
LoCoDL: COMMUNICATION-EFFICIENT DISTRIBUTED LEARNING WITH LOCAL TRAINING AND COMPRESSION

Anonymous authors

Paper under double-blind review

ABSTRACT

In Distributed optimization and Learning, and even more in the modern framework of federated learning, communication, which is slow and costly, is critical. We introduce LoCoDL, a communication-efficient algorithm that leverages the two popular and effective techniques of Local training, which reduces the communication frequency, and Compression, in which short bitstreams are sent instead of full-dimensional vectors of floats. LoCoDL works with a large class of unbiased compressors that includes widely-used sparsification and quantization methods. LoCoDL provably benefits from local training and compression and enjoys a doubly-accelerated communication complexity, with respect to the condition number of the functions and the model dimension, in the general heterogenous regime with strongly convex functions. This is confirmed in practice, with LoCoDL outperforming existing algorithms.

1 INTRODUCTION

Performing distributed computations is now pervasive in all areas of science. Notably, Federated Learning (FL) consists in training machine learning models in a distributed and collaborative way (Konečný et al., 2016a;b; McMahan et al., 2017; Bonawitz et al., 2017). The key idea in this rapidly growing field is to exploit the wealth of information stored on distant devices, such as mobile phones or hospital workstations. The many challenges to face in FL include data privacy and robustness to adversarial attacks, but communication-efficiency is likely to be the most critical (Kairouz et al., 2021; Li et al., 2020a; Wang et al., 2021). Indeed, in contrast to the centralized setting in a datacenter, in FL the clients perform parallel computations but also communicate back and forth with a distant orchestrating server. Communication typically takes place over the internet or cell phone network, and can be slow, costly, and unreliable. It is the main bottleneck that currently prevents large-scale deployment of FL in mass-market applications.

Two strategies to reduce the communication burden have been popularized by the pressing needs of FL: 1) **Local Training (LT)**, which consists in reducing the communication frequency. That is, instead of communicating the output of every computation step involving a (stochastic) gradient call, several such steps are performed between successive communication rounds. 2) **Communication Compression (CC)**, in which compressed information is sent instead of full-dimensional vectors. We review the literature of LT and CC in Section 1.2.

We propose a new randomized algorithm named LoCoDL, which features LT and unbiased CC for communication-efficient FL and distributed optimization. It is variance-reduced (Hanzely & Richtárik, 2019; Gorbunov et al., 2020a; Gower et al., 2020), so that it converges to an exact solution. It provably benefits from the two mechanisms of LT and CC: the communication complexity is doubly accelerated, with a better dependency on the condition number of the functions and on the dimension of the model.

1.1 PROBLEM AND MOTIVATION

We study distributed optimization problems of the form

$$\min_{x \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n f_i(x) + g(x), \quad (1)$$

where $d \geq 1$ is the model dimension and the functions $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$ and $g : \mathbb{R}^d \rightarrow \mathbb{R}$ are smooth. We consider the server-client model in which $n \geq 1$ clients do computations in parallel and communicate back and forth with a server. The private function f_i is owned by and stored on client $i \in [n] := \{1, \dots, n\}$. Problem (1) models empirical risk minimization, of utmost importance in machine learning (Sra et al., 2011; Shalev-Shwartz & Ben-David, 2014). More generally, minimizing a sum of functions appears in virtually all areas of science and engineering. Our goal is to solve Problem (1) in a communication-efficient way, in the general **heterogeneous** setting in which the functions f_i , as well as g , can be *arbitrarily different*: we do not make any assumption on their similarity whatsoever.

We consider in this work the strongly convex setting — an analysis with nonconvex functions would certainly require very different proof techniques, which we currently do not know how to derive. That is, the following holds:

Assumption 1.1 (strongly convex functions). The functions f_i and g are all L -smooth and μ -strongly convex, for some $0 < \mu \leq L$.¹ Then we denote by x^* the solution of the strongly convex problem (1), which exists and is unique. We define the condition number $\kappa := \frac{L}{\mu}$.

Problem (1) can be viewed as the minimization of the average of the n functions $(f_i + g)$, which can be performed using calls to $\nabla(f_i + g) = \nabla f_i + \nabla g$. We do not use this straightforward interpretation. Instead, let us illustrate the interest of having the **additional function** g in (1), using 4 different viewpoints. We stress that we can handle the case $g = 0$, as discussed in Section 3.1.

- **Viewpoint 1: regularization.** The function g can be a regularizer. For instance, if the functions f_i are convex, adding $g = \frac{\mu}{2} \|\cdot\|^2$ for a small $\mu > 0$ makes the problem μ -strongly convex.
- **Viewpoint 2: shared dataset.** The function g can model the cost of a common dataset, or a piece thereof, that is known to all clients.
- **Viewpoint 3: server-aided training.** The function g can model the cost of a core dataset, known only to the server, which makes calls to ∇g . This setting has been investigated in several works, with the idea that using a small auxiliary dataset representative of the global data distribution, the server can correct for the deviation induced by partial participation (Zhao et al., 2018; Yang et al., 2021; 2024). We do not focus on this setting, because we deal with the general heterogeneous setting in which g and the f_i are not meant to be similar in any sense, and in our work g is handled by the clients, not by the server.
- **Viewpoint 4: a new mathematical and algorithmic principle.** This is the idea that led to the construction of **LoCoDL**, and we detail it in Section 2.1.

In **LoCoDL**, the clients make all gradient calls; that is, Client i makes calls to ∇f_i and ∇g .

1.2 STATE OF THE ART

We review the latest developments on communication-efficient algorithms for distributed learning, making use of LT, CC, or both. Before that, we note that we should distinguish uplink, or clients-to-server, from downlink, or server-to-clients, communication. Uplink is usually slower than downlink communication, since the clients uploading *different* messages in parallel to the server is slower than the clients downloading *the same* message in parallel from the server. This can be due to cache memory and aggregation speed constraints of the server, as well as asymmetry of the service provider’s systems or protocols used on the internet or cell phone network. In this work, we focus on the **uplink communication complexity**, which is often the bottleneck in practice. Indeed, the goal is to exploit parallelism to obtain better performance when n increases. Precisely, with **LoCoDL**, the uplink communication complexity decreases from $\mathcal{O}(d\sqrt{\kappa} \log \epsilon^{-1})$ when n is small to $\mathcal{O}(\sqrt{d}\sqrt{\kappa} \log \epsilon^{-1})$ when n is large, where the condition number κ is defined in Assumption 1.1, see Corollary 3.2. Many works have considered bidirectional compression, which consists in compressing the messages sent both ways (Gorbunov et al., 2020b; Philippenko & Dieuleveut, 2020; Liu et al., 2020; Philippenko & Dieuleveut, 2021; Condat & Richtárik, 2022; Gruntkowska et al., 2023; Tyurin

¹A differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is said to be L -smooth if ∇f is L -Lipschitz continuous; that is, for every $x \in \mathbb{R}^d$ and $y \in \mathbb{R}^d$, $\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|$ (the norm is the Euclidean norm throughout the paper). f is said to be μ -strongly convex if $f - \frac{\mu}{2} \|\cdot\|^2$ is convex.

