and 4 in §6.

(iv) **Empirical study.** We perform diverse numerical experiments on synthetic and real datasets by using both convex and non-convex problems (using 4 DNN models) and invoking various compression techniques; see details in §7, Table 1. In training larger DNN models, to obtain the same global Top-1 test accuracy, compressed L2GD reduces the communicated data volume (bits normalized by the number of local devices or clients, #bits/n), from 10^{15} to 10^{11} , rendering approximately 10^4 times improvement compared to FedAvg; see §7.2. Moreover, L2GD with natural compressor (that by design has smaller variance) empirically behaves the best and converges approximately 5 times faster, and reaches the best accuracy on both train and the test sets compared to no-compression FedOpt (Reddi et al., 2020) baseline; see §7.2 and §A.2. These experiments validate the effect of the parameters used and, the effect of compressors, and show the efficiency of our algorithm in practice.

2 Related Work

Numerous studies are proposed to reduce communication but not all of them are in the FL setting. In this scope, for completeness, we quote a few representatives from each class of communication-efficient SGD.

Smith et al. (2017) proposed a communication-efficient primal-dual optimization that learns separate but related models for each participating device. FedAvg by McMahan et al. (2017) performs local steps on a subset of participating devices in an FL setting. Similar to FedAvg, but without data and device heterogeneity, (Haddadpour et al., 2019b; Stich, 2019; Wang & Joshi, 2019; Zhou & Cong, 2018; Lin et al., 2020) independently proposed local SGD, where several local steps are taken on the participating devices before periodic communication and averaging the local models. While FedProx by Li et al. (2020) is a generalization of FedAvg, SCAFFOLD uses a variance reduction to correct local updates occurring from non-i.i.d data in FedAvg. From the system’s perspective, on `TensorFlow`, Bonawitz et al. (2019) built a FL system on mobile devices.

Compression has also been introduced in the FL setup. Shlezinger et al. (2020) combined universal vector quantization with FL for throughput limited uplink channel. In FedPAQ by Reiszadeh et al. (2020), each local device sends a compressed difference between its input and output model to the central server, after computing the local updates for a fixed number of iterations. While Amiri et al. (2020) used a bidirectional compression in FL set-up, Philippenko & Dieuleveut (2020) combined it with a memory mechanism or error feedback (Stich et al., 2018). For a unified analysis of compression in FL we refer to (Haddadpour et al., 2021).

Among other proposed communication-efficient SGDs, parallel restarted SGD (Yu et al., 2019a) reduces the number of communication rounds compared to the baseline SGD. Haddadpour et al. (2019a) showed that redundancy reduces residual error as compared with the baseline SGD where all nodes can sample from the complete data and this leads to lower communication overheads. CoCoA by (Jaggi et al., 2014), and Dane by Shamir et al. (2014) perform several local steps and hence fewer communication rounds before communicating with the other workers. Lazily aggregated gradient (LAG) algorithm by Chen et al. (2018a) selects a subgroup of workers and uses their gradients, instead of obtaining a fresh gradient from each worker in each iteration. For communication-efficient local SGD see (Gao et al., 2021).

In decentralized training, where the nodes only communicate with their neighbors, Koloskova et al. (2019) implemented an *average consensus* where the nodes can communicate to their neighbors via a fixed communication graph. Li et al. (2018) proposed Pipe-SGD—a framework with decentralized pipelined training and balanced communication.

Personalization in FL is a growing research area. Arivazhagan et al. (2019) proposed FedPer to mitigate statistical heterogeneity of data; also, see adaptive personalized FL algorithm in (Deng et al., 2020). Mei et al. (2021) proposed to obtain personalization in FL by using layer-wise parameters, and two-stage training; also, see (Ma et al., 2022) and model personalization in (Shen et al., 2022). Shamsian et al. (2021) trained a central hypernetwork model to generate a set of personalized models for the local devices. Li et al. (2021) proposed Hermes—a communication-efficient personalized FL, where each local device identifies a small subnetwork by applying the structured pruning, communicates these subnetworks to the server and the devices, the server performs the aggregation on only overlapped parameters across each subnetwork; also, see Pillutla et al. (2022) for partial model personalization in FL. DispFL is another communication-efficient personalized FL algorithm proposed by Dai et al. (2022). In recent work, Zhang et al. (2021) introduces personalization by calculating optimal weighted model combinations for each client without assuming any data distribution. For a connection between personalization in FL and model-agnostic-meta-learning (MAML), see (Fallah et al., 2020). Additionally, we refer to the surveys (Kulkarni et al., 2020; Tan et al., 2021) for an overview of personalization in FL.

3 Background and Preliminaries

Notation. For a given vector, $x \in \mathbb{R}^{nd}$, by x_i we denote the i^{th} subvector of x , and write $x = (x_1^\top, \dots, x_n^\top)^\top$, where $x_i \in \mathbb{R}^d$. We denote the i^{th} component of x by $x_{(i)}$ and $\|x\|$ represents its Euclidean norm. By $[n]$ we denote the set of indexes, $\{1, \dots, n\}$. By $\mathbb{E}_\xi(\cdot)$ we define the expectation over the randomness of ξ conditional to all the other potential random variables. The operator, $\mathcal{C}(\cdot) := (\mathcal{C}_1(\cdot)^\top, \dots, \mathcal{C}_n(\cdot)^\top)^\top : \mathbb{R}^{nd} \rightarrow \mathbb{R}^{nd}$ denotes a compression operator with each $\mathcal{C}_i(\cdot)$ being compatible with the size of x_i . Denote $Q := [I, I, \dots, I]^\top \in \mathbb{R}^{nd \times d}$, where I denotes the identity matrix of $\mathbb{R}^{d \times d}$. With our Assumptions that we will introduce later in the paper, the problem in equation 1 has a unique solution, which we denote by x^* and we define \bar{x}^* as $\bar{x}^* = \frac{1}{n} \sum_{i=1}^n x_i^*$. By $|S|$ we denote the cardinality of a set, S .

Loopless local gradient descent (L2GD). We give a brief overview of the loopless local gradient descent (L2GD) algorithm by Hanzely & Richtárik (2020) to solve equation 1 as a two sum problem. At each iteration, to estimate the gradient of F , L2GD samples either the gradient of f or the gradient of h and updates the local models via:

$$x_i^{k+1} = x_i^k - \alpha G_i(x^k), \quad i = 1, \dots, n,$$

where $G_i(x^k)$ for $i = 1, \dots, n$, is the i^{th} block of the vector

$$G(x^k) = \begin{cases} \frac{\nabla f(x^k)}{1-p} & \text{with probability } 1-p, \\ \text{(Local gradient step)} \\ \frac{\nabla h(x^k)}{p} & \text{with probability } p, \\ \text{(Aggregation step)} \end{cases}$$

where $0 < p < 1$, $\nabla f(x^k)$ is the gradient of f at x^k , and $\nabla_i h(x^k) = \frac{\lambda}{n} (x_i^k - \bar{x}^k)$ is the i^{th} block of the gradient of h at x^k .

In this approach, there is a hidden communication between the local devices because in aggregation steps they need the average of the local models. That is, the communication occurs when the devices switch from a local gradient step to an aggregation step. Note that there is no need for communication between the local devices when they switch from an aggregation step to a local gradient step. There is also no need for communication after two consecutive aggregation steps since the average of the local models does not change in this case. If k and $k+1$ are both aggregation steps, we have $\bar{x}^{k+1} = \frac{1}{n} \sum_{i=1}^n x_i^{k+1} = \frac{1}{n} \sum_{i=1}^n x_i^k - \frac{\alpha\lambda}{n} \frac{1}{n} \sum_{i=1}^n (x_i^k - \bar{x}^k) = \bar{x}^k$.

4 Compressed L2GD

Now, we are all set to describe the compressed L2GD algorithm for solving (1). We start by defining how the compression operates in our set-up.

4.1 Compressed communication

Recall that the original L2GD algorithm has a probabilistic communication protocol—the devices do not communicate after every fixed number of local steps. The communication occurs when the devices switch from a local gradient step to an aggregation. Therefore, instead of using the compressors in a fixed time stamp (after every $T > 0$ iterations, say), each device i requires to compress its local model x_i when it needs to communicate it to the master, based on the probabilistic protocol. We assume that device i uses the compression operator, $\mathcal{C}_i(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^d$. Moreover, another compression happens when the master needs to communicate with the devices. We assume that the master uses the compression operator, $\mathcal{C}_M(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^d$. Therefore, the compression is used in uplink and downlink channels similar to Dutta et al. (2020); Horváth et al. (2019), but occurs in a probabilistic fashion. There exists another subtlety—although the model parameters (either from the local devices or the global aggregated model) are communicated in the network for training the FL model via compressed L2GD, the compressors that we use in this work are the compressors used for gradient compression in distributed DNN training; see (Xu et al., 2021a).

4.2 The algorithm

Note that, in each iteration $k \geq 0$, there exists a random variable, $\xi_k \in \{0, 1\}$ with $P(\xi_k = 1) = p$ and $P(\xi_k = 0) = 1 - p$. If $\xi_k = 0$, all local devices at iteration k perform one local gradient step. Otherwise (if $\xi_k = 1$), all local devices perform an aggregation step. However, to perform an aggregation step, the local devices need to know the average of the local models. If the previous iteration (i.e., $k - 1^{\text{th}}$ iteration) was an aggregation step (i.e., $\xi_{k-1} = 1$) then at the current iteration the local devices can use the same average as the one at iteration $k - 1$ (recall, the average of the local models does not change after two consecutive aggregation steps). Otherwise, a communication happens with the master to compute the average. In this case, each local device i compresses its local model x_i^k to $\mathcal{C}_i(x_i^k)$ and communicates the result to the master. The master computes the average based on the compressed values of local models:

$$\bar{y}^k := \frac{1}{n} \sum_{j=1}^n \mathcal{C}_j(x_j^k),$$

then it compresses \bar{y}^k to $\mathcal{C}_M(\bar{y}^k)$ by using a compression operator at the master's end and communicates it back to the local devices. The local devices further perform an aggregation step by using $\mathcal{C}_M(\bar{y}^k)$ instead of

the exact average. This process continues until convergence. From Algorithm 1, we have, for $i = 1, \dots, n$:

$$x_i^{k+1} = x_i^k - \eta G_i(x^k),$$

where

$$G_i(x^k) = \begin{cases} \frac{\nabla f_i(x_i^k)}{n(1-p)} & \text{if } \xi_k = 0 \\ \frac{\lambda}{np} (x_i^k - \mathcal{C}_M(\bar{y}^k)) & \text{if } \xi_k = 1 \ \& \ \xi_{k-1} = 0, \\ \frac{\lambda}{np} (x_i^k - \bar{x}^k) & \text{if } \xi_k = 1 \ \& \ \xi_{k-1} = 1. \end{cases}$$

We give the pseudo code in Algorithm 1.

Algorithm 1 Compressed L2GD

Input: $\{x_i^0\}_{i=1,\dots,n}$, stepsize $\eta > 0$, probability p , $\xi_{-1} = 1$, $\bar{x}^{-1} = \frac{1}{n} \sum_{i=1}^n x_i^0$.
for $k = 0, 1, 2, \dots$ **do**
 Draw: $\xi_k = 1$ with probability p
 if $\xi_k = 0$ **then**
 on all devices: $x_i^{k+1} = x_i^k - \frac{\eta}{n(1-p)} \nabla f_i(x_i^k)$ for $i \in [n]$
 else
 if $\xi_{k-1} = 0$ **then**
 on all devices: Compress x_i^k to $\mathcal{C}_i(x_i^k)$ and communicate $\mathcal{C}_i(x_i^k)$ to the master
 on master: receive $\mathcal{C}_i(x_i^k)$ from the device i , for all $i \in [n]$ compute $\bar{y}^k := \frac{1}{n} \sum_{j=1}^n \mathcal{C}_j(x_j^k)$
 compress \bar{y}^k to $\mathcal{C}_M(\bar{y}^k)$ and communicate it to all devices
 on all devices: Perform aggregation step $x_i^{k+1} = x_i^k - \frac{\eta\lambda}{np} (x_i^k - \mathcal{C}_M(\bar{y}^k))$
 else
 on all devices: $\bar{x}^k = \bar{x}^{k-1}$, Perform aggregation step $x_i^{k+1} = x_i^k - \frac{\eta\lambda}{np} (x_i^k - \bar{x}^k)$
 end if
 end if
end for

Remark 1 The initialization, ξ_{-1} is not important, as it does not impact the Algorithm. In Algorithm 1, we start the for loop (for $k = 0, 1, \dots$) by drawing ξ_k . At iteration k , we perform an aggregation step if $\xi_k = 1$. That is why to be consistent in the algorithm, we initialized ξ at iteration “-1” by 1 and initialized \bar{x}_{-1} by the average of the initial models, x_i^0 .

Remark 2 All local devices have access to the same value of ξ . We assumed that ξ is drawn at random at the server side and broadcasted to the local devices.

5 Convergence Analysis

With the above setup, we now prove the convergence of Algorithm 1; see detailed proofs in §A.1.

5.1 Assumptions

We make the following general assumptions in this paper.

Assumption 1 For $i = 1, \dots, n$:

- The compression operator, $\mathcal{C}_i(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is unbiased,

$$\mathbb{E}_{\mathcal{C}_i} [\mathcal{C}_i(x)] = x, \quad \forall x \in \mathbb{R}^d.$$

- There exists constant, $\omega_i > 0$ such that the variance of \mathcal{C}_i is bounded as follows:

$$\mathbb{E}_{\mathcal{C}_i} [\|\mathcal{C}_i(x) - x\|^2] \leq \omega_i \|x\|^2, \forall x \in \mathbb{R}^d.$$

- The operators, $\{\mathcal{C}_i(\cdot)\}_{i=1}^n$ are independent from each other, and independent from ξ_k , for all $k \geq 0$.
- The compression operator, $\mathcal{C}_M(\cdot)$ is unbiased, independent from $\{\mathcal{C}_i\}_{i=1}^n$ and has compression factor, ω_M .

From the above assumption we conclude that for all $x \in \mathbb{R}^d$, we have

$$\mathbb{E}_{\mathcal{C}_i} [\|\mathcal{C}_i(x)\|^2] \leq (1 + \omega_i)\|x\|^2.$$

The following lemma characterizes the compression factor, ω of the joint compression operator, $\mathcal{C}(\cdot) = (\mathcal{C}_1(\cdot)^\top, \dots, \mathcal{C}_n(\cdot)^\top)^\top$ as a function of $\omega_1, \dots, \omega_n$.