& Richtárik, 2023b) but to the best of our knowledge, this has no impact on the downlink complexity, which cannot be reduced further than $\mathcal{O}(d\sqrt{\kappa} \log \epsilon^{-1})$, just because there is no parallelism to exploit in this direction. Thus, we focus our analysis on theoretical and algorithmic techniques to reduce the uplink communication complexity, which we call communication complexity in short, and we ignore downlink communication.

Communication Compression (CC) consists in applying some lossy scheme that compresses vectors into messages of small bit size, which are communicated. For instance, the well-known `rand-k` compressor selects k coordinates of the vector uniformly at random, for some $k \in [d] := \{1, \dots, d\}$. k can be as small as 1, in which case the compression factor is d , which can be huge. Some compressors, such as `rand-k`, are unbiased, whereas others are biased; we refer to Beznosikov et al. (2020); Albasyoni et al. (2020); Horváth et al. (2022); Condat et al. (2022b) for several examples and a discussion of their properties. The introduction of **DIANA** by Mishchenko et al. (2019) was a major milestone, as this algorithm converges linearly with the large class of unbiased compressors defined in Section 1.3 and also considered in **LoCoDL**. The communication complexity $\mathcal{O}(d\kappa \log \epsilon^{-1})$ of the basic Gradient Descent (**GD**) algorithm is reduced with **DIANA** to $\mathcal{O}((\kappa + d) \log \epsilon^{-1})$ when n is large, see Table 2. **DIANA** was later extended in several ways (Horváth et al., 2022; Gorbunov et al., 2020a; Condat & Richtárik, 2022). An accelerated version of **DIANA** called **ADIANA** based on Nesterov Accelerated GD has been proposed (Li et al., 2020b) and further analyzed in He et al. (2023); it has the state-of-the-art theoretical complexity.

Algorithms converging linearly with biased compressors have also been proposed, such as **EF21** (Richtárik et al., 2021; Fatkhullin et al., 2021; Condat et al., 2022b), but the acceleration potential is less understood than with unbiased compressors. Algorithms with CC such as **MARINA** (Gorbunov et al., 2021) and **DASHA** (Tyurin & Richtárik, 2023a) have been proposed for nonconvex optimization, but their analysis requires a different approach and there is a gap in the achievable performance: their complexity depends on $\frac{\omega\kappa}{\sqrt{n}}$ instead of $\frac{\omega\kappa}{n}$ with **DIANA**, where ω characterizes the compression error variance, see (2). Therefore, we focus on the convex setting and leave the nonconvex study for future work.

Local Training (LT) is a simple but remarkably efficient idea: the clients perform multiple Gradient Descent (**GD**) steps, instead of only one, between successive communication rounds. The intuition behind is that this leads to the communication of richer information, so that the number of communication rounds to reach a given accuracy is reduced. We refer to Mishchenko et al. (2022) for a comprehensive review of LT-based algorithms, which include the popular **FedAvg** and **Scaffold** algorithms of McMahan et al. (2017) and Karimireddy et al. (2020), respectively. Mishchenko et al. (2022) made a breakthrough by proposing **Scaffnew**, the first LT-based variance-reduced algorithm that not only converges linearly to the exact solution in the strongly convex setting, but does so with accelerated communication complexity $\mathcal{O}(d\sqrt{\kappa} \log \epsilon^{-1})$. In **Scaffnew**, communication can occur randomly after every iteration, but occurs only with a small probability p . Thus, there are in average p^{-1} local steps between successive communication rounds. The optimal dependency on $\sqrt{\kappa}$ (Scaman et al., 2019) is obtained with $p = 1/\sqrt{\kappa}$. **LoCoDL** has the same probabilistic LT mechanism as **Scaffnew** but does not revert to it when compression is disabled, because of the additional function g and tracking variables y and v . A different approach to LT was developed by Sadiev et al. (2022a) with the **APDA-Inexact** algorithm, and generalized to handle partial participation by Grudzień et al. (2023) with the **5GCS** algorithm: in both algorithms, the local GD steps form an inner loop in order to compute a proximity operator inexactly.

Combining LT and CC while retaining their benefits is very challenging. In our strongly convex and heterogeneous setting, the methods **Qsparse-local-SGD** (Basu et al., 2020) and **FedPAQ** (Reisizadeh et al., 2020) do not converge linearly. **FedCOMGATE** features LT + CC and converges linearly (Haddadpour et al., 2021), but its complexity $\mathcal{O}(d\kappa \log \epsilon^{-1})$ does not show any acceleration. We can mention that random reshuffling, a technique that can be seen as a type of LT, has been combined with CC in Sadiev et al. (2022b); Malinovsky & Richtárik (2022). Recently, Condat et al. (2022a) managed to design a specific compression technique compatible with the LT mechanism of **Scaffnew**, leading to **CompressedScaffnew**, the first LT + CC algorithm exhibiting a doubly-accelerated complexity, namely $\mathcal{O}((\sqrt{d}\sqrt{\kappa} + \frac{d\sqrt{\kappa}}{\sqrt{n}} + d) \log \epsilon^{-1})$, as reported in Table 2. However, **CompressedScaffnew** uses a specific linear compression scheme that requires shared randomness; that is, all clients have to agree on a random permutation of the columns of the global compression pattern. No other compressor can be used, which notably rules out any type of quantization.

Table 1: Communication complexity in number of communication rounds to reach ϵ -accuracy for linearly-converging algorithms allowing for CC with independent compressors in $\mathbb{U}(\omega)$ for any $\omega \geq 0$. Since the compressors are independent, $\omega_{\text{av}} = \frac{\omega}{n}$. We provide the leading asymptotic factor and ignore log factors such as $\log \epsilon^{-1}$. The state of the art is highlighted in green.

Algorithm	Com. complexity in # rounds	case $\omega = \mathcal{O}(n)$	case $\omega = \Theta(n)$
DIANA	$(1 + \frac{\omega}{n})\kappa + \omega$	$\kappa + \omega$	$\kappa + \omega$
EF21	$(1 + \omega)\kappa$	$(1 + \omega)\kappa$	$(1 + \omega)\kappa$
5GCS-CC	$(1 + \sqrt{\omega} + \frac{\omega}{\sqrt{n}})\sqrt{\kappa} + \omega$	$(1 + \sqrt{\omega})\sqrt{\kappa} + \omega$	$(1 + \sqrt{\omega})\sqrt{\kappa} + \omega$
ADIANA ¹	$(1 + \frac{\omega^{3/4}}{n^{1/4}} + \frac{\omega}{\sqrt{n}})\sqrt{\kappa} + \omega$	$(1 + \frac{\omega^{3/4}}{n^{1/4}})\sqrt{\kappa} + \omega$	$(1 + \sqrt{\omega})\sqrt{\kappa} + \omega$
ADIANA ²	$(1 + \frac{\omega}{\sqrt{n}})\sqrt{\kappa} + \omega$	$(1 + \frac{\omega}{\sqrt{n}})\sqrt{\kappa} + \omega$	$(1 + \sqrt{\omega})\sqrt{\kappa} + \omega$
lower bound ²	$(1 + \frac{\omega}{\sqrt{n}})\sqrt{\kappa} + \omega$	$(1 + \frac{\omega}{\sqrt{n}})\sqrt{\kappa} + \omega$	$(1 + \sqrt{\omega})\sqrt{\kappa} + \omega$
LoCoDL	$(1 + \sqrt{\omega} + \frac{\omega}{\sqrt{n}})\sqrt{\kappa} + \omega(1 + \frac{\omega}{n})$	$(1 + \sqrt{\omega})\sqrt{\kappa} + \omega$	$(1 + \sqrt{\omega})\sqrt{\kappa} + \omega$

¹This is the complexity derived in the original paper Li et al. (2020b).

²This is the complexity derived by a refined analysis in the preprint He et al. (2023), where a matching lower bound is also derived.

1.3 A GENERAL CLASS OF UNBIASED RANDOM COMPRESSORS

For every $\omega \geq 0$, we define $\mathbb{U}(\omega)$ as the set of random compression operators $\mathcal{C} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ that are unbiased, i.e. $\mathbb{E}[\mathcal{C}(x)] = x$, and satisfy, for every $x \in \mathbb{R}^d$,

$$\mathbb{E}[\|\mathcal{C}(x) - x\|^2] \leq \omega \|x\|^2. \quad (2)$$

In addition, given a collection $(\mathcal{C}_i)_{i=1}^n$ of compression operators in $\mathbb{U}(\omega)$ for some $\omega \geq 0$, in order to characterize their joint variance, we introduce the constant $\omega_{\text{av}} \geq 0$ such that, for every $x_i \in \mathbb{R}^d$, $i \in [n]$, we have

$$\mathbb{E} \left[\left\| \frac{1}{n} \sum_{i=1}^n (\mathcal{C}_i(x_i) - x_i) \right\|^2 \right] \leq \frac{\omega_{\text{av}}}{n} \sum_{i=1}^n \|x_i\|^2. \quad (3)$$

The inequality (3) is not an additional assumption: it is satisfied with $\omega_{\text{av}} = \omega$ by convexity of the squared norm. But the convergence rate will depend on ω_{av} , which is typically much smaller than ω . In particular, if the compressors \mathcal{C}_i are mutually independent, the variance of their sum is the sum of their variances, and (3) is satisfied with $\omega_{\text{av}} = \frac{\omega}{n}$.

1.4 CHALLENGE AND CONTRIBUTIONS

This work addresses the following question: *Can we combine LT and CC with any compressors in the generic class $\mathbb{U}(\omega)$ defined in the previous section, and fully benefit from both techniques by obtaining a doubly-accelerated communication complexity?*

We answer this question in the affirmative. **LoCoDL** has the same probabilistic LT mechanism as **Scaffnew** and features CC with compressors in $\mathbb{U}(\omega)$ with arbitrarily large $\omega \geq 0$, with proved linear convergence under Assumption 1.1, without further requirements. By choosing the communication probability and the variance ω appropriately, double acceleration is obtained. Thus, **LoCoDL** achieves the same theoretical complexity as **CompressedScaffnew**, but allows for a large class of compressors instead of the cumbersome permutation-based compressor of the latter. In particular, with compressors performing sparsification and quantization, **LoCoDL** outperforms existing algorithms, as we show by experiments in Section 4. This is remarkable, since **ADIANA**, based on Nesterov acceleration and not LT, has an even better theoretical complexity when n is larger than d , see Table 2, but this is not reflected in practice: **ADIANA** is clearly behind **LoCoDL** in our experiments. Thus, our experiments indicate that **LoCoDL** sets new standards in terms of communication efficiency.

Table 2: (Uplink) communication complexity in number of reals to reach ϵ -accuracy for linearly-converging algorithms allowing for CC, with an optimal choice of unbiased compressors. We provide the leading asymptotic factor and ignore log factors such as $\log \epsilon^{-1}$. The state of the art is highlighted in green.

Algorithm	complexity in # reals	case $n = \mathcal{O}(d)$
DIANA	$(1 + \frac{d}{n})\kappa + d$	$\frac{d}{n}\kappa + d$
EF21	$\frac{d\kappa}{d\kappa}$	$\frac{d\kappa}{d\kappa}$
5GCS-CC	$(\sqrt{d} + \frac{d}{\sqrt{n}})\sqrt{\kappa} + d$	$\frac{d}{\sqrt{n}}\sqrt{\kappa} + d$
ADIANA	$(1 + \frac{d}{\sqrt{n}})\sqrt{\kappa} + d$	$\frac{d}{\sqrt{n}}\sqrt{\kappa} + d$
CompressedScaffnew	$(\sqrt{d} + \frac{d}{\sqrt{n}})\sqrt{\kappa} + d$	$\frac{d}{\sqrt{n}}\sqrt{\kappa} + d$
FedCOMGATE	$\frac{d\kappa}{d\kappa}$	$\frac{d\kappa}{d\kappa}$
LoCoDL	$(\sqrt{d} + \frac{d}{\sqrt{n}})\sqrt{\kappa} + d$	$\frac{d}{\sqrt{n}}\sqrt{\kappa} + d$

2 PROPOSED ALGORITHM LoCoDL

2.1 PRINCIPLE: DOUBLE LIFTING OF THE PROBLEM TO A CONSENSUS PROBLEM

In **LoCoDL**, every client stores and updates *two* local model estimates. They will all converge to the same solution x^* of (1). This construction comes from two ideas.