Lemma 1 *Let $x \in \mathbb{R}^{nd}$, then*

$$\mathbb{E}_{\mathcal{C}} [\|\mathcal{C}(x)\|^2] \leq (1 + \omega)\|x\|^2,$$

where $\omega = \max_{i=1, \dots, n} \{\omega_i\}$.

For the convergence of strongly convex functions, we require an additional assumption on the function f as follows.

Assumption 2 *We assume that f is L_f -smooth and μ -strongly convex.*

5.2 Auxiliary results

Before we state our main convergence theorem, we state several intermediate results needed for the convergence. In the following two lemmas, we show that based on the randomness of the compression operators, in expectation, we recover the exact average of the local models and the exact gradients for all iterations.

Lemma 2 *Let Assumption 1 hold, then for all $k \geq 0$, $\mathbb{E}_{\mathcal{C}, \mathcal{C}_M} [\mathcal{C}_M(\bar{y}^k)] = \bar{x}^k$.*

Lemma 3 *Let Assumptions 1 hold. Then for all $k \geq 0$, knowing x^k , $G(x^k)$ is an unbiased estimator of the gradient of function F at x^k .*

Our next lemma gives an upper bound on the iterate at each iteration. This bound is composed of two terms—the optimality gap, $F(x^k) - F(x^*)$, and the norm at the optimal point, x^* .

Lemma 4 *Let Assumption 2 hold, then*

$$\|x^k\|^2 \leq \frac{4}{\mu} (F(x^k) - F(x^*)) + 2\|x^*\|^2.$$

Lemma 5 helps us to prove the expected smoothness property (Gower et al., 2019). The bound in Lemma 5 is composed of—the optimality gap, the difference between the gradients of h at x^k and x^* , and an extra constant, β that depends on the used compressors.

Lemma 5 *Let Assumptions 1 and 2 hold. Then*

$$\mathcal{A} := \mathbb{E}_{\mathcal{C}_M, \mathcal{C}} \|x^k - Q\mathcal{C}_M(\bar{y}^k) - x^* + Q\mathcal{C}_M(\bar{y}^*)\|^2 \leq \frac{4n^2}{\lambda^2} \|\nabla h(x^k) - \nabla h(x^*)\|^2 + \alpha (F(x^k) - F(x^*)) + \beta,$$

where $\bar{y}^* := \frac{1}{n} \sum_{j=1}^n \mathcal{C}_j(x_j^*)$, $\alpha := \frac{4(4\omega + 4\omega_M(1+\omega))}{\mu}$, and $\beta := 2(4\omega + 4\omega_M(1+\omega))\|x^*\|^2 + 4\mathbb{E}_{\mathcal{C}_M, \mathcal{C}} \|Q\mathcal{C}_M(\bar{y}^*) - Q\bar{x}^*\|^2$.

Lemma 6 is the final result that (together with Lemma 5) shows the expected smoothness property and gives an upper bound on the stochastic gradient. This bound is composed of three terms—the optimality gap, $F(x^k) - F(x^*)$, the expected norm of the stochastic gradient at the optimal point, $\mathbb{E}\|G(x^*)\|^2$, and some other quantity that involves interplay between the parameters used in Algorithm 1 and the used compressor.

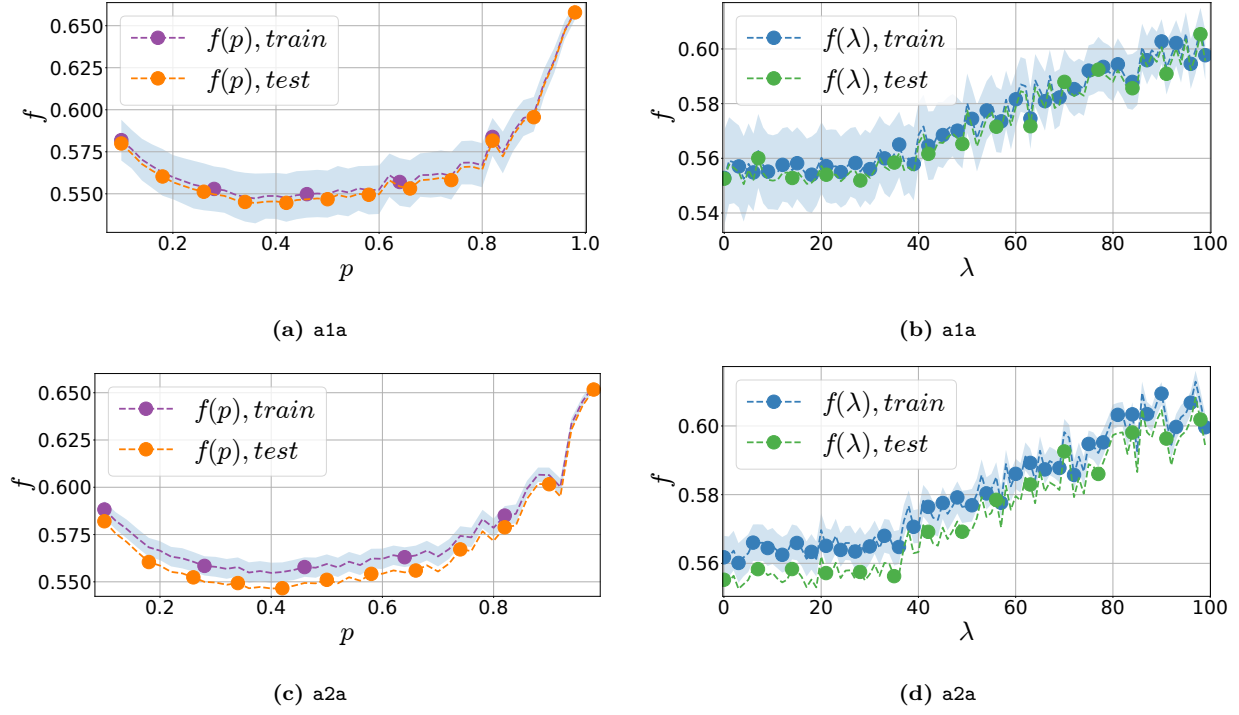


Figure 2: Uncompressed L2GD on $n = 5$ workers for $f_i(x)$ to be local empirical risk minimization for logistic regression loss with ℓ_2 regularization for local data D_i . We show the loss, f as a function of p and λ obtained after $K = 100$ iterations of Algorithm 1 with C an identity compressor (no compression). (a) a1a dataset, $d = 124, \lambda = 10$, (b) a1a dataset, $d = 124, p = 0.65$, (c) a2a dataset, $d = 124, \lambda = 10$, (d) a2a dataset, $d = 124, p = 0.65$.

Lemma 6 (Expected Smoothness) *Let Assumptions 1 and 2 hold, then*

$$\mathbb{E} [\|G(x^k)\|^2 | x^k] \leq 4\gamma (F(x^k) - F(x^*)) + \delta, \quad (2)$$

where

$$\gamma := \frac{\alpha\lambda^2(1-p)}{2n^2p} + \max \left\{ \frac{L_f}{(1-p)}, \frac{\lambda}{n} \left(1 + \frac{4(1-p)}{p} \right) \right\}$$

and

$$\delta := \frac{2\beta\lambda^2(1-p)}{n^2p} + 2\mathbb{E}\|G(x^*)\|^2.$$

Remark 3 *If there is no compression, the operators, $C_i(\cdot)$, for $i \in [n]$, and $C_M(\cdot)$ are equal to identity. The compression constants, ω_i , for $i \in [n]$, and ω_M are equal to zero. Therefore, $\alpha = \beta = 0$, and the factor 4 in the formula of γ can be replaced by 1 and thus*

$$\delta = \frac{2\beta\lambda^2(1-p)}{n^2p} + 2\mathbb{E}\|G(x^*)\|^2 = 2\mathbb{E}\|G(x^*)\|^2,$$

with

$$\gamma = \max \left\{ \frac{L_f}{(1-p)}, \frac{\lambda}{n} \left(1 + \frac{(1-p)}{p} \right) \right\} = \max \left\{ \frac{L}{n(1-p)}, \frac{\lambda}{np} \right\},$$

where $L = nL_f$. Same constants arise in the expected smoothness property in Hanzely & Richtárik (2020).

For nonconvex convergence of our algorithm, we cannot use the expected smoothness of Lemma 6, as it requires μ -strong convexity of the loss function. Therefore, similar to (Sahu et al., 2021; Stich & Karimireddy, 2020) we require a new assumption to bound the stochastic gradient, $G(x^k)$. Let $G(x^k)$ is of the form: $G(x^k) = \nabla F(x^k) + \zeta_k$, where ζ_k is the stochastic noise on the gradient.

Assumption 3 *There exist $M, \sigma^2 \geq 0$, such that $\mathbb{E}[\|G(x^k)\|^2 | x^k] \leq M\|\nabla F(x^k)\|^2 + \sigma^2$.*

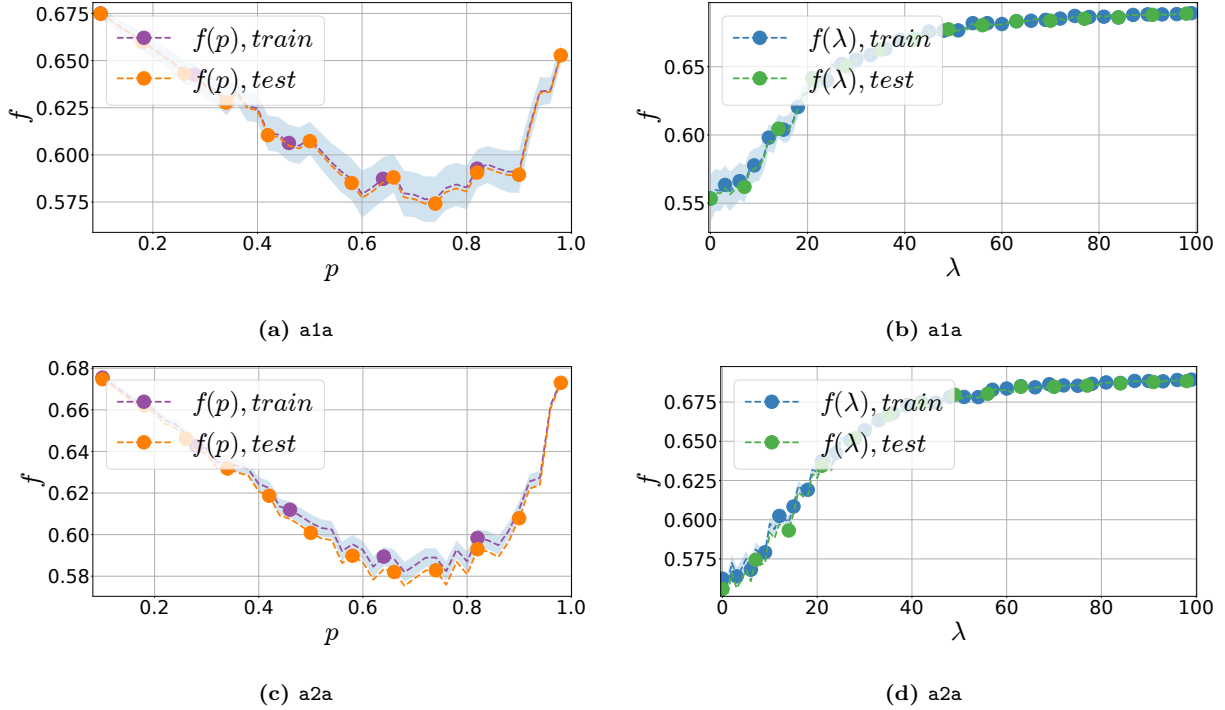


Figure 3: Compressed L2GD on $n = 5$ workers for $f_i(x)$ to be local empirical risk minimization for logistic regression loss with ℓ_2 regularization for local data D_i . We show the loss, f as a function of p and λ obtained after $K = 100$ iterations of Algorithm 1 with with natural compressor at the workers and identity compressor (no compression) at the master. (a) a1a dataset, $d = 124, \lambda = 10$, (b) a1a dataset, $d = 124, p = 0.65$, (c) a2a dataset, $d = 124, \lambda = 10$, (d) a2a dataset, $d = 124, p = 0.65$.

5.3 Main result

We now state the convergence result for Algorithm 1 for both strongly convex and nonconvex cases.

Theorem 1 (Strongly convex case) *Let Assumptions 1 and 2 hold. If $\eta \leq \frac{1}{2\gamma}$, then*

$$\mathbb{E} \|x^k - x^*\|^2 \leq \left(1 - \frac{\eta\mu}{n}\right)^k \|x^0 - x^*\|^2 + \frac{n\eta\delta}{\mu}.$$

Proof 1 *The proof follows directly from Lemma 3, 6, and Theorem 3.1 from Gower et al. (2019).*

Theorem 2 (Non convex case) *Let Assumptions 1 and 3 hold. Assume also that F is L_f -smooth, bounded from below by $F(x^*)$. Then to reach a precision, $\epsilon > 0$, set the stepsize, $\eta = \min\{\frac{1}{L_f M}, \frac{\epsilon^2}{2L_f\sigma^2}\}$, such that for $K \geq \frac{4L_f M(F(x^0) - F(x^*))}{\epsilon^2}$, we have $\min_{k=0,1,\dots,K} \mathbb{E} \|\nabla F(x^k)\|_2 \leq \epsilon$.*

Remark 4 *For smooth non-convex problems, we recover the optimal $O(\epsilon^4)$ classical rate as vanilla SGD.*

6 Optimal Rate and Communication

In this section, we provide the “optimal” setting of our algorithm that is obtained by optimizing the complexity bounds of our algorithm as a function of the parameters involved. The analysis on this section is based on

Table 1: Gradient compression methods used in this work. Note that $\|\tilde{g}\|_0$ and $\|g\|_0$ are the number of elements in the compressed and uncompressed gradient, respectively; nature of operator \mathcal{C} is random or deterministic. We implement that mechanisms for FedML.ai framework.

Compression	Ref.	Similar Methods	$\ \tilde{g}\ _0$	Nature of \mathcal{C}
QSGD	Alistarh et al. (2017)	Horváth et al. (2019); Wang et al. (2018); Wen et al. (2017) Wu et al. (2018); Yu et al. (2019b); Zhang et al. (2017)	$\ g\ _0$	Rand, unbiased
Natural	Horváth et al. (2019)	Alistarh et al. (2017); Yu et al. (2019b); Zhang et al. (2017)	$\ g\ _0$	Rand, unbiased
TernGrad	Wen et al. (2017)	Alistarh et al. (2017); Wang et al. (2018); Yu et al. (2019b)	$\ g\ _0$	Rand, unbiased
Bernoulli	Khairat et al. (2018)	—	—	Rand, unbiased
Top- k	Aji & Heafield (2017)	Alistarh et al. (2018); Stich et al. (2018)	k	Det, Biased

the following upper bound of γ . We recall that

$$\begin{aligned}\gamma &= \frac{\alpha\lambda^2(1-p)}{2n^2p} + \max\left\{\frac{L_f}{(1-p)}, \frac{\lambda}{n}\left(1 + \frac{4(1-p)}{p}\right)\right\} \\ &\leq \frac{\alpha\lambda^2(1-p)}{2n^2p} + \max\left\{\frac{L_f}{(1-p)}, \frac{4\lambda}{np}\right\} := \gamma_u.\end{aligned}$$

Note that the number of iterations is linearly dependent on γ . Therefore, to minimize the total number of iterations, it suffices to minimize γ . Define $L := nL_f$.

Theorem 3 (Optimal rate) *The probability p^* minimizing γ is equal to $\max\{p_e, p_A\}$, where $p_e = \frac{7\lambda+L-\sqrt{\lambda^2+14\lambda L+L^2}}{6\lambda}$ and p_A is the optimizer of the function $A(p) = \frac{\alpha\lambda^2}{2n^2p} + \frac{L}{n(1-p)}$ in $(0, 1)$.*

Remark 5 *If we maximize the upper bound, γ_u instead of γ then p_e simplifies to $\frac{4\lambda}{L+4\lambda}$.*

Lemma 7 *The optimizer probability p_A of the function $A(p) = \frac{\alpha\lambda^2}{2n^2p} + \frac{L}{n(1-p)}$ in $(0, 1)$ is equal to*

$$p_A = \begin{cases} \frac{1}{2} & \text{if } 2nL = \alpha\lambda^2 \\ \frac{-2\alpha\lambda^2+2\lambda\sqrt{2\alpha nL}}{2(2nL-\alpha\lambda^2)} & \text{if } 2nL > \alpha\lambda^2 \\ \frac{-2\alpha\lambda^2-2\lambda\sqrt{2\alpha nL}}{2(2nL-\alpha\lambda^2)} & \text{otherwise.} \end{cases}$$

Note that the number of communication rounds is linearly proportional to $C := p(1-p)\gamma$. Therefore, minimizing the total number of communication rounds suffices to minimize C or nC .

Theorem 4 (Optimal communication) *The probability p^* optimizing C is equal to $\max\{p_e, p_A\}$, where $p_e = \frac{7\lambda+L-\sqrt{\lambda^2+14\lambda L+L^2}}{6\lambda}$ and $p_A = 1 - \frac{Ln}{\alpha\lambda^2}$.*

Remark 6 *As in Remark 5, we note that, if we use the upper bound, γ_u instead of γ then p_e simplifies to $\frac{4\lambda}{L+4\lambda}$.*

We note that $\lambda \rightarrow 0$ implies $p^* \rightarrow 0$. This means that the optimal strategy, in this case, is *no communication at all*. This result is intuitive since for $\lambda = 0$, we deal with pure local models which can be computed without any communication. As $\lambda \rightarrow \infty$ implies $p^* \rightarrow 1$ denoting that the optimal strategy is to communicate often to find the global model.

7 Empirical Study

We conducted diverse numerical experiments with L2GD algorithm that includes: (i) Analysis of algorithm meta-parameters (with and without compression) for logistic regression in strongly convex setting; see §7.1; (ii) analysis of compressed L2GD algorithm on image classification with DNNs; see §7.2.

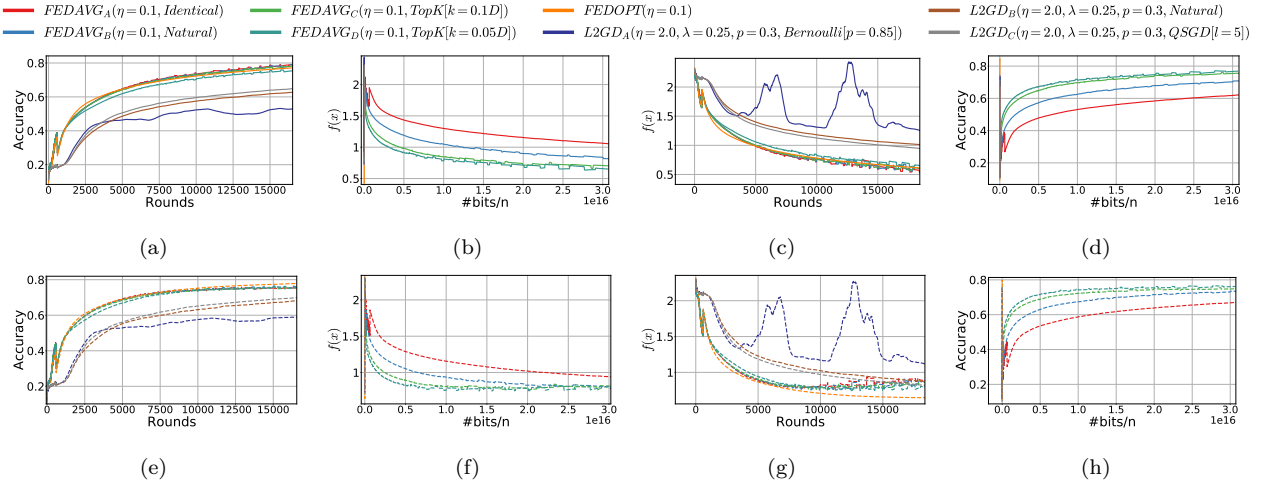


Figure 4: Training ResNet-18 on CIFAR-10 with $n = 10$ workers. The top row presents the Top-1 accuracy vs. rounds in (a), loss functional value vs. communicated bits in (b), loss functional value vs. rounds in (c), and Top-1 accuracy vs. communicated bits in (d) on the train set. The bottom row presents the similar plots on the Test set in (e)–(h).

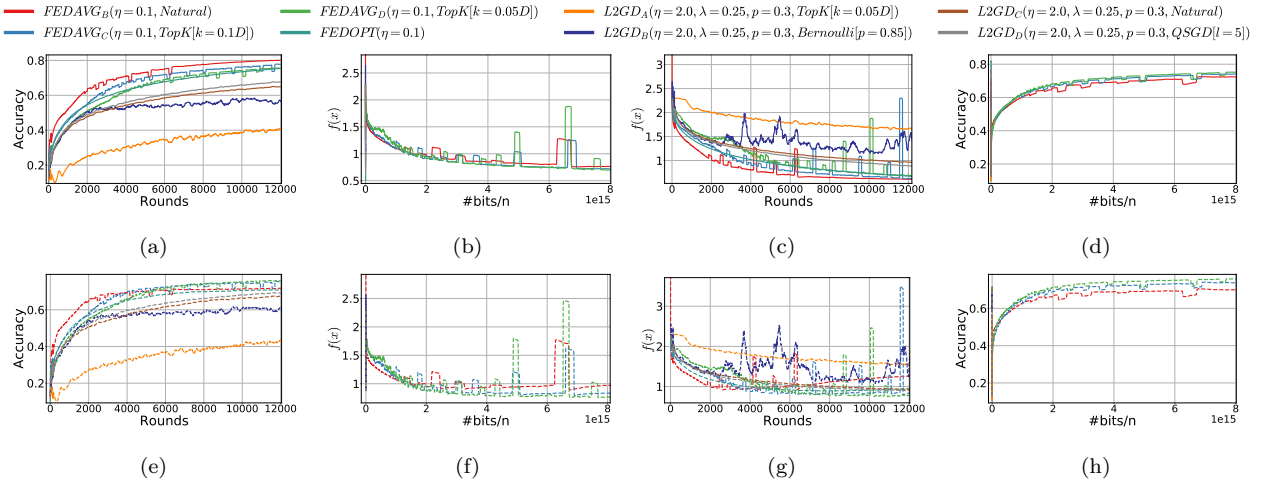


Figure 5: Training DenseNet-121 on CIFAR-10 with $n = 10$ workers. The top row represents the Top-1 accuracy vs. rounds in (a), loss functional value vs. communicated bits in (b), loss functional value vs. rounds in (c), and Top-1 accuracy vs. communicated bits in (d) on the train set. The bottom row presents the similar quantities on the Test set in (e)–(h).

Computing environment. We performed experiments on server-grade machines running Ubuntu 18.04 and Linux Kernel v5.4.0, equipped with 8-cores 3.3 GHz Intel Xeon and a single NVIDIA GeForce RTX 2080 Ti. Tesla-V100-SXM2 GPU with 32GB of GPU memory. The computation backend for Logistics Regression experiments was NumPy library with leveraging MPI4PY for inter-node communication. For DNNs we used recent version of FedML He et al. (2020) benchmark⁴ and patched it with: (i) distributed and standalone version of Algorithm 1; (ii) serializing and plotting mechanism; (iii) modifications in standalone, distributed version of FedAvg McMahan et al. (2017) and FedOpt Reddi et al. (2020) to be consistent with equation 1; (iv) not to drop the last batch while processing the dataset.

⁴FedML.AI

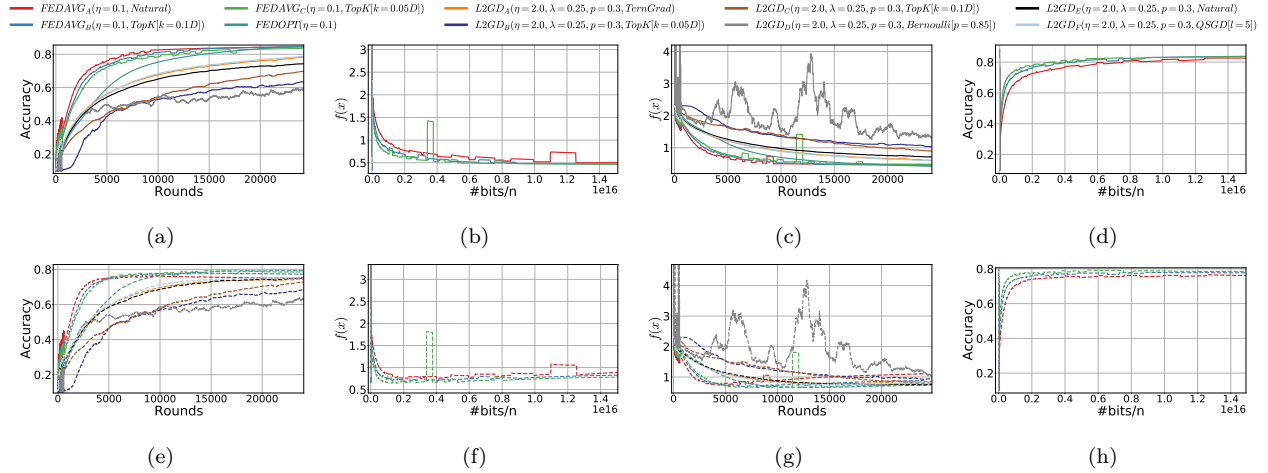


Figure 6: Training MobileNet on CIFAR-10 with $n = 10$ workers. The top row represents the Top-1 accuracy vs. rounds in (a), loss functional value vs. communicated bits in (b), loss functional value vs. rounds in (c), and Top-1 accuracy vs. communicated bits in (d) on the train set. The bottom row presents the similar plots on the Test set in (e)–(h).

7.1 Meta-parameter study

The purpose of these experiments is to study the meta-parameters involved in uncompressed and compressed L2GD algorithm. We used L2GD algorithm with and without compression for solving ℓ_2 regularized logistic regression on LIBSVM a1a and a2a datasets Chang & Lin (2011). Both datasets contain shuffled examples in the train set, and we did not perform any extra shuffling. To simulate the FL settings, we divided both datasets into 5 parts. After splitting, each worker has 321 and 453 records for a1a, and a2a, respectively.

Setup and results. We define $f_i(x)$ to be local empirical risk minimization for logistic loss with additive regularization term for local data D_i and of the form:

$$f_i(x) = \frac{1}{n_i} \sum_{j=1}^{n_i} \log(1 + \exp(-b^{(j)} x^\top a^{(j)})) + \frac{L_2}{2} \|x\|^2,$$

where $a^{(j)} \in \mathbb{R}^{124}$, $b^{(j)} \in \{+1, -1\}$, $n_i = |D_i|$. We set $L_2 = 1$, and varied meta-parameters p and λ . For each parameter, we performed 100 iterations of Algorithm 1. Note that, as meta-parameter λ decreases the models will fit more to its local data, while p provides stochastic balance between local gradient steps with probability, $1 - p$ and aggregation with probability, p . For compressed L2GD, we used natural compressor at the workers and identity compressor at the master.

Takeaway message. The results in Figure 2 support the theoretical finding—there exists an optimal choice of (p, λ) , where the loss function, f achieves the least value. Nevertheless, this choice is problem dependent. Additionally, we find small p is not good due to lack of samples in a single node compared to samples available at other nodes. There is a trade-off for each node in learning from the other nodes' data and spending time to learn from its own data. For uncompressed L2GD, the optimal setting of our algorithm is attained for $p = 0.4$ and λ in $[0, 25]$. The same observations hold for compressed L2GD; see Figure 3. For compressed L2GD with natural compressor, the optimal setting of our algorithm is attained for $p = 0.8$ and λ in $[0, 5]$. Finally, we observe that *to get the smallest errors on the training and validation sets, it is better not to perform the averaging step too often.*

7.2 Training DNN models

We choose four practically important DNN models used for image classification, and other down-streaming tasks, such as feature extractions for image segmentation, object detection, image embedding, image captioning, to name a few.

- **ResNet** (He et al., 2016). The overwhelmingly popular **ResNet** architecture exploits residual connections to remedy vanishing gradients. The network supported the trend toward smaller filters and deeper architectures, more curated towards FL training. We use **ResNet-18** and **ResNet-56** architectures (He et al., 2016).
- **DenseNet** (Huang et al., 2017) contains a short connection between layers via connecting each layer to every other layer in a feed-forward fashion. Dense connection allows propagating information to the final classifier via concatenating all feature maps. Each layer in **DenseNet** is narrow and contains only 12 filters—another practical model for FL training.
- **MobileNet** (Howard et al., 2017). DNN architecture has a trade off between computational complexity and accuracy (Bianco et al., 2018, p.3, Fig.1). For mobile devices that appear in cross-device FL, the computation cost and energy consumption are both important. The energy consumption is mostly driven by memory movement (Chen et al., 2018b; Horowitz, 2014). In **MobileNet** architecture standard convolution blocks performs depth-wise convolution followed by 1×1 convolution. This is computationally less expensive in flops during inference time (see (Bianco et al., 2018, Fig.1, p.3)) and is $\sim 3.5 \times$ more power efficient compare to **DenseNet** (García-Martín et al., 2019, p.85, Table 7). This makes **MobileNet** an attractive model for FL training.