Local steps with local models. In algorithms making use of LT, such as **FedAvg**, **Scaffold** and **Scaffnew**, the clients store and update local model estimates x_i . When communication occurs, an estimate of their average is formed by the server and broadcast to all clients. They all resume their computations with this new model estimate.

Compressing the difference between two estimates. To implement CC, a powerful idea is to compress not the vectors themselves, but *difference vectors* that converge to zero. This way, the algorithm is variance-reduced; that is, the compression error vanishes at convergence. The technique of compressing the difference between a gradient vector and a control variate is at the core of algorithms such as **DIANA** and **EF21**. Here, we want to compress differences between model estimates, not gradient estimates. That is, we want Client i to compress the difference between x_i and another model estimate that converges to the solution x^* as well. We see the need of an additional model estimate that plays the role of an anchor for compression. This is the variable y common to all clients in **LoCoDL**, which compress $x_i - y$ and send these compressed differences to the server.

Combining the two ideas. Accordingly, an equivalent reformulation of (1) is the consensus problem with $n + 1$ variables

$$\min_{x_1, \dots, x_n, y} \frac{1}{n} \sum_{i=1}^n f_i(x_i) + g(y) \quad \text{s.t.} \quad x_1 = \dots = x_n = y.$$

The primal–dual optimality conditions are $x_1 = \dots = x_n = y$, $0 = \nabla f_i(x_i) - u_i \forall i \in [n]$, $0 = \nabla g(y) - v$, and $0 = u_1 + \dots + u_n + nv$ (dual feasibility), for some dual variables u_1, \dots, u_n, v introduced in **LoCoDL**, that always satisfy the dual feasibility condition.

2.2 DESCRIPTION OF LoCoDL

LoCoDL is a randomized primal–dual algorithm, shown as Algorithm 1. At every iteration, for every $i \in [n]$ in parallel, Client i first constructs a prediction \hat{x}_i^t of its updated local model estimate, using a GD step with respect to f_i corrected by the dual variable u_i^t . It also constructs a prediction \hat{y}^t of the updated model estimate, using a GD step with respect to g corrected by the dual variable v^t . Since g is known by all clients, they all maintain and update identical copies of the variables y and v . If there is no communication, which is the case with probability $1 - p$, x_i and y are updated with these predicted estimates, and the dual variables u_i and v are unchanged. If communication occurs, which is the case with probability p , the clients compress the differences $\hat{x}_i^t - \hat{y}^t$ and send these compressed vectors to the server, which forms \bar{d}^t equal to one half of their average. Then the

Algorithm 1 LoCoDL

1: **input:** stepsizes $\gamma > 0, \chi > 0, \rho > 0$; probability $p \in (0, 1]$; variance factor $\omega \geq 0$; local initial estimates $x_1^0, \dots, x_n^0 \in \mathbb{R}^d$, initial estimate $y^0 \in \mathbb{R}^d$, initial control variates $u_1^0, \dots, u_n^0 \in \mathbb{R}^d$ and $v \in \mathbb{R}^d$ such that $\frac{1}{n} \sum_{i=1}^n u_i^0 + v^0 = 0$.
 2: **for** $t = 0, 1, \dots$ **do**
 3: **for** $i = 1, \dots, n$, at clients in parallel, **do**
 4: $\hat{x}_i^t := x_i^t - \gamma \nabla f_i(x_i^t) + \gamma u_i^t$
 5: $\hat{y}^t := y^t - \gamma \nabla g(y^t) + \gamma v^t$ // the clients store and update identical copies of y^t, v^t, \hat{y}^t
 6: flip a coin $\theta^t \in \{0, 1\}$ with $\text{Prob}(\theta^t = 1) = p$
 7: **if** $\theta^t = 1$ **then**
 8: $d_i^t := \mathcal{C}_i^t(\hat{x}_i^t - \hat{y}^t)$
 9: send d_i^t to the server
 10: at server: aggregate $\bar{d}^t := \frac{1}{2n} \sum_{j=1}^n d_j^t$ and broadcast \bar{d}^t to all clients
 11: $x_i^{t+1} := (1 - \rho)\hat{x}_i^t + \rho(\hat{y}^t + \bar{d}^t)$
 12: $u_i^{t+1} := u_i^t + \frac{p\chi}{\gamma(1+2\omega)}(\bar{d}^t - d_i^t)$
 13: $y^{t+1} := \hat{y}^t + \rho\bar{d}^t$
 14: $v^{t+1} := v^t + \frac{p\chi}{\gamma(1+2\omega)}\bar{d}^t$
 15: **else**
 16: $x_i^{t+1} := \hat{x}_i^t, y^{t+1} = \hat{y}^t, u_i^{t+1} := u_i^t, v^{t+1} := v^t$
 17: **end if**
 18: **end for**
 19: **end for**

variables x_i are updated using a convex combination of the local predicted estimates \hat{x}_i^t and the global but noisy estimate $\hat{y}^t + \bar{d}^t$. y is updated similarly. Finally, the dual variables are updated using the compressed differences minus their weighted average, so that the dual feasibility condition remains satisfied. The model estimates $x_i^t, \hat{x}_i^t, y^t, \hat{y}^t$ all converge to x^* , so that their differences, as well as the compressed differences as a consequence of (2), converge to zero. This is the key property that makes the algorithm variance-reduced. We consider the following assumption.

Assumption 2.1 (class of compressors). In **LoCoDL** the compressors \mathcal{C}_i^t are all in $\mathbb{U}(\omega)$ for some $\omega \geq 0$. Moreover, for every $i \in [n], i' \in [n], t \geq 0, t' \geq 0$, \mathcal{C}_i^t and $\mathcal{C}_{i'}^{t'}$ are independent if $t \neq t'$ (\mathcal{C}_i^t and $\mathcal{C}_{i'}^t$ at the same iteration t need not be independent). We define $\omega_{\text{av}} \geq 0$ such that for every $t \geq 0$, the collection $(\mathcal{C}_i^t)_{i=1}^n$ satisfies (3).