Dataset and setup. We consider **CIFAR-10** dataset (Krizhevsky & Hinton, 2009) for image classification. It contains color images of resolution 28×28 from 10 classes. The training and the test set are of size, 5×10^4 and 10^4 , respectively. The training set is partitioned heterogeneously across 10 clients. The proportion of samples of each class stored at each local node is drawn by using the Dirichlet distribution ($\alpha = 0.5$). In our experiments, all clients are involved in each communication round. Additionally, we added a linear head in all CNN models for **CIFAR-10**, as they are originally designed for classification task with 1000 output classes.

Loss function. Denote $f_i(x) = w_i \cdot \frac{1}{|D_i|} \sum_{(a_i, b_i) \in D_i} l(a_i, b_i, x)$ to be a weighted local empirical risk associated with the local data, D_i stored in node, i . We note that $l(a_i, b_i, x)$ is a standard unweighted cross-entropy loss, $a_i \in \mathbb{R}^{28 \times 28 \times 3}$, $b_i \in \{0, 1\}^{10}$ with only one component equal to 1, the ground truth value, and the weight is set to $w_i = |D_i| / |D_1 \cup \dots \cup D_n|$.

Metrics. To measure the performance, we examine the loss function value, $f(x)$, and the Top-1 accuracy on both train and the test set. We use the weighted average of the local models, where the weight, w_i for each client is defined above. In our experiments, we wanted to assess the efficiency of both models — the local models and the global model (or the average model), $\frac{1}{n} \sum x_i$. To do this, we used two metrics:

- **Loss.** We compute the average loss over all the losses of the local models, that is, $f(x_1, \dots, x_n) = \frac{1}{n} \sum_i f_i(x_i)$. This allows us to assess the efficiency of the local models.
- **Accuracy.** We compute the accuracy of the average model, that is, we compute the accuracy using the model. As the experiments demonstrate both local models and the average model perform well.

We use state-of-the-art FedML benchmarking for our experiments, and it does not support personalization. Due to this limitation, although we wanted to present the average accuracy over all the accuracies of the local models, we were unable to do so — changing the FedML benchmarking for personalization is an involved task. Nevertheless, we compared with compressed communication approaches in the FedML framework that supports a single global model to see how despite personalization, our compressed L2GD performs. Additionally, we measure the number of rounds, and bits/ n —communicated bits normalized by the number of local clients, n . The intuition behind using the last metric is to measure the communicated data-volume; it is widely *hypothesized* that the reduced data-volume translates to a faster training in a constant speed network in distributed setup (Gajjala et al., 2020; Xu et al., 2021a).

Compressors used. The theoretical results of compressed L2GD are attributed to unbiased compressors. We used 4 different unbiased compressors at the clients: Bernoulli (Khirirat et al., 2018), natural compressor (Horváth et al., 2019), random dithering a.k.a. QSGD (Alistarh et al., 2017), and Terngrad (Wen et al., 2017); see Table 1 for details. Additionally, we note that biased compressors (mostly sparsifiers) are popular in DNN

Model	Training parameters	L2GD bits/n	Baseline bits/n
DenseNet-121	79×10^5	8×10^{11}	$4 \cdot 10^{15}$
MobileNet	32×10^5	1.7×10^{11}	1×10^{15}
ResNet-18	11×10^6	1.1×10^{12}	1.5×10^{16}

Table 2: Summary of the benchmarks. The measured quantity is bits/n to achieve 0.7 Top-1 test accuracy, with $n = 10$ clients. For DenseNet-121, MobileNet, Resnet-18 the baseline is FedAvg with natural compressor with 1 local epoch; L2GD also uses natural compressor.

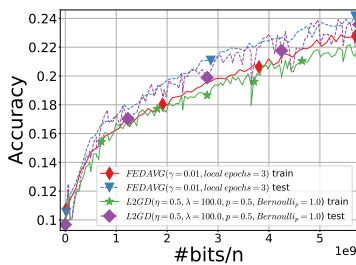


Figure 7: FedAvg as a particular case of L2GD: Test and train accuracy for ResNet-56 on CIFAR-10.

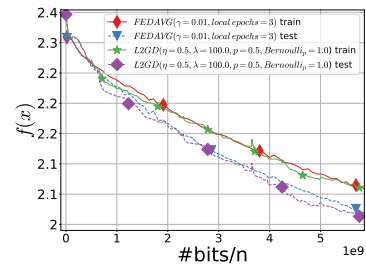


Figure 8: FedAvg as a particular case of L2GD: Test and train loss for ResNet-56 on CIFAR-10.

training. Therefore, out of scientific curiosity, we used a popular sparsifier: Top- k (Aji & Heafield, 2017; Sahu et al., 2021) as a proof of concept. We note that extending the compressed L2GD theory for biased compressors (with or without error-feedback (Xu et al., 2021a)) is nontrivial and mathematically involved, and left for future work.

Algorithms used for comparison. We used state-of-the-art FL training algorithms, FedAvg (McMahan et al., 2017) and FedOpt (Reddi et al., 2020) as no compression baseline to compare against our L2GD. However, the performance of FedAvg is not stable but improves with the compression mechanism. The original FedAvg algorithm does not contain any compression mechanism, but for comparison, we incorporated compressors into FedAvg via the following schema which is similar to the classic error feedback (Xu et al., 2021a): (i) After local steps, client estimates change of current iterate from the previous round and formulates direction, $g_{c,\text{computed}}^i$; (ii) client sends compressed difference between previous gradient estimator from previous round and currently computed gradient estimator, $\mathcal{C}(g_{c,\text{computed}}^i - g_c^{i-1})$ to the master; (iii) both master and client updating g_c^i via the following schema: $g_c^i = g_c^{i-1} + \mathcal{C}(g_{c,\text{computed}}^i - g_c^{i-1})$. We provide the details about step size and batch size in Appendix.

7.2.1 Results

We show the results for training ResNet-18, DenseNet-121, and MobileNet with compressed L2GD and other state-of-the-art FL algorithms in Figure 4–6. For these experiments, the communication rounds are set to 12×10^3 , 25×10^3 , and 20×10^3 , respectively. For the FedAvg algorithm, each client performs one epoch over the local data. We empirically tried 1, 2, 3, and 4 epochs over the local data as local steps, but one epoch is empirically the best choice.

For training ResNet-18, from Figure 4 we observe that FedAvg with compression has albeit better convergence than no compression FedAvg⁵. At the same time, compressed FedAvg affects the convergence as a function of communicated rounds only negligibly (see Figure 4 (d),(b)). Therefore, for training other DNN models we use FedAvg with compression and FedOpt without any compressors to enjoy the best of both baselines.

Take away message. Compressed L2GD with natural compressor sends the least data and drives the loss down the most in these experiments. At the same time, L2GD with natural compressor (by design it has smaller variance) reaches the best accuracy for both train and test sets. Compressed L2GD outperforms FedAvg by a huge margin—For all DNN experiments, to reach the desired Top-1 test accuracy, compressed L2GD reduces the communicated data-volume, #bits/n, from 10^{15} to 10^{11} , rendering approximately a 10^4 times improvement compared to FedAvg; see Table 2.

⁵We have observed that batch normalization (Ioffe & Szegedy, 2015) in ResNet is sensitive for aggregation; see our discussion in §A.2.

Interestingly, in training **MobileNet**, the performance of biased Top- k compressor degrades only about 10% compared to natural compressor, while approximately degrades 35% in training **DenseNet**. Additionally, see discussion in §A.2, Figures 9–11. This phenomena may lead the researchers to design unbiased compressors with smaller variance to empirically harvest the best behavior of compressed L2GD in personalized FL training.

Nevertheless, we also observe that compressed L2GD converges slower compared to other FL algorithms without compression in all cases. What follows, it can be argued, is that when we compare the communicated data volume for all DNN models, the convergence of compressed L2GD is much better. Additionally, the gain in terms of lowering the loss function value is significant—*by sending the same amount of data, L2GD lowers the loss the most compared to the other no-compression FL baseline algorithms*. These experiments also demonstrate that when communication is a bottleneck, FedAvg is not comparable with L2GD. The only comparable baseline for L2GD is FedOpt; see Table II, also, see discussion in §A.2, Figures 9–11. A similar observation holds for the Top-1 test and train accuracy. Taken together, these indicate that for training larger DNN models in a personalized FL settings, with resource constrained and geographically remote devices, compressed L2GD could be the preferred algorithm because its probabilistic communication protocol sends less data but obtains better test accuracy than no compression FedAvg and FedOpt.

Additionally, we observe that when $\frac{\eta\lambda}{np} \in [0.5, 0.95]$, compressed L2GD incurs a significant variance in objective function during training. Empirically, the best behavior was observed for $\frac{\eta\lambda}{np} \approx 1$ or $\frac{\eta\lambda}{np} \in (0, 0.17]$.

FedAvg as a particular case of L2GD. We note that if $\eta\lambda/np = 1$, then the aggregation step of Algorithm 1 reduces to $x_i^{k+1} = \bar{x}^k$, for all devices. Thus, in this regime L2GD works similarly as FedAvg with random number of local steps. E.g., if $p = 0.5$, then Algorithm 1 reduces to randomized version of FedAvg with an average of 3 local steps. Figures 7 and 8, confirm this observation numerically, where we see that both algorithms exhibit similar performance. For that experiment, we trained **ResNet-56** on **CIFAR10** with $n = 100$ workers, and for 600 rounds. For L2GD, we set $\frac{\eta\lambda}{pn} = 1$.

8 Conclusion and Future Direction

In this paper, we equipped the loopless gradient descent (L2GD) algorithm with a compression mechanism to reduce the communication bottleneck between local devices and the server in an FL context. We showed that the new algorithm enjoys similar convergence properties as the uncompressed L2GD with a natural increase in the stochastic gradient variance due to compression. This phenomenon is similar to classical convergence bounds for compressed SGD algorithms. We also show that in a personalized FL setting, there is a trade-off that must be considered by devices between learning from other devices’ data and spending time learning from their own data. However, a particular parameterization of our algorithm recovers the well-known FedAvg Algorithm. We assessed the performance of the new algorithm compared to the state-of-the-art and validated our theoretical insights through a large set of experiments.

Several questions remain open and merit further investigation in the future. For example, we plan on including compression when devices calculate their local updates, especially in an FL setting, as the devices might not be powerful, and the computing energy is limited, and examine how the algorithm behaves. Additionally, we observed the efficacy of compressed L2GD with a biased compressor, such as Top- k . Nevertheless, extending the compressed L2GD theory for biased compressors (with or without error-feedback (Xu et al., 2021a)) is nontrivial and challenging. In the future, we plan to prove a more general theory for compressed L2GD that include both biased and unbiased compressor operating in a bidirectional fashion. A more detailed meta-parameter study covering different network bandwidths, diverse ML tasks with different DNN architectures, and deploying the models on real-life, geographically remote servers will be our future empirical quest.

Acknowledgments

Aritra Dutta acknowledges being an affiliated researcher at the Pioneer Centre for AI, Denmark. The authors acknowledge many fruitful discussions with Md. Patel on this project while he was a remote undergraduate intern at KAUST.

References

- Alham Fikri Aji and Kenneth Heafield. Sparse communication for distributed gradient descent. In Martha Palmer, Rebecca Hwa, and Sebastian Riedel (eds.), *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, EMNLP 2017, Copenhagen, Denmark, September 9-11, 2017*, pp. 440–445. Association for Computational Linguistics, 2017. doi: 10.18653/V1/D17-1045. URL <https://doi.org/10.18653/v1/d17-1045>.
- Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. QSGD: communication-efficient SGD via gradient quantization and encoding. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pp. 1709–1720, 2017. URL <https://proceedings.neurips.cc/paper/2017/hash/6c340f25839e6acdc73414517203f5f0-Abstract.html>.
- Dan Alistarh, Torsten Hoeffler, Mikael Johansson, Nikola Konstantinov, Sarit Khirirat, and Cédric Renggli. The convergence of sparsified gradient methods. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pp. 5977–5987, 2018. URL <https://proceedings.neurips.cc/paper/2018/hash/314450613369e0ee72d0da7f6fee773c-Abstract.html>.
- Mohammad Mohammadi Amiri, Deniz Gündüz, Sanjeev R. Kulkarni, and H. Vincent Poor. Federated learning with quantized global model updates. *CoRR*, abs/2006.10672, 2020. URL <https://arxiv.org/abs/2006.10672>.
- Manoj Ghuhari Arivazhagan, Vinay Aggarwal, Aaditya Kumar Singh, and Sunav Choudhary. Federated learning with personalization layers. *CoRR*, abs/1912.00818, 2019. URL <http://arxiv.org/abs/1912.00818>.
- Debraj Basu, Deepesh Data, Can Karakus, and Suhas N. Diggavi. Qsparse-local-sgd: Distributed SGD with quantization, sparsification and local computations. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 14668–14679, 2019. URL <https://proceedings.neurips.cc/paper/2019/hash/d202ed5bcfa858c15a9f383c3e386ab2-Abstract.html>.
- Jeremy Bernstein, Yu-Xiang Wang, Kamyar Azizzadenesheli, and Animashree Anandkumar. SIGNSGD: compressed optimisation for non-convex problems. In Jennifer G. Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pp. 559–568. PMLR, 2018. URL <http://proceedings.mlr.press/v80/bernstein18a.html>.
- Aleksandr Beznosikov, Samuel Horváth, Peter Richtárik, and Mher Safaryan. On biased compression for distributed learning. *CoRR*, abs/2002.12410, 2020. URL <https://arxiv.org/abs/2002.12410>.
- Simone Bianco, Rémi Cadène, Luigi Celona, and Paolo Napoletano. Benchmark analysis of representative deep neural network architectures. *IEEE Access*, 6:64270–64277, 2018. doi: 10.1109/ACCESS.2018.2877890. URL <https://doi.org/10.1109/ACCESS.2018.2877890>.
- Kallista A. Bonawitz, Hubert Eichner, Wolfgang Grieskamp, Dzmitry Huba, Alex Ingerman, Vladimir Ivanov, Chloé Kiddon, Jakub Konečný, Stefano Mazzocchi, Brendan McMahan, Timon Van Overveldt, David Petrou, Daniel Ramage, and Jason Roselander. Towards federated learning at scale: System design. In Ameet Talwalkar, Virginia Smith, and Matei Zaharia (eds.), *Proceedings of Machine Learning and Systems 2019, MLSys 2019, Stanford, CA, USA, March 31 - April 2, 2019*. mlsys.org, 2019. URL <https://proceedings.mlsys.org/book/271.pdf>.