Remark 2.2 (partial participation). **LoCoDL** allows for a form of partial participation if we set $\rho = 1$. Indeed, in that case, at steps 11 and 13 of the algorithm, all local variables x_i as well as the common variable y are overwritten by the same up-to-date model $\hat{y}^t + \bar{d}^t$. So, it does not matter that for a non-participating client i with $d_i^t = 0$, the $\hat{x}_i^{t'}$ were not computed for the $t' \leq t$ since its last participation, as they are not used in the process. However, a non-participating client should still update its local copy of y at every iteration. This can be done when ∇g is much cheaper to compute than ∇f_i , as is the case with $g = \frac{L}{2} \|\cdot\|^2$. A non-participating client can be completely idle for a certain period of time, but when it resumes participating, it should receive the last estimates of x, y and v from the server as it lost synchronization.

3 CONVERGENCE AND COMPLEXITY OF **LoCoDL**

Theorem 3.1 (linear convergence of **LoCoDL**). Suppose that Assumptions 1.1 and 2.1 hold. In **LoCoDL**, suppose that $0 < \gamma < \frac{2}{L}, 2\rho - \rho^2(1 + \omega_{\text{av}}) - \chi \geq 0$. For every $t \geq 0$, define the Lyapunov function

$$\Psi^t := \frac{1}{\gamma} \left(\sum_{i=1}^n \|x_i^t - x^*\|^2 + n \|y^t - x^*\|^2 \right) + \frac{\gamma(1+2\omega)}{p^2\chi} \left(\sum_{i=1}^n \|u_i^t - u_i^*\|^2 + n \|v^t - v^*\|^2 \right), \quad (4)$$

where $v^* := \nabla g(x^*)$ and $u_i^* := \nabla f_i(x^*)$. Then **LoCoDL** converges linearly: for every $t \geq 0$,

$$\mathbb{E}[\Psi^t] \leq \tau^t \Psi^0, \quad \text{where } \tau := \max \left((1 - \gamma\mu)^2, (1 - \gamma L)^2, 1 - \frac{p^2\chi}{1 + 2\omega} \right) < 1. \quad (5)$$

In addition, for every $i \in [n]$, $(x_i^t)_{t \in \mathbb{N}}$ and $(y^t)_{t \in \mathbb{N}}$ converge to x^* , $(u_i^t)_{t \in \mathbb{N}}$ converges to u_i^* , and $(v^t)_{t \in \mathbb{N}}$ converges to v^* , almost surely.

We place ourselves in the conditions of Theorem 3.1. We observe that in (5), the larger χ , the better, so given ρ we should set $\chi = 2\rho - \rho^2(1 + \omega_{\text{av}})$. Then, choosing ρ to maximize χ yields

$$\chi = \rho = \frac{1}{1 + \omega_{\text{av}}}. \quad (6)$$

We now study the complexity of **LoCoDL** with χ and ρ chosen as in (6) and $\gamma = \Theta(\frac{1}{L})$. We remark that **LoCoDL** has the same rate $\tau^\sharp := \max(1 - \gamma\mu, \gamma L - 1)^2$ as mere distributed gradient descent, as long as p^{-1} , ω and ω_{av} are small enough to have $1 - \frac{p^2\chi}{1+2\omega} \leq \tau^\sharp$. This is remarkable: communicating with a low frequency and compressed vectors does not harm convergence at all, until some threshold.

The iteration complexity of **LoCoDL** to reach ϵ -accuracy, i.e. $\mathbb{E}[\Psi^t] \leq \epsilon\Psi^0$, is

$$\mathcal{O}\left(\left(\kappa + \frac{(1 + \omega_{\text{av}})(1 + \omega)}{p^2}\right) \log \epsilon^{-1}\right). \quad (7)$$

By choosing

$$p = \min\left(\sqrt{\frac{(1 + \omega_{\text{av}})(1 + \omega)}{\kappa}}, 1\right), \quad (8)$$

the iteration complexity becomes $\mathcal{O}\left((\kappa + \omega(1 + \omega_{\text{av}})) \log \epsilon^{-1}\right)$ and the communication complexity in number of communication rounds is p times the iteration complexity, that is

$$\mathcal{O}\left(\left(\sqrt{\kappa(1 + \omega_{\text{av}})(1 + \omega)} + \omega(1 + \omega_{\text{av}})\right) \log \epsilon^{-1}\right).$$

If the compressors are mutually independent, $\omega_{\text{av}} = \frac{\omega}{n}$ and the communication complexity can be equivalently written as

$$\mathcal{O}\left(\left(\left(1 + \sqrt{\omega} + \frac{\omega}{\sqrt{n}}\right) \sqrt{\kappa} + \omega \left(1 + \frac{\omega}{n}\right)\right) \log \epsilon^{-1}\right),$$

as shown in Table 1.

Let us consider the example of independent **rand- k** compressors, for some $k \in [d]$. We have $\omega = \frac{d}{k} - 1$. Therefore, the communication complexity in numbers of reals is k times the complexity in number of rounds; that is, $\mathcal{O}\left(\left(\left(\sqrt{\kappa d} + \frac{d}{\sqrt{n}}\right) \sqrt{\kappa} + d \left(1 + \frac{d}{kn}\right)\right) \log \epsilon^{-1}\right)$. We can now choose k to minimize this complexity: with $k = \lceil \frac{d}{n} \rceil$, it becomes $\mathcal{O}\left(\left(\left(\sqrt{d} + \frac{d}{\sqrt{n}}\right) \sqrt{\kappa} + d\right) \log \epsilon^{-1}\right)$, as shown in Table 2. Let us state this result:

Corollary 3.2. *In the conditions of Theorem 3.1, suppose in addition that the compressors \mathcal{C}_i^t are independent **rand- k** compressors with $k = \lceil \frac{d}{n} \rceil$. Suppose that $\gamma = \Theta(\frac{1}{L})$, $\chi = \rho = \frac{n}{n-1+d/k}$, and*

$$p = \min\left(\sqrt{\frac{dk(n-1) + d^2}{nk^2\kappa}}, 1\right). \quad (9)$$

Then the uplink communication complexity in number of reals of **LoCoDL** is

$$\mathcal{O}\left(\left(\sqrt{d}\sqrt{\kappa} + \frac{d\sqrt{\kappa}}{\sqrt{n}} + d\right) \log \epsilon^{-1}\right). \quad (10)$$

This is the same complexity as **CompressedScaffnew** (Condat et al., 2022a). However, it is obtained with simple independent compressors, which is much more practical than the permutation-based compressors with shared randomness of **CompressedScaffnew**. Moreover, this complexity can be obtained with other types of compressors, and further reduced, when reasoning in number of bits and not only reals, by making use of quantization (Albasyoni et al., 2020), as we illustrate by experiments in the next section.