- Chih-Chung Chang and Chih-Jen Lin. LIBSVM: A library for support vector machines. *ACM Trans. Intell. Syst. Technol.*, 2(3):27:1–27:27, 2011. doi: 10.1145/1961189.1961199. URL <https://doi.org/10.1145/1961189.1961199>.
- Tianyi Chen, Georgios B. Giannakis, Tao Sun, and Wotao Yin. LAG: lazily aggregated gradient for communication-efficient distributed learning. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pp. 5055–5065, 2018a. URL <https://proceedings.neurips.cc/paper/2018/hash/feecee9f1643651799ede2740927317a-Abstract.html>.
- Y Chen, Tien-Ju Yang, Joel Emer, and Vivienne Sze. Understanding the limitations of existing energy-efficient design approaches for deep neural networks. *Energy*, 2(L1):L3, 2018b. URL <https://people.csail.mit.edu/emergemedia/papers/2018.02.sysml.energy-efficient-dnn.pdf>.
- Sélim Chraïbi, Ahmed Khaled, Dmitry Kovalev, Peter Richtárik, Adil Salim, and Martin Takác. Distributed fixed point methods with compressed iterates. *CoRR*, abs/1912.09925, 2019. URL <http://arxiv.org/abs/1912.09925>.
- Rong Dai, Li Shen, Fengxiang He, Xinmei Tian, and Dacheng Tao. Dispf: Towards communication-efficient personalized federated learning via decentralized sparse training. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato (eds.), *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pp. 4587–4604. PMLR, 2022. URL <https://proceedings.mlr.press/v162/dai22b.html>.
- Jeffrey Dean, Greg Corrado, Rajat Monga, Kai Chen, Matthieu Devin, Quoc V. Le, Mark Z. Mao, Marc’Aurelio Ranzato, Andrew W. Senior, Paul A. Tucker, Ke Yang, and Andrew Y. Ng. Large scale distributed deep networks. In Peter L. Bartlett, Fernando C. N. Pereira, Christopher J. C. Burges, Léon Bottou, and Kilian Q. Weinberger (eds.), *Advances in Neural Information Processing Systems 25: 26th Annual Conference on Neural Information Processing Systems 2012. Proceedings of a meeting held December 3-6, 2012, Lake Tahoe, Nevada, United States*, pp. 1232–1240, 2012. URL <https://proceedings.neurips.cc/paper/2012/hash/6aca97005c68f1206823815f66102863-Abstract.html>.
- Yuyang Deng, Mohammad Mahdi Kamani, and Mehrdad Mahdavi. Adaptive personalized federated learning. *CoRR*, abs/2003.13461, 2020. URL <https://arxiv.org/abs/2003.13461>.
- Aritra Dutta, El Houcine Bergou, Ahmed M. Abdelmoniem, Chen-Yu Ho, Atal Narayan Sahu, Marco Canini, and Panos Kalnis. On the discrepancy between the theoretical analysis and practical implementations of compressed communication for distributed deep learning. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pp. 3817–3824. AAAI Press, 2020. doi: 10.1609/AAAI.V34I04.5793. URL <https://doi.org/10.1609/aaai.v34i04.5793>.
- Alireza Fallah, Aryan Mokhtari, and Asuman E. Ozdaglar. Personalized federated learning with theoretical guarantees: A model-agnostic meta-learning approach. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/24389bfe4fe2eba8bf9aa9203a44cdad-Abstract.html>.
- Rishikesh R. Gajjala, Shashwat Banchhor, Ahmed M. Abdelmoniem, Aritra Dutta, Marco Canini, and Panos Kalnis. Huffman coding based encoding techniques for fast distributed deep learning. In *DistributedML@CoNEXT 2020: Proceedings of the 1st Workshop on Distributed Machine Learning, Barcelona, Spain, December 1, 2020*, pp. 21–27. ACM, 2020. doi: 10.1145/3426745.3431334. URL <https://doi.org/10.1145/3426745.3431334>.

- Hongchang Gao, An Xu, and Heng Huang. On the convergence of communication-efficient local SGD for federated learning. In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021*, pp. 7510–7518. AAAI Press, 2021. doi: 10.1609/AAAI.V35I9.16920. URL <https://doi.org/10.1609/aaai.v35i9.16920>.
- Eva García-Martín, Crefeda Faviola Rodrigues, Graham D. Riley, and Håkan Grahn. Estimation of energy consumption in machine learning. *J. Parallel Distributed Comput.*, 134:75–88, 2019. doi: 10.1016/J.JPDC.2019.07.007. URL <https://doi.org/10.1016/j.jpdc.2019.07.007>.
- Eduard Gorbunov, Dmitry Kovalev, Dmitry Makarenko, and Peter Richtárik. Linearly converging error compensated SGD. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/ef9280fbc5317f17d480e4d4f61b3751-Abstract.html>.
- Eduard Gorbunov, Filip Hanzely, and Peter Richtárik. Local SGD: unified theory and new efficient methods. In Arindam Banerjee and Kenji Fukumizu (eds.), *The 24th International Conference on Artificial Intelligence and Statistics, AISTATS 2021, April 13-15, 2021, Virtual Event*, volume 130 of *Proceedings of Machine Learning Research*, pp. 3556–3564. PMLR, 2021. URL <http://proceedings.mlr.press/v130/gorbunov21a.html>.
- Robert Mansel Gower, Nicolas Loizou, Xun Qian, Alibek Sailanbayev, Egor Shulgin, and Peter Richtárik. SGD: general analysis and improved rates. *CoRR*, abs/1901.09401, 2019. URL <http://arxiv.org/abs/1901.09401>.
- Yunhui Guo. A survey on methods and theories of quantized neural networks. *CoRR*, abs/1808.04752, 2018. URL <http://arxiv.org/abs/1808.04752>.
- Farzin Haddadpour, Mohammad Mahdi Kamani, Mehrdad Mahdavi, and Viveck R. Cadambe. Trading redundancy for communication: Speeding up distributed SGD for non-convex optimization. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pp. 2545–2554. PMLR, 2019a. URL <http://proceedings.mlr.press/v97/haddadpour19a.html>.
- Farzin Haddadpour, Mohammad Mahdi Kamani, Mehrdad Mahdavi, and Viveck R. Cadambe. Local SGD with periodic averaging: Tighter analysis and adaptive synchronization. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 11080–11092, 2019b. URL <https://proceedings.neurips.cc/paper/2019/hash/c17028c9b6e0c5deaad29665d582284a-Abstract.html>.
- Farzin Haddadpour, Mohammad Mahdi Kamani, Aryan Mokhtari, and Mehrdad Mahdavi. Federated learning with compression: Unified analysis and sharp guarantees. In Arindam Banerjee and Kenji Fukumizu (eds.), *The 24th International Conference on Artificial Intelligence and Statistics, AISTATS 2021, April 13-15, 2021, Virtual Event*, volume 130 of *Proceedings of Machine Learning Research*, pp. 2350–2358. PMLR, 2021. URL <http://proceedings.mlr.press/v130/haddadpour21a.html>.
- Filip Hanzely and Peter Richtárik. Federated learning of a mixture of global and local models. *CoRR*, abs/2002.05516, 2020. URL <https://arxiv.org/abs/2002.05516>.
- Chaoyang He, Songze Li, Jinhyun So, Mi Zhang, Hongyi Wang, Xiaoyang Wang, Praneeth Vepakomma, Abhishek Singh, Hang Qiu, Li Shen, Peilin Zhao, Yan Kang, Yang Liu, Ramesh Raskar, Qiang Yang, Murali Annavaram, and Salman Avestimehr. Fedml: A research library and benchmark for federated machine learning. *CoRR*, abs/2007.13518, 2020. URL <https://arxiv.org/abs/2007.13518>.

- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pp. 770–778. IEEE Computer Society, 2016. doi: 10.1109/CVPR.2016.90. URL <https://doi.org/10.1109/CVPR.2016.90>.
- Mark Horowitz. 1.1 computing’s energy problem (and what we can do about it). In *2014 IEEE International Conference on Solid-State Circuits Conference, ISSCC 2014, Digest of Technical Papers, San Francisco, CA, USA, February 9-13, 2014*, pp. 10–14. IEEE, 2014. doi: 10.1109/ISSCC.2014.6757323. URL <https://doi.org/10.1109/ISSCC.2014.6757323>.
- Samuel Horváth, Chen-Yu Ho, Ludovít Horváth, Atal Narayan Sahu, Marco Canini, and Peter Richtárik. Natural compression for distributed deep learning. *CoRR*, abs/1905.10988, 2019. URL <http://arxiv.org/abs/1905.10988>.
- Samuel Horváth, Dmitry Kovalev, Konstantin Mishchenko, Sebastian Stich, and Peter Richtárik. Stochastic Distributed Learning with Gradient Quantization and Variance Reduction. *arXiv:1904.05115*, 2019. URL <https://arxiv.org/abs/1904.05115>.
- Andrew G. Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *CoRR*, abs/1704.04861, 2017. URL <http://arxiv.org/abs/1704.04861>.
- Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger. Densely connected convolutional networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pp. 2261–2269. IEEE Computer Society, 2017. doi: 10.1109/CVPR.2017.243. URL <https://doi.org/10.1109/CVPR.2017.243>.
- Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In Francis R. Bach and David M. Blei (eds.), *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, volume 37 of *JMLR Workshop and Conference Proceedings*, pp. 448–456. JMLR.org, 2015. URL <http://proceedings.mlr.press/v37/ioffe15.html>.
- Martin Jaggi, Virginia Smith, Martin Takáč, Jonathan Terhorst, Sanjay Krishnan, Thomas Hofmann, and Michael I. Jordan. Communication-efficient distributed dual coordinate ascent. In Zoubin Ghahramani, Max Welling, Corinna Cortes, Neil D. Lawrence, and Kilian Q. Weinberger (eds.), *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, pp. 3068–3076, 2014. URL <https://proceedings.neurips.cc/paper/2014/hash/894b77f805bd94d292574c38c5d628d5-Abstract.html>.
- Peter Kairouz, H. Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Kallista A. Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, Rafael G. L. D’Oliveira, Salim El Rouayheb, David Evans, Josh Gardner, Zachary Garrett, Adrià Gascón, Badih Ghazi, Phillip B. Gibbons, Marco Gruteser, Zaïd Harchaoui, Chaoyang He, Lie He, Zhouyuan Huo, Ben Hutchinson, Justin Hsu, Martin Jaggi, Tara Javidi, Gauri Joshi, Mikhail Khodak, Jakub Konečný, Aleksandra Korolova, Farinaz Koushanfar, Sanmi Koyejo, Tancrede Lepoint, Yang Liu, Prateek Mittal, Mehryar Mohri, Richard Nock, Ayfer Özgür, Rasmus Pagh, Mariana Raykova, Hang Qi, Daniel Ramage, Ramesh Raskar, Dawn Song, Weikang Song, Sebastian U. Stich, Ziteng Sun, Ananda Theertha Suresh, Florian Tramèr, Praneeth Vepakomma, Jianyu Wang, Li Xiong, Zheng Xu, Qiang Yang, Felix X. Yu, Han Yu, and Sen Zhao. Advances and open problems in federated learning. *CoRR*, abs/1912.04977, 2019. URL <http://arxiv.org/abs/1912.04977>.
- Ahmed Khaled, Konstantin Mishchenko, and Peter Richtárik. First analysis of local GD on heterogeneous data. *CoRR*, abs/1909.04715, 2019. URL <http://arxiv.org/abs/1909.04715>.
- Ahmed Khaled, Konstantin Mishchenko, and Peter Richtárik. Tighter theory for local SGD on identical and heterogeneous data. In Silvia Chiappa and Roberto Calandra (eds.), *The 23rd International Conference on Artificial Intelligence and Statistics, AISTATS 2020, 26-28 August 2020, Online [Palermo, Sicily,*

- Italy*], volume 108 of *Proceedings of Machine Learning Research*, pp. 4519–4529. PMLR, 2020. URL <http://proceedings.mlr.press/v108/bayoumi20a.html>.
- Sarit Khirirat, Hamid Reza Feyzmahdavian, and Mikael Johansson. Distributed learning with compressed gradients. *arXiv preprint arXiv:1806.06573*, 2018. URL <https://arxiv.org/abs/1806.06573>.
- Anastasia Koloskova, Sebastian U. Stich, and Martin Jaggi. Decentralized stochastic optimization and gossip algorithms with compressed communication. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pp. 3478–3487. PMLR, 2019. URL <http://proceedings.mlr.press/v97/koloskova19a.html>.
- Jakub Konečný and Peter Richtárik. Randomized distributed mean estimation: Accuracy vs. communication. *Frontiers Appl. Math. Stat.*, 4:62, 2018. doi: 10.3389/FAMS.2018.00062. URL <https://doi.org/10.3389/fams.2018.00062>.
- Jakub Konečný, H. Brendan McMahan, Felix X. Yu, Peter Richtárik, Ananda Theertha Suresh, and Dave Bacon. Federated learning: Strategies for improving communication efficiency. *CoRR*, abs/1610.05492, 2016. URL <http://arxiv.org/abs/1610.05492>.
- Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. Technical Report 0, University of Toronto, Toronto, Ontario, 2009.
- Viraj Kulkarni, Milind Kulkarni, and Aniruddha Pant. Survey of personalization techniques for federated learning. *CoRR*, abs/2003.08673, 2020. URL <https://arxiv.org/abs/2003.08673>.
- Ang Li, Jingwei Sun, Pengcheng Li, Yu Pu, Hai Li, and Yiran Chen. Hermes: an efficient federated learning framework for heterogeneous mobile clients. In *ACM MobiCom '21: The 27th Annual International Conference on Mobile Computing and Networking, New Orleans, Louisiana, USA, October 25-29, 2021*, pp. 420–437. ACM, 2021. doi: 10.1145/3447993.3483278. URL <https://doi.org/10.1145/3447993.3483278>.
- Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. In Inderjit S. Dhillon, Dimitris S. Papailiopoulos, and Vivienne Sze (eds.), *Proceedings of Machine Learning and Systems 2020, MLSys 2020, Austin, TX, USA, March 2-4, 2020*. mlsys.org, 2020. URL <https://proceedings.mlsys.org/book/316.pdf>.
- Youjie Li, Mingchao Yu, Songze Li, Salman Avestimehr, Nam Sung Kim, and Alexander G. Schwing. Pipe-sgd: A decentralized pipelined SGD framework for distributed deep net training. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pp. 8056–8067, 2018. URL <https://proceedings.neurips.cc/paper/2018/hash/2c6a0bae0f071cbbf0bb3d5b11d90a82-Abstract.html>.
- Tao Lin, Sebastian U. Stich, Kumar Kshitij Patel, and Martin Jaggi. Don’t use large mini-batches, use local SGD. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. URL <https://openreview.net/forum?id=B1ey01BFPr>.
- Xiaosong Ma, Jie Zhang, Song Guo, and Wenchao Xu. Layer-wised model aggregation for personalized federated learning. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pp. 10082–10091. IEEE, 2022. doi: 10.1109/CVPR52688.2022.00985. URL <https://doi.org/10.1109/CVPR52688.2022.00985>.
- Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. Communication-efficient learning of deep networks from decentralized data. In Aarti Singh and Xiaojin (Jerry) Zhu (eds.), *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics, AISTATS 2017, 20-22 April 2017, Fort Lauderdale, FL, USA*, volume 54 of *Proceedings of Machine Learning Research*, pp. 1273–1282. PMLR, 2017. URL <http://proceedings.mlr.press/v54/mcmahan17a.html>.