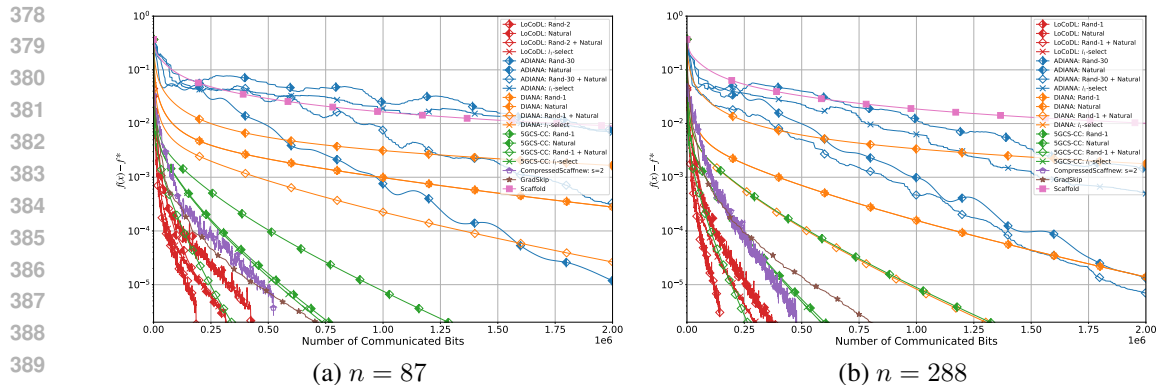


Figure 1: Comparison of several algorithms with several compressors on logistic regression with the ‘a5a’ dataset from the LibSVM, which has $d = 122$ and 6,414 data points. We chose different values of n to illustrate the two regimes $n < d$ and $n > d$, as discussed at the end of Section 3.

We can distinguish 2 regimes:

1. In the “large d small n ” regime, i.e. $n = \mathcal{O}(d)$, the communication complexity of **LoCoDL** in (10) becomes $\mathcal{O}\left(\left(\frac{d\sqrt{\kappa}}{\sqrt{n}} + d\right) \log \epsilon^{-1}\right)$. This is the state of the art, as reported in Table 2.

2. In the “large n small d ” regime, i.e. $n = \Omega(d)$, the communication complexity of **LoCoDL** in (10) becomes $\mathcal{O}\left(\left(\sqrt{d}\sqrt{\kappa} + d\right) \log \epsilon^{-1}\right)$. If n is even larger with $n = \Omega(d^2)$, **ADIANA** achieves the even better complexity $\mathcal{O}\left(\left(\sqrt{\kappa} + d\right) \log \epsilon^{-1}\right)$.

Yet, in the experiments we ran with different datasets and values of d, n, κ , **LoCoDL** outperforms the other algorithms, including **ADIANA**, in all cases.

3.1 THE CASE $g = 0$

We have assumed the presence of a function g in Problem (1), whose gradient is called by all clients. In this section, we show that we can handle the case where such a function is not available. So, let us assume that we want to minimize $\frac{1}{n} \sum_{i=1}^n f_i$, with the functions f_i satisfying Assumption 1.1. We now define the functions $\tilde{f}_i := f_i - \frac{\mu}{4} \|\cdot\|^2$ and $\tilde{g} := \frac{\mu}{4} \|\cdot\|^2$. They are all \tilde{L} -smooth and $\tilde{\mu}$ -strongly convex, with $\tilde{L} := L - \frac{\mu}{2}$ and $\tilde{\mu} := \frac{\mu}{2}$. Moreover, it is equivalent to minimize $\frac{1}{n} \sum_{i=1}^n \tilde{f}_i$ or $\frac{1}{n} \sum_{i=1}^n \tilde{f}_i + \tilde{g}$. We can then apply **LoCoDL** to the latter problem. At Step 5, we simply have $y^t - \gamma \nabla \tilde{g}(y^t) = (1 - \frac{\gamma\mu}{2})y^t$. The rate in (5) applies with L and μ replaced by \tilde{L} and $\tilde{\mu}$, respectively. Since $\kappa \leq \tilde{\kappa} := \frac{\tilde{L}}{\tilde{\mu}} \leq 2\kappa$, the asymptotic complexities derived above also apply to this setting. Thus, the presence of g in Problem (1) is not restrictive at all, as the only property of g that matters is that it has the same amount of strong convexity as the f_i s.

4 EXPERIMENTS

We evaluate the performance of our proposed method **LoCoDL** and compare it with several other methods that also allow for CC and converge linearly to x^* . We also include **GradSkip** (Maranjyan et al., 2022) and **Scaffold** (McMahan et al., 2017) in our comparisons. We focus on a regularized logistic regression problem, which has the form (1) with

$$f_i(x) = \frac{1}{m} \sum_{s=1}^m \log\left(1 + \exp\left(-b_{i,s} a_{i,s}^\top x\right)\right) + \frac{\mu}{2} \|x\|^2 \quad (11)$$

and $g = \frac{\mu}{2} \|x\|^2$, where n is the number of clients, m is the number of data points per client, $a_{i,s} \in \mathbb{R}^d$ and $b_{i,s} \in \{-1, +1\}$ are the data samples, and μ is the regularization parameter, set so that $\kappa = 10^4$.

For all algorithms other than **LoCoDL**, for which there is no function g , the functions f_i in (11) have a twice higher μ , so that the problem remains the same.

We considered several datasets from the LibSVM library (Chang & Lin, 2011) (3-clause BSD license). We show the results with the ‘a5a’ dataset in Figure 1 and with other datasets in the Appendix. We prepared each dataset by first shuffling it, then distributing it equally among the n clients (since m in (11) is an integer, the remaining datapoints were discarded). We used four different compression operators in the class $\mathbb{U}(\omega)$, for some $\omega \geq 0$:

- **rand- k** for some $k \in [d]$, which communicates $32k + k\lceil\log_2(d)\rceil$ bits. Indeed, the k randomly chosen values are sent in the standard 32-bits IEEE floating-point format, and their locations are encoded with $k\lceil\log_2(d)\rceil$ additional bits. We have $\omega = \frac{d}{k} - 1$.
- **Natural Compression** (Horváth et al., 2022), a form of quantization in which floats are encoded into 9 bits instead of 32 bits. We have $\omega = \frac{1}{8}$.
- A combination of **rand- k** and **Natural Compression**, in which the k chosen values are encoded into 9 bits, which yields a total of $9k + k\lceil\log_2(d)\rceil$ bits. We have $\omega = \frac{9d}{8k} - 1$.
- The l_1 -selection compressor, defined as $C(x) = \text{sign}(x_j)\|x\|_1 e_j$, where j is chosen randomly in $[d]$, with the probability of choosing $j' \in [d]$ equal to $|x_{j'}|/\|x\|_1$, and e_j is the j -th standard unit basis vector in \mathbb{R}^d . $\text{sign}(x_j)\|x\|_1$ is sent as a 32-bits float and the location of j is indicated with $\lceil\log_2(d)\rceil$, so that this compressor communicates $32 + \lceil\log_2(d)\rceil$ bits. Like with **rand-1**, we have $\omega = d - 1$.

The compressors at different clients are independent, so that $\omega_{\text{av}} = \frac{\omega}{n}$ in (3).

We can see that **LoCoDL**, when combined with **rand- k** and **Natural Compression**, converges faster than all other algorithms, with respect to the total number of communicated bits per client. We chose two different numbers n of clients, one with $n < d$ and another one with $n > 2d$, since the compressor of **CompressedScaffnew** is different in the two cases $n < 2d$ and $n > 2d$ (Condat et al., 2022a). **LoCoDL** outperforms **CompressedScaffnew** in both cases. As expected, all methods exhibit faster convergence with larger n . Remarkably, **ADIANA**, which has the best theoretical complexity for large n , improves upon **DIANA** but is not competitive with the LT-based methods **CompressedScaffnew**, **5GCS-CC**, and **LoCoDL**. This illustrates the power of doubly-accelerated methods based on a successful combination of LT and CC. In this class, our new proposed **LoCoDL** algorithm shines. For all algorithms, we used the theoretical parameter values given in their available convergence results (Corollary 3.2 for **LoCoDL**). We tried to tune the parameter values, such as k in **rand- k** and the (average) number of local steps per round, but this only gave minor improvements. For instance, **ADIANA** in Figure 1 was a bit faster with the best value of $k = 20$ than with $k = 30$. Increasing the learning rate γ led to inconsistent results, with sometimes divergence.

5 CONCLUSION

We have proposed **LoCoDL**, which combines a probabilistic Local Training mechanism similar to the one of **Scaffnew** and Communication Compression with a large class of unbiased compressors. This successful combination makes **LoCoDL** highly communication-efficient, with a doubly accelerated complexity with respect to the model dimension d and the condition number of the functions. In practice, **LoCoDL** outperforms other algorithms, including **ADIANA**, which has an even better complexity in theory obtained from Nesterov acceleration and not Local Training. This again shows the relevance of the popular mechanism of Local Training, which has been widely adopted in Federated Learning. A venue for future work is to implement bidirectional compression (Liu et al., 2020; Philippenko & Dieuleveut, 2021; Dorfman et al., 2023). We will also investigate extensions of our method with calls to stochastic gradient estimates, with or without variance reduction, as well as partial participation. These two features have been proposed for **Scaffnew** in Malinovsky et al. (2022) and Condat et al. (2023), but they are challenging to combine with generic compression.

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

REFERENCES

- A. Albasyoni, M. Safaryan, L. Condat, and P. Richtárik. Optimal gradient compression for distributed and federated learning. preprint arXiv:2010.03246, 2020.
- D. Basu, D. Data, C. Karakus, and S. N. Diggavi. Qsparse-Local-SGD: Distributed SGD With Quantization, Sparsification, and Local Computations. *IEEE Journal on Selected Areas in Information Theory*, 1(1):217–226, 2020.
- D. P. Bertsekas. *Convex optimization algorithms*. Athena Scientific, Belmont, MA, USA, 2015.
- A. Beznosikov, S. Horváth, P. Richtárik, and M. Safaryan. On biased compression for distributed learning. preprint arXiv:2002.12410, 2020.
- K. Bonawitz, V. Ivanov, B. Kreuter, A. Marcedone, H. B. McMahan, S. Patel, D. Ramage, A. Segal, and K. Seth. Practical secure aggregation for privacy-preserving machine learning. In *Proc. of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, pp. 1175–1191, 2017.
- C.-C. Chang and C.-J. Lin. LIBSVM: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, 2:27:1–27:27, 2011. Software available at <http://www.csie.ntu.edu.tw/%7Ecjlin/libsvm>.
- L. Condat and P. Richtárik. MURANA: A generic framework for stochastic variance-reduced optimization. In *Proc. of the conference Mathematical and Scientific Machine Learning (MSML), PMLR 190*, 2022.
- L. Condat and P. Richtárik. RandProx: Primal-dual optimization algorithms with randomized proximal updates. In *Proc. of International Conference on Learning Representations (ICLR)*, 2023.
- L. Condat, I. Agarský, and P. Richtárik. Provably doubly accelerated federated learning: The first theoretically successful combination of local training and compressed communication. preprint arXiv:2210.13277, 2022a.
- L. Condat, K. Li, and P. Richtárik. EF-BV: A unified theory of error feedback and variance reduction mechanisms for biased and unbiased compression in distributed optimization. In *Proc. of Conf. Neural Information Processing Systems (NeurIPS)*, 2022b.
- L. Condat, I. Agarský, G. Malinovsky, and P. Richtárik. TAMUNA: Doubly accelerated federated learning with local training, compression, and partial participation. preprint arXiv:2302.09832 presented at the *Int. Workshop on Federated Learning in the Age of Foundation Models in Conjunction with NeurIPS 2023*, 2023.
- R. Dorfman, S. Vargaftik, Y. Ben-Itzhak, and K. Y. Levy. DoCoFL: Downlink compression for cross-device federated learning. In *Proc. of Int. Conf. Machine Learning (ICML)*, 2023.
- I. Fatkhullin, I. Sokolov, E. Gorbunov, Z. Li, and P. Richtárik. EF21 with bells & whistles: Practical algorithmic extensions of modern error feedback. preprint arXiv:2110.03294, 2021.
- E. Gorbunov, F. Hanzely, and P. Richtárik. A unified theory of SGD: Variance reduction, sampling, quantization and coordinate descent. In *Proc. of 23rd Int. Conf. Artificial Intelligence and Statistics (AISTATS), PMLR 108*, 2020a.
- E. Gorbunov, D. Kovalev, D. Makarenko, and P. Richtárik. Linearly converging error compensated SGD. In *Proc. of Conf. Neural Information Processing Systems (NeurIPS)*, 2020b.
- E. Gorbunov, K. Burlachenko, Z. Li, and P. Richtárik. MARINA: Faster non-convex distributed learning with compression. In *Proc. of 38th Int. Conf. Machine Learning (ICML)*, pp. 3788–3798, 2021.
- R. M. Gower, M. Schmidt, F. Bach, and P. Richtárik. Variance-reduced methods for machine learning. *Proc. of the IEEE*, 108(11):1968–1983, November 2020.
- M. Grudzień, G. Malinovsky, and P. Richtárik. Can 5th Generation Local Training Methods Support Client Sampling? Yes! In *Proc. of Int. Conf. Artificial Intelligence and Statistics (AISTATS)*, 2023.