- Yuan Mei, Binbin Guo, Danyang Xiao, and Weigang Wu. Fedvf: Personalized federated learning based on layer-wise parameter updates with variable frequency. In *IEEE International Performance, Computing, and Communications Conference, IPCCC 2021, Austin, TX, USA, October 29-31, 2021*, pp. 1–9. IEEE, 2021. doi: 10.1109/IPCCC51483.2021.9679416. URL <https://doi.org/10.1109/IPCCC51483.2021.9679416>.
- Constantin Philippenko and Aymeric Dieuleveut. Artemis: tight convergence guarantees for bidirectional compression in federated learning. *CoRR*, abs/2006.14591, 2020. URL <https://arxiv.org/abs/2006.14591>.
- Krishna Pillutla, Kshitiz Malik, Abdelrahman Mohamed, Michael G. Rabbat, Maziar Sanjabi, and Lin Xiao. Federated learning with partial model personalization. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato (eds.), *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pp. 17716–17758. PMLR, 2022. URL <https://proceedings.mlr.press/v162/pillutla22a.html>.
- Sashank J. Reddi, Zachary Charles, Manzil Zaheer, Zachary Garrett, Keith Rush, Jakub Konečný, Sanjiv Kumar, and H. Brendan McMahan. Adaptive federated optimization. *CoRR*, abs/2003.00295, 2020. URL <https://arxiv.org/abs/2003.00295>.
- Amirhossein Reiszadeh, Aryan Mokhtari, Hamed Hassani, Ali Jadbabaie, and Ramin Pedarsani. Fedpaq: A communication-efficient federated learning method with periodic averaging and quantization. In Silvia Chiappa and Roberto Calandra (eds.), *The 23rd International Conference on Artificial Intelligence and Statistics, AISTATS 2020, 26-28 August 2020, Online [Palermo, Sicily, Italy]*, volume 108 of *Proceedings of Machine Learning Research*, pp. 2021–2031. PMLR, 2020. URL <http://proceedings.mlr.press/v108/reiszadeh20a.html>.
- Mher Safaryan, Egor Shulgin, and Peter Richtárik. Uncertainty principle for communication compression in distributed and federated learning and the search for an optimal compressor. *CoRR*, abs/2002.08958, 2020. URL <https://arxiv.org/abs/2002.08958>.
- Atal Narayan Sahu, Aritra Dutta, Ahmed M. Abdelmoniem, Trambak Banerjee, Marco Canini, and Panos Kalnis. Rethinking gradient sparsification as total error minimization. In Marc’Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp. 8133–8146, 2021. URL <https://proceedings.neurips.cc/paper/2021/hash/447b0408b80078338810051bb38b177f-Abstract.html>.
- Ohad Shamir, Nathan Srebro, and Tong Zhang. Communication-efficient distributed optimization using an approximate newton-type method. In *Proceedings of the 31th International Conference on Machine Learning, ICML 2014, Beijing, China, 21-26 June 2014*, volume 32 of *JMLR Workshop and Conference Proceedings*, pp. 1000–1008. JMLR.org, 2014. URL <http://proceedings.mlr.press/v32/shamir14.html>.
- Aviv Shamsian, Aviv Navon, Ethan Fetaya, and Gal Chechik. Personalized federated learning using hypernetworks. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pp. 9489–9502. PMLR, 2021. URL <http://proceedings.mlr.press/v139/shamsian21a.html>.
- Yiqing Shen, Yuyin Zhou, and Lequan Yu. Cd2-pfed: Cyclic distillation-guided channel decoupling for model personalization in federated learning. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pp. 10031–10040. IEEE, 2022. doi: 10.1109/CVPR52688.2022.00980. URL <https://doi.org/10.1109/CVPR52688.2022.00980>.
- Nir Shlezinger, Mingzhe Chen, Yonina C. Eldar, H. Vincent Poor, and Shuguang Cui. Federated learning with quantization constraints. In *2020 IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2020, Barcelona, Spain, May 4-8, 2020*, pp. 8851–8855. IEEE, 2020. doi: 10.1109/ICASSP40776.2020.9054168. URL <https://doi.org/10.1109/ICASSP40776.2020.9054168>.

- Virginia Smith, Chao-Kai Chiang, Maziar Sanjabi, and Ameet Talwalkar. Federated multi-task learning. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pp. 4424–4434, 2017. URL <https://proceedings.neurips.cc/paper/2017/hash/6211080fa89981f66b1a0c9d55c61d0f-Abstract.html>.
- Sebastian U. Stich. Local SGD converges fast and communicates little. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019. URL <https://openreview.net/forum?id=Sig2JnRcFX>.
- Sebastian U. Stich and Sai Praneeth Karimireddy. The error-feedback framework: Better rates for sgd with delayed gradients and compressed updates. *The Journal of Machine Learning Research*, 21(1):9613–9648, 2020. URL <https://arxiv.org/abs/1909.05350>.
- Sebastian U. Stich, Jean-Baptiste Cordonnier, and Martin Jaggi. Sparsified SGD with memory. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pp. 4452–4463, 2018. URL <https://proceedings.neurips.cc/paper/2018/hash/b440509a0106086a67bc2ea9df0a1dab-Abstract.html>.
- Nikko Strom. Scalable distributed DNN training using commodity GPU cloud computing. In *INTERSPEECH 2015, 16th Annual Conference of the International Speech Communication Association, Dresden, Germany, September 6-10, 2015*, pp. 1488–1492. ISCA, 2015. doi: 10.21437/INTERSPEECH.2015-354. URL <https://doi.org/10.21437/Interspeech.2015-354>.
- Ananda Theertha Suresh, Felix X. Yu, Sanjiv Kumar, and H. Brendan McMahan. Distributed mean estimation with limited communication. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pp. 3329–3337. PMLR, 2017. URL <http://proceedings.mlr.press/v70/suresh17a.html>.
- Alysa Ziyang Tan, Han Yu, Lizhen Cui, and Qiang Yang. Towards personalized federated learning. *CoRR*, abs/2103.00710, 2021. URL <https://arxiv.org/abs/2103.00710>.
- Hanlin Tang, Chen Yu, Xiangru Lian, Tong Zhang, and Ji Liu. Doublesqueeze: Parallel stochastic gradient descent with double-pass error-compensated compression. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pp. 6155–6165. PMLR, 2019. URL <http://proceedings.mlr.press/v97/tang19d.html>.
- Thijs Vogels, Sai Praneeth Karimireddy, and Martin Jaggi. Powersgd: Practical low-rank gradient compression for distributed optimization. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 14236–14245, 2019. URL <https://proceedings.neurips.cc/paper/2019/hash/d9fbed9da256e344c1fa46bb46c34c5f-Abstract.html>.
- Hongyi Wang, Scott Sievert, Shengchao Liu, Zachary Charles, Dimitris S. Papailiopoulos, and Stephen J. Wright. ATOMO: communication-efficient learning via atomic sparsification. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pp. 9872–9883, 2018. URL <https://proceedings.neurips.cc/paper/2018/hash/33b3214d792caf311e1f00fd22b392c5-Abstract.html>.

- Jianyu Wang and Gauri Joshi. Adaptive communication strategies to achieve the best error-runtime trade-off in local-update SGD. In Ameet Talwalkar, Virginia Smith, and Matei Zaharia (eds.), *Proceedings of Machine Learning and Systems 2019, MLSys 2019, Stanford, CA, USA, March 31 - April 2, 2019*. mlsys.org, 2019. URL <https://proceedings.mlsys.org/book/257.pdf>.
- Wei Wen, Cong Xu, Feng Yan, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. Terngrad: Ternary gradients to reduce communication in distributed deep learning. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pp. 1509–1519, 2017. URL <https://proceedings.neurips.cc/paper/2017/hash/89fcd07f20b6785b92134bd6c1d0fa42-Abstract.html>.
- Jiaxiang Wu, Weidong Huang, Junzhou Huang, and Tong Zhang. Error compensated quantized SGD and its applications to large-scale distributed optimization. In Jennifer G. Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pp. 5321–5329. PMLR, 2018. URL <http://proceedings.mlr.press/v80/wu18d.html>.
- Hang Xu, Chen-Yu Ho, Ahmed M. Abdelmoniem, Aritra Dutta, El Houcine Bergou, Konstantinos Karatsenidis, Marco Canini, and Panos Kalnis. GRACE: A compressed communication framework for distributed machine learning. In *41st IEEE International Conference on Distributed Computing Systems, ICDCS 2021, Washington DC, USA, July 7-10, 2021*, pp. 561–572. IEEE, 2021a. doi: 10.1109/ICDCS51616.2021.00060. URL <https://doi.org/10.1109/ICDCS51616.2021.00060>.
- Hang Xu, Kelly Kostopoulou, Aritra Dutta, Xin Li, Alexandros Ntoulas, and Panos Kalnis. Deepreduce: A sparse-tensor communication framework for federated deep learning. In Marc’Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp. 21150–21163, 2021b. URL <https://proceedings.neurips.cc/paper/2021/hash/b0ab42fcb7133122b38521d13da7120b-Abstract.html>.
- Hao Yu, Rong Jin, and Sen Yang. On the linear speedup analysis of communication efficient momentum SGD for distributed non-convex optimization. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pp. 7184–7193. PMLR, 2019a. URL <http://proceedings.mlr.press/v97/yu19d.html>.
- Yue Yu, Jiaxiang Wu, and Junzhou Huang. Exploring fast and communication-efficient algorithms in large-scale distributed networks. In Kamalika Chaudhuri and Masashi Sugiyama (eds.), *The 22nd International Conference on Artificial Intelligence and Statistics, AISTATS 2019, 16-18 April 2019, Naha, Okinawa, Japan*, volume 89 of *Proceedings of Machine Learning Research*, pp. 674–683. PMLR, 2019b. URL <http://proceedings.mlr.press/v89/yu19a.html>.
- Hantian Zhang, Jerry Li, Kaan Kara, Dan Alistarh, Ji Liu, and Ce Zhang. Zipml: Training linear models with end-to-end low precision, and a little bit of deep learning. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pp. 4035–4043. PMLR, 2017. URL <http://proceedings.mlr.press/v70/zhang17e.html>.
- Michael Zhang, Karan Sapra, Sanja Fidler, Serena Yeung, and José M. Álvarez. Personalized federated learning with first order model optimization. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. URL <https://openreview.net/forum?id=ehJqJQk9cw>.
- Sixin Zhang, Anna Choromanska, and Yann LeCun. Deep learning with elastic averaging SGD. In Corinna Cortes, Neil D. Lawrence, Daniel D. Lee, Masashi Sugiyama, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems*

2015, December 7-12, 2015, Montreal, Quebec, Canada, pp. 685–693, 2015. URL <https://proceedings.neurips.cc/paper/2015/hash/d18f655c3fce66ca401d5f38b48c89af-Abstract.html>.

Fan Zhou and Guojing Cong. On the convergence properties of a k-step averaging stochastic gradient descent algorithm for nonconvex optimization. In Jérôme Lang (ed.), *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI 2018, July 13-19, 2018, Stockholm, Sweden*, pp. 3219–3227. ijcai.org, 2018. doi: 10.24963/IJCAI.2018/447. URL <https://doi.org/10.24963/ijcai.2018/447>.

A Appendix

Contents

1	Introduction	1
1.1	Personalized FL	2
1.2	Contributions	4
2	Related Work	4
3	Background and Preliminaries	5
4	Compressed L2GD	6
4.1	Compressed communication	6
4.2	The algorithm	6
5	Convergence Analysis	7
5.1	Assumptions	7
5.2	Auxiliary results	8
5.3	Main result	10
6	Optimal Rate and Communication	10
7	Empirical Study	11
7.1	Meta-parameter study	13
7.2	Training DNN models	13
7.2.1	Results	15
8	Conclusion and Future Direction	16
A	Appendix	26
A.1	Convergence Analysis—Proofs of the Lemmas and the Theorems	26
A.1.1	Technical results used for convergence	27
A.1.2	Main convergence results	28
A.1.3	Nonconvex convergence	30
A.1.4	Optimal rate and communication	31
A.2	Addendum to the Experimental Results	32

A.1 Convergence Analysis—Proofs of the Lemmas and the Theorems

In this section, we provide the proofs of convex and non-convex convergence results of the compressed L2GD algorithm.

Overview of results. In §A.1.1, we provide the technical lemmas necessary for the analyses. §A.1.2 contains the auxiliary results pertaining to both convex and nonconvex convergence. In §A.1.3 we provide the non-convex convergence results, and §A.1.4 contains the proofs for optimal rate and communication.

A.1.1 Technical results used for convergence

The following two Lemmas are instrumental in proving other compression related results.

Lemma 1 *Let $x \in R^{nd}$, then*

$$\mathbb{E}_{\mathcal{C}} [\|\mathcal{C}(x)\|^2] \leq (1 + \omega)\|x\|^2,$$

where $\omega = \max_{i=1, \dots, n} \{\omega_i\}$.

Proof 2 *By using Assumption 1, we have*

$$\mathbb{E}_{\mathcal{C}} [\|\mathcal{C}(x)\|^2] = \mathbb{E}_{\mathcal{C}} \left[\sum_{i=1}^n \|\mathcal{C}_i(x_i)\|^2 \right] = \sum_{i=1}^n \mathbb{E}_{\mathcal{C}_i} \|\mathcal{C}_i(x_i)\|^2 \leq \sum_{i=1}^n (1 + \omega_i) \|x_i\|^2 \leq (1 + \omega) \|x\|^2.$$

Hence the result.

Lemma 2 *Let Assumption 1 hold, then for all $k \geq 0$, $\mathbb{E}_{\mathcal{C}, \mathcal{C}_M} [\mathcal{C}_M(\bar{y}^k)] = \bar{x}^k$.*

Proof 3 *We have*

$$\mathbb{E}_{\mathcal{C}, \mathcal{C}_M} [\mathcal{C}_M(\bar{y}^k)] = \mathbb{E}_{\mathcal{C}} [\mathbb{E}_{\mathcal{C}_M} [\mathcal{C}_M(\bar{y}^k)]] = \mathbb{E}_{\mathcal{C}} \left[\frac{1}{n} \sum_{j=1}^n \mathcal{C}_j(x_j^k) \right] = \frac{1}{n} \sum_{j=1}^n \mathbb{E}_{\mathcal{C}_j} [\mathcal{C}_j(x_j^k)] = \bar{x}^k.$$

Hence the result.

In the following Lemma, we show that based on the randomness of the compression operators, in expectation, we recover the exact average of the local models and the exact gradient for all iterations.

Lemma 3 *Let Assumptions 1 hold. Then for all $k \geq 0$, knowing x^k , $G(x^k)$ is an unbiased estimator of the gradient of function F at x^k .*

Proof 4 *We have*

$$\begin{aligned} \mathbb{E}_{\mathcal{C}, \mathcal{C}_M} [G_i(x^k)] &= \begin{cases} \frac{\nabla f_i(x_i^k)}{n(1-p)} & \text{if } \xi_k = 0 \\ \frac{\lambda}{np} (x_i^k - \mathbb{E}_{\mathcal{C}, \mathcal{C}_M} [\mathcal{C}_M(\bar{y}^k)]) & \text{if } \xi_k = 1 \text{ \& } \xi_{k-1} = 0, \\ \frac{\lambda}{np} (x_i^k - \bar{x}^k) & \text{if } \xi_k = 1 \text{ \& } \xi_{k-1} = 1, \end{cases} \\ \stackrel{\text{By Lemma 2}}{=} & \begin{cases} \frac{\nabla f_i(x_i^k)}{n(1-p)} & \text{if } \xi_k = 0, \\ \frac{\lambda}{np} (x_i^k - \bar{x}^k) & \text{if } \xi_k = 1. \end{cases} \end{aligned}$$

Therefore,

$$\begin{aligned} \mathbb{E}[G_i(x^k)|x^k] &= \mathbb{E}_{\xi_k} [\mathbb{E}_{\mathcal{C}, \mathcal{C}_M} [G_i(x^k)]] \\ &= (1-p) \frac{\nabla f_i(x_i^k)}{n(1-p)} + p \frac{\lambda}{np} (x_i^k - \bar{x}^k) \\ &= \nabla_{x_i} f(x^k) + \nabla_{x_i} h(x^k) = \nabla_{x_i} F(x^k). \end{aligned}$$

Hence the result.

A.1.2 Main convergence results

Based on the results given in the previous section, we are now set to quote our key convergence results. Our next lemma gives an upper bound on the iterate at each iteration. This bound is composed of two terms—the optimality gap, $F(x^k) - F(x^*)$, and the norm of the optimal point, $\|x^*\|$.

Lemma 8 *For a μ -strongly convex function F , we have $\|x - x^*\| \leq \frac{2}{\mu}(F(x) - F(x^*))$.*

Proof 5 *For a μ -strongly convex function, F , for all x, y we have*

$$F(x) \geq F(y) + \nabla F(y)^\top (x - y) + \frac{\mu}{2} \|x - y\|^2.$$

At the optimal point $y = x^$, we have $\nabla F(x^*) = 0$, and we obtain the desired result.*

Lemma 4 *Let Assumption 2 hold, then*

$$\|x^k\|^2 \leq \frac{4}{\mu} (F(x^k) - F(x^*)) + 2 \|x^*\|^2.$$

Proof 6 *We have*

$$\|x^k\|^2 \stackrel{\|a+b\|^2 \leq 2\|a\|^2 + 2\|b\|^2}{\leq} 2 \|x^k - x^*\|^2 + 2 \|x^*\|^2 \stackrel{\text{By Lemma 8}}{\leq} \frac{4}{\mu} (F(x^k) - F(x^*)) + 2 \|x^*\|^2.$$

Hence the result.

Recall that, inspired by the expected smoothness property Gower et al. (2019), we use a similar idea in our convergence proofs. The next lemma is a technical Lemma that helps us to prove the expected smoothness property Gower et al. (2019). The bound in Lemma 5 is composed of the optimality gap, $F(x^k) - F(x^*)$, the difference between the gradients of h at x^k and x^* , that is, $\|\nabla h(x^k) - \nabla h(x^*)\|$, and an extra constant, β , which depends on the used compressors.

Lemma 5 *Let Assumptions 1 and 2 hold, then*

$$\mathcal{A} := \mathbb{E}_{\mathcal{C}_M, \mathcal{C}} \|x^k - \mathcal{QC}_M(\bar{y}^k) - x^* + \mathcal{QC}_M(\bar{y}^*)\|^2 \leq \frac{4n^2}{\lambda^2} \|\nabla h(x^k) - \nabla h(x^*)\|^2 + \alpha (F(x^k) - F(x^*)) + \beta,$$

where $\bar{y}^* := \frac{1}{n} \sum_{j=1}^n \mathcal{C}_j(x_j^*)$, $\alpha := \frac{4(4\omega + 4\omega_M(1+\omega))}{\mu}$, and

$$\beta := 2(4\omega + 4\omega_M(1+\omega)) \|x^*\|^2 + 4\mathbb{E}_{\mathcal{C}_M, \mathcal{C}} \|\mathcal{QC}_M(\bar{y}^*) - Q\bar{x}^*\|^2.$$

Proof 7 We have

$$\begin{aligned}
\mathcal{A} &= \mathbb{E}_{\mathcal{C}_M, \mathcal{C}} \|x^k - Q\bar{x}^k + Q\bar{x}^k - QC_M(\bar{y}^k) - x^* + Q\bar{x}^* - Q\bar{x}^* + QC_M(\bar{y}^*)\|^2 \\
&= \mathbb{E}_{\mathcal{C}_M, \mathcal{C}} \|(x^k - Q\bar{x}^k - x^* + Q\bar{x}^*) + (Q\bar{x}^k - Q\bar{y}^k) + (Q\bar{y}^k - QC_M(\bar{y}^k)) + (QC_M(\bar{y}^*) - Q\bar{x}^*)\|^2 \\
&\leq 4\|x^k - Q\bar{x}^k - x^* + Q\bar{x}^*\|^2 + 4\mathbb{E}_{\mathcal{C}} \|Q\bar{x}^k - Q\bar{y}^k\|^2 + 4\mathbb{E}_{\mathcal{C}_M, \mathcal{C}} \|Q\bar{y}^k - QC_M(\bar{y}^k)\|^2 \\
&+ 4\mathbb{E}_{\mathcal{C}_M, \mathcal{C}} \|QC_M(\bar{y}^*) - Q\bar{x}^*\|^2 \\
&= 4\|x^k - Q\bar{x}^k - x^* + Q\bar{x}^*\|^2 + 4n\mathbb{E}_{\mathcal{C}} \|\bar{x}^k - \bar{y}^k\|^2 + 4n\mathbb{E}_{\mathcal{C}_M, \mathcal{C}} \|\bar{y}^k - \mathcal{C}_M(\bar{y}^k)\|^2 \\
&+ 4\mathbb{E}_{\mathcal{C}_M, \mathcal{C}} \|QC_M(\bar{y}^*) - Q\bar{x}^*\|^2 \\
&\leq 4\|x^k - Q\bar{x}^k - x^* + Q\bar{x}^*\|^2 + 4\sum_{i=1}^n \mathbb{E}_{\mathcal{C}} \|x_i^k - \mathcal{C}_i(x_i^k)\|^2 + 4n\omega_M \mathbb{E}_{\mathcal{C}_M} \|\bar{y}^k\|^2 \\
&+ 4\mathbb{E}_{\mathcal{C}_M, \mathcal{C}} \|QC_M(\bar{y}^*) - Q\bar{x}^*\|^2 \\
&\leq 4\|x^k - Q\bar{x}^k - x^* + Q\bar{x}^*\|^2 + 4\sum_{i=1}^n \omega_i \|x_i^k\|^2 + 4\omega_M \sum_{i=1}^n (1 + \omega_i) \|x_i^k\|^2 + 4\mathbb{E}_{\mathcal{C}_M, \mathcal{C}} \|QC_M(\bar{y}^*) - Q\bar{x}^*\|^2 \\
&\leq 4\frac{n^2}{\lambda^2} \|\nabla h(x^k) - \nabla h(x^*)\|^2 + (4\omega + 4\omega_M(1 + \omega)) \|x^k\|^2 + 4\mathbb{E}_{\mathcal{C}_M, \mathcal{C}} \|QC_M(\bar{y}^*) - Q\bar{x}^*\|^2 \\
&\stackrel{\text{By Lemma 4}}{\leq} 4\frac{n^2}{\lambda^2} \|\nabla h(x^k) - \nabla h(x^*)\|^2 + (4\omega + 4\omega_M(1 + \omega)) \left(\frac{4}{\mu} (F(x^k) - F(x^*)) + 2\|x^*\|^2 \right) \\
&\quad + 4\mathbb{E}_{\mathcal{C}_M, \mathcal{C}} \|QC_M(\bar{y}^*) - Q\bar{x}^*\|^2 \\
&\leq 4\frac{n^2}{\lambda^2} \|\nabla h(x^k) - \nabla h(x^*)\|^2 + \alpha (F(x^k) - F(x^*)) + \beta.
\end{aligned}$$

Hence the result.

Now we are all set to prove the expected smoothness property in our setup.

Lemma 6 (Expected Smoothness) Let Assumptions 1 and 2 hold, then

$$\mathbb{E} [\|G(x^k)\|^2 | x^k] \leq 4\gamma (F(x^k) - F(x^*)) + \delta, \quad (3)$$

where

$$\gamma := \frac{\alpha\lambda^2(1-p)}{2n^2p} + \max \left\{ \frac{L_f}{(1-p)}, \frac{\lambda}{n} \left(1 + \frac{4(1-p)}{p} \right) \right\}$$

and

$$\delta := \frac{2\beta\lambda^2(1-p)}{n^2p} + 2\mathbb{E}\|G(x^*)\|^2.$$

Proof 8 We have

$$\|G(x^k) - G(x^*)\|^2 = \begin{cases} \frac{\|\nabla f(x^k) - \nabla f(x^*)\|^2}{(1-p)^2} & \text{if } \xi_k = 0 \\ \frac{\lambda^2}{n^2p^2} \|x^k - QC_M(\bar{y}^k) - x^* + QC_M(\bar{y}^*)\|^2 & \text{if } \xi_k = 1 \text{ \& } \xi_{k-1} = 0, \\ \frac{1}{p^2} \|\nabla h(x^k) - \nabla h(x^*)\|^2 & \text{if } \xi_k = 1 \text{ \& } \xi_{k-1} = 1. \end{cases}$$

Finally,

$$\begin{aligned}
\mathbb{E}_{\xi_k, \xi_{k-1}} \|G(x^k) - G(x^*)\|^2 &= (1-p) \frac{\|\nabla f(x^k) - \nabla f(x^*)\|^2}{(1-p)^2} + p^2 \frac{1}{p^2} \|\nabla h(x^k) - \nabla h(x^*)\|^2 \\
&+ p(1-p) \frac{\lambda^2}{n^2p^2} \|x^k - QC_M(\bar{y}^k) - x^* + QC_M(\bar{y}^*)\|^2 \\
&= \frac{\|\nabla f(x^k) - \nabla f(x^*)\|^2}{(1-p)} + \|\nabla h(x^k) - \nabla h(x^*)\|^2 \\
&+ \frac{\lambda^2(1-p)}{n^2p} \|x^k - QC_M(\bar{y}^k) - x^* + QC_M(\bar{y}^*)\|^2.
\end{aligned}$$

Therefore, by using Lemma 5 we get

$$\begin{aligned}
\mathbb{E}\|G(x^k) - G(x^*)|x^k\|^2 &= \frac{\|\nabla f(x^k) - \nabla f(x^*)\|^2}{(1-p)} + \|\nabla h(x^k) - \nabla h(x^*)\|^2 + \frac{\lambda^2(1-p)}{n^2p} \mathcal{A} \\
&\leq \frac{\|\nabla f(x^k) - \nabla f(x^*)\|^2}{(1-p)} + \|\nabla h(x^k) - \nabla h(x^*)\|^2 \\
&+ \frac{\lambda^2(1-p)}{n^2p} \left(4 \frac{n^2}{\lambda^2} \|\nabla h(x^k) - \nabla h(x^*)\|^2 + \alpha(F(x^k) - F(x^*)) + \beta \right) \\
&= \frac{\|\nabla f(x^k) - \nabla f(x^*)\|^2}{(1-p)} + \left(1 + \frac{4(1-p)}{p} \right) \|\nabla h(x^k) - \nabla h(x^*)\|^2 \\
&+ \frac{\alpha\lambda^2(1-p)}{n^2p} (F(x^k) - F(x^*)) + \frac{\beta\lambda^2(1-p)}{n^2p} \\
&\leq \frac{2L_f}{(1-p)} (f(x^k) - f(x^*)) + \frac{2\lambda}{n} \left(1 + \frac{4(1-p)}{p} \right) (h(x^k) - h(x^*)) \\
&+ \frac{\alpha\lambda^2(1-p)}{n^2p} (F(x^k) - F(x^*)) + \frac{\beta\lambda^2(1-p)}{n^2p} \\
&\leq 2\gamma (F(x^k) - F(x^*)) + \frac{\beta\lambda^2(1-p)}{n^2p}.
\end{aligned}$$

Finally, we obtain

$$\begin{aligned}
\mathbb{E}\|G(x^k)|x^k\|^2 &\leq 2\mathbb{E}\|G(x^k) - G(x^*)|x^k\|^2 + 2\mathbb{E}\|G(x^*)\|^2 \\
&\leq 4\gamma (F(x^k) - F(x^*)) + \frac{2\beta\lambda^2(1-p)}{n^2p} + 2\mathbb{E}\|G(x^*)\|^2 \\
&\leq 4\gamma (F(x^k) - F(x^*)) + \delta.
\end{aligned}$$

Hence the result.

Based on the above results, the convergence of Algorithm 1 for strongly convex functions follows directly from Lemmas 3, 6 and Theorem 3.1 from Gower et al. (2019).

A.1.3 Nonconvex convergence

Theorem 5 (Non convex case) *Let Assumptions 1 and 3 hold. Assume also that F is L_f -smooth, bounded from below by $F(x^*)$. Then to reach a precision, $\epsilon > 0$, set the stepsize, $\eta = \min\{\frac{1}{L_f M}, \frac{\epsilon^2}{2L_f\sigma^2}\}$, such that for $K \geq \frac{4L_f M(F(x^0) - F(x^*))}{\epsilon^2}$, we have $\min_{k=0,1,\dots,K} \mathbb{E}\|\nabla F(x^k)\|_2 \leq \epsilon$.*

Proof 9 *From L_f -smoothness of F we have*

$$F(x^{k+1}) \leq F(x^k) - \eta_k \nabla F(x^k)^\top G(x^k) + \frac{L_f}{2} \eta_k^2 \|G(x^k)\|^2.$$

By taking the expectation in the above inequality, conditional on x^k , we get

$$\mathbb{E}[F(x^{k+1}) | x^k] \stackrel{\text{By Lemma 3}}{\leq} F(x^k) - \eta_k \|\nabla F(x^k)\|_2^2 + \frac{L_f \eta_k^2}{2} \mathbb{E}(\|G(x^k)\|^2 | x^k),$$

which by using Assumption 3 reduces to

$$\begin{aligned}
\mathbb{E}[F(x^{k+1}) | x^k] &\leq F(x^k) - \eta_k \|\nabla F(x^k)\|_2^2 + \frac{L_f \eta_k^2}{2} (M \|\nabla F(x^k)\|^2 + \sigma^2) \\
&\leq F(x^k) - \eta_k \left(1 - \frac{L_f M \eta_k}{2} \right) \|\nabla F(x^k)\|_2^2 + \frac{L_f \eta_k^2 \sigma^2}{2}.
\end{aligned}$$

After rearranging, we have

$$\eta_k \left(1 - \frac{L_f M \eta_k}{2}\right) \|\nabla F(x^k)\|_2^2 \leq F(x^k) - \mathbb{E}[F(x^{k+1}) | x^k] + \frac{L_f \eta_k^2 \sigma^2}{2}.$$

Setting $\eta_k = \eta > 0$ in the above, taking expectation, using the tower property of expectation, and finally summing over the iterates $k = 0, 1, \dots, K-1$ we have

$$\eta \left(1 - \frac{L_f M \eta}{2}\right) \sum_{k=0}^{K-1} \mathbb{E} [\|\nabla F(x^k)\|_2^2] \leq (F(x^0) - F(x^*)) + \frac{K L_f \eta \sigma^2}{2}.$$

If $\eta \leq \frac{1}{L_f M}$ then

$$\sum_{k=0}^{K-1} \mathbb{E} [\|\nabla F(x^k)\|_2^2] \leq \frac{2}{\eta} (F(x^0) - F(x^*)) + L_f K \eta \sigma^2.$$

Dividing throughout by K , we get

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} [\|\nabla F(x^k)\|_2^2] \leq \frac{2}{\eta K} (F(x^0) - F(x^*)) + L_f \eta \sigma^2.$$

Finally, setting $\eta = \frac{1}{L_f M}$ we have

$$\min_{k=0,1,\dots,K-1} \mathbb{E} \|\nabla F(x^k)\|_2^2 \leq \frac{2L_f M}{K} (F(x^0) - F(x^*)) + \frac{\sigma^2}{M}. \quad (4)$$

For a given precision, $\epsilon > 0$, to make $\min_{k=0,1,\dots,K-1} \mathbb{E} \|\nabla F(x^k)\|_2^2 \leq \epsilon^2$, we require $\frac{2L_f M (F(x^0) - F(x^*))}{K} \leq \frac{\epsilon^2}{2}$ and $L_f \eta \sigma^2 \leq \frac{\epsilon^2}{2}$, resulting in

$$K \geq \frac{4L_f M (F(x^0) - F(x^*))}{\epsilon^2} \text{ and } \eta \leq \frac{\epsilon^2}{2L_f \sigma^2}.$$

Hence the result.

A.1.4 Optimal rate and communication

The following proofs are related to optimal rate and communication as given in §6.

Theorem 1 (Optimal rate) *The probability p^* minimizing γ is equal to $\max\{p_e, p_A\}$, where $p_e = \frac{7\lambda + L - \sqrt{\lambda^2 + 14\lambda L + L^2}}{6\lambda}$ and p_A is the optimizer of the function $A(p) = \frac{\alpha\lambda^2}{2n^2p} + \frac{L}{n(1-p)}$ in $(0, 1)$.*

Proof 10 *We can rewrite γ as follows*

$$\gamma = -\frac{\alpha\lambda^2}{2n^2} + \max\{A(p), B(p)\},$$

where $A(p) = \frac{\alpha\lambda^2}{2n^2p} + \frac{L}{n(1-p)}$ and $B(p) = \frac{\alpha\lambda^2}{2n^2p} + \frac{4\lambda}{np} - \frac{3\lambda}{n}$. The function B is monotonically decreasing as a function of p . The function A goes to ∞ as p goes to zero or one, and it has one stationary point between zero and one hence it is convex in the interval $(0, 1)$. Thus it admits an optimizer p_A in $(0, 1)$. Note that $p_e = \frac{7\lambda + L - \sqrt{\lambda^2 + 14\lambda L + L^2}}{6\lambda}$ is the point for which $A(p)$ is equal to $B(p)$. Note also that near to zero $B(p) \geq A(p)$. Therefore if $p_e \leq p_A$ then the optimizer of γ is p_A otherwise it is equal to p_e . Thus the probability p^* optimizing γ is equal to $\max\{p_e, p_A\}$.

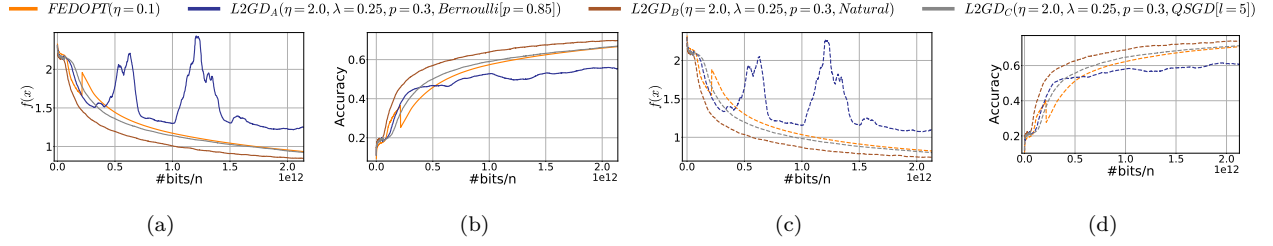


Figure 9: Training ResNet-18 on CIFAR-10, with $n = 10$ workers. Loss and Top-1 accuracy on train (a) - (b) and test data (c) - (d).

Lemma 7 The optimizer probability p_A of the function $A(p) = \frac{\alpha\lambda^2}{2n^2p} + \frac{L}{n(1-p)}$ in $(0, 1)$ is equal to

$$p_A = \begin{cases} \frac{1}{2} & \text{if } 2nL = \alpha\lambda^2 \\ \frac{-2\alpha\lambda^2 + 2\lambda\sqrt{2\alpha nL}}{2(2nL - \alpha\lambda^2)} & \text{if } 2nL > \alpha\lambda^2 \\ \frac{-2\alpha\lambda^2 - 2\lambda\sqrt{2\alpha nL}}{2(2nL - \alpha\lambda^2)} & \text{otherwise} \end{cases}$$

Proof 11 If $2nL \neq \alpha\lambda^2$, then the function A has the following two stationary points $\frac{-2\alpha\lambda^2 + 2\lambda\sqrt{2\alpha nL}}{2(2nL - \alpha\lambda^2)}$ and $\frac{-2\alpha\lambda^2 - 2\lambda\sqrt{2\alpha nL}}{2(2nL - \alpha\lambda^2)}$. If $2nL = \alpha\lambda^2$, then the function A has one stationary point equal to $\frac{1}{2}$.

Theorem 2 (Optimal communication) The probability p^* optimizing C is equal to $\max\{p_e, p_A\}$, where $p_e = \frac{7\lambda + L - \sqrt{\lambda^2 + 14\lambda L + L^2}}{6\lambda}$ and $p_A = 1 - \frac{Ln}{\alpha\lambda^2}$.

Proof 12 We can rewrite nC as follows

$$nC = \max\{A(p), B(p)\},$$

where $A(p) = \frac{\alpha\lambda^2 p(1-p)}{2n} + \frac{\alpha\lambda^2(1-p)}{2n} + Lp$ and $B(p) = \frac{\alpha\lambda^2 p(1-p)}{2n} + \frac{\alpha\lambda^2(1-p)}{2n} + 4\lambda(1-p) - 3\lambda p(1-p)$. The function B is monotonically decreasing as a function of p in $[0, 1]$. Note that $B(0) = \frac{\alpha\lambda^2}{2n} + 4\lambda$ and $B(1) = 0$. The function A admits a minimizer equal to $p_A = 1 - \frac{Ln}{\alpha\lambda^2}$. Of course p_A is a probability under the condition that $Ln \leq \alpha\lambda^2$. Thus we consider the following 2 scenarios

1. If $Ln > \alpha\lambda^2$ ($p_A < 0$) then $p^* = p_e$
2. else $p^* = \max\{p_e, p_A\}$.

We conclude in both cases that $p^* = \max\{p_e, p_A\}$.

A.2 Addendum to the Experimental Results

Batch Normalization. Beside the trainable parameters, the ResNet models contain batch normalization Ioffe & Szegedy (2015) layers that are crucial for training. The logic of batch normalization depends on the estimation of running mean and variance, and these statistics can be pretty personalized for each client in a heterogeneous data regime. The implementation of FedAvg and FedOpt in FedML considers the averaging of these statistics during the aggregation phase. In our implementation, the batch normalization statistics are included into aggregation.

Step-size. The step-sizes for FedAvg and FedOpt tuned via selecting step sizes from the following set $\{0.01, 0.1, 0.2, 0.5, 1.0, 2.0, 4.0\}$. We consider the step size for both algorithms to be 0.1. Starting with step size 0.2 algorithms diverge; we also did not use step size schedulers. Additionally, we have tuned number of local epochs for FedAvg from the following set $\{1, 2, 3, 4\}$, and pick 1 local epoch. The batch size is set to 256.

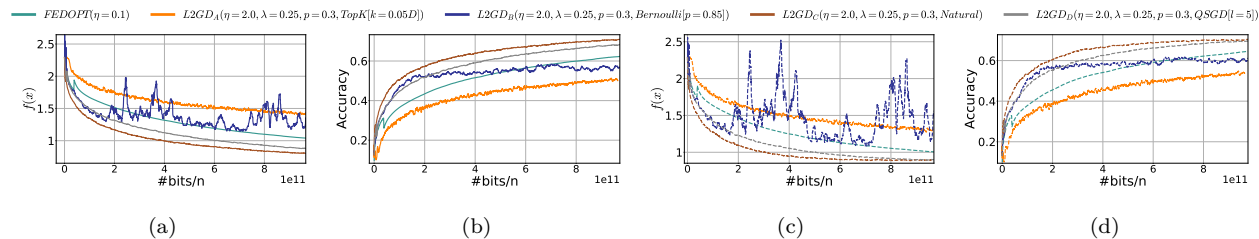


Figure 10: Training DenseNet-121 on CIFAR-10, with $n = 10$ workers. Loss and Top-1 accuracy on train (a) - (b), and test data (c) - (d).

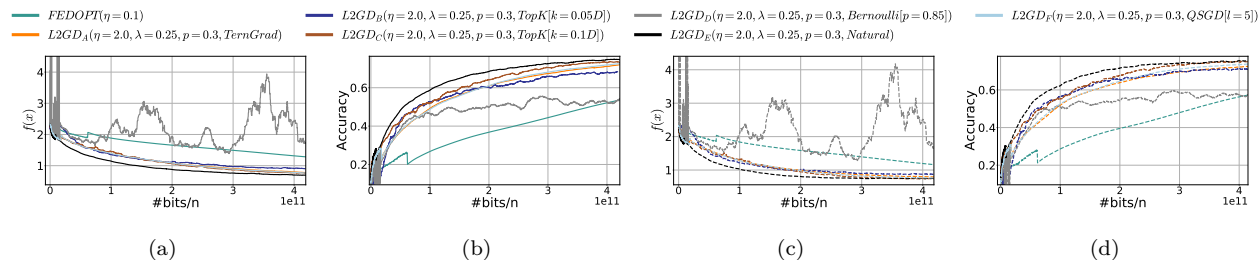


Figure 11: Training MobileNet on CIFAR-10, with $n = 10$ workers. Loss and Top-1 accuracy on train (a) - (b), and test data (c) - (d).

Compressed L2Gd vs. FedOpt. From the experiments in Section 7.2, Figures 4–6, we realized that FedAvg is not a competitive no-compression baseline for L2GD; see Table 2. FedOpt, on the other hand, remains a competitive no-compression baseline comparable to compressed L2GD. Therefore, we separately measure the performance of compressed L2GD and non-compression FedOpt for training ResNet-18, DenseNet-121, and MobileNet. Figures 9–10 demonstrate that L2GD with natural compressor (that by design has small variance) empirically behaves the best and converges approximately 5 times faster compare to FedOpt. They also show that compressed L2GD with natural compressor sends the least data and drives the loss down the most. At the same time, L2GD with natural compressor reaches the best accuracy for both train and test sets.

Reproducible research. See our repository online: <https://github.com/burlachenkok/compressed-f1-l2gd-code>. Our source codes have been constructed on top of the popular federated learning repository: FedML.ai; see <https://github.com/FedML-AI/FedML/commit/3b9b68764d922ce239e0b84aceda986cfa977f96>.