-
- 540 K. Gruntkowska, A. Tyurin, and P. Richtárik. EF21-P and friends: Improved theoretical communica-
541 tion complexity for distributed optimization with bidirectional compression. In *Proc. of 40th Int.*
542 *Conf. Machine Learning (ICML)*, 2023.
- 543 F. Haddadpour, M. M. Kamani, A. Mokhtari, and M. Mahdavi. Federated learning with compression:
544 Unified analysis and sharp guarantees. In *Proc. of Int. Conf. Artificial Intelligence and Statistics*
545 *(AISTATS)*, PMLR 130, pp. 2350–2358, 2021.
- 546 F. Hanzely and P. Richtárik. One method to rule them all: Variance reduction for data, parameters
547 and many new methods. preprint arXiv:1905.11266, 2019.
- 548 Y. He, X. Huang, and K. Yuan. Unbiased compression saves communication in distributed optimiza-
549 tion: When and how much? In *Proc. of Conf. Neural Information Processing Systems (NeurIPS)*,
550 2023.
- 551 S. Horváth, C.-Y. Ho, L. Horváth, A. N. Sahu, M. Canini, and P. Richtárik. Natural compression
552 for distributed deep learning. In *Proc. of the conference Mathematical and Scientific Machine*
553 *Learning (MSML)*, PMLR 190, 2022.
- 554 S. Horváth, D. Kovalev, K. Mishchenko, S. Stich, and P. Richtárik. Stochastic distributed learning
555 with gradient quantization and variance reduction. *Optimization Methods and Software*, 2022.
- 556 P. Kairouz et al. Advances and open problems in federated learning. *Foundations and Trends in*
557 *Machine Learning*, 14(1–2), 2021.
- 558 S. P. Karimireddy, S. Kale, M. Mohri, S. Reddi, S. U. Stich, and A. T. Suresh. SCAFFOLD: Stochastic
559 controlled averaging for federated learning. In *Proc. of 37th Int. Conf. Machine Learning (ICML)*,
560 pp. 5132–5143, 2020.
- 561 J. Konečný, H. B. McMahan, D. Ramage, and P. Richtárik. Federated optimization: distributed
562 machine learning for on-device intelligence. arXiv:1610.02527, 2016a.
- 563 J. Konečný, H. B. McMahan, F. X. Yu, P. Richtárik, A. T. Suresh, and D. Bacon. Federated learning:
564 Strategies for improving communication efficiency. In *NIPS Private Multi-Party Machine Learning*
565 *Workshop*, 2016b. arXiv:1610.05492.
- 566 Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to
567 document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 568 T. Li, A. K. Sahu, A. Talwalkar, and V. Smith. Federated learning: Challenges, methods, and future
569 directions. *IEEE Signal Processing Magazine*, 3(37):50–60, 2020a.
- 570 Z. Li, D. Kovalev, X. Qian, and P. Richtárik. Acceleration for compressed gradient descent in
571 distributed and federated optimization. In *Proc. of 37th Int. Conf. Machine Learning (ICML)*,
572 volume PMLR 119, 2020b.
- 573 Xiaorui Liu, Yao Li, Jiliang Tang, and Ming Yan. A double residual compression algorithm for
574 efficient distributed learning. In *Proc. of Int. Conf. Artificial Intelligence and Statistics (AISTATS)*,
575 PMLR 108, pp. 133–143, 2020.
- 576 G. Malinovsky and P. Richtárik. Federated random reshuffling with compression and variance
577 reduction. preprint arXiv:arXiv:2205.03914, 2022.
- 578 G. Malinovsky, K. Yi, and P. Richtárik. Variance reduced ProxSkip: Algorithm, theory and application
579 to federated learning. In *Proc. of Conf. Neural Information Processing Systems (NeurIPS)*, 2022.
- 580 A. Maranjyan, M. Safaryan, and P. Richtárik. GradSkip: Communication-accelerated local gradient
581 methods with better computational complexity. preprint arXiv:2210.16402, 2022.
- 582 H. B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. Agüera y Arcas. Communication-
583 efficient learning of deep networks from decentralized data. In *Proc. of Int. Conf. Artificial*
584 *Intelligence and Statistics (AISTATS)*, PMLR 54, 2017.
- 585 K. Mishchenko, E. Gorbunov, M. Takáč, and P. Richtárik. Distributed learning with compressed
586 gradient differences. arXiv:1901.09269, 2019.
- 587
- 588
- 589
- 590
- 591
- 592
- 593

594 K. Mishchenko, G. Malinovsky, S. Stich, and P. Richtárik. ProxSkip: Yes! Local Gradient Steps
595 Provably Lead to Communication Acceleration! Finally! In *Proc. of the 39th International*
596 *Conference on Machine Learning (ICML)*, July 2022.

597 C. Philippenko and A. Dieuleveut. Artemis: tight convergence guarantees for bidirectional compres-
598 sion in federated learning. preprint arXiv:2006.14591, 2020.

600 C. Philippenko and A. Dieuleveut. Preserved central model for faster bidirectional compression in
601 distributed settings. In *Proc. of Conf. Neural Information Processing Systems (NeurIPS)*, 2021.

602 A. Reiszadeh, A. Mokhtari, H. Hassani, A. Jadbabaie, and R. Pedarsani. FedPAQ: A communication-
603 efficient federated learning method with periodic averaging and quantization. In *Proc. of Int. Conf.*
604 *Artificial Intelligence and Statistics (AISTATS)*, pp. 2021–2031, 2020.

605 P. Richtárik, I. Sokolov, and I. Fatkhullin. EF21: A new, simpler, theoretically better, and practically
606 faster error feedback. In *Proc. of 35th Conf. Neural Information Processing Systems (NeurIPS)*,
607 2021.

609 A. Sadiev, D. Kovalev, and P. Richtárik. Communication acceleration of local gradient methods via
610 an accelerated primal-dual algorithm with an inexact prox. In *Proc. of Conf. Neural Information*
611 *Processing Systems (NeurIPS)*, 2022a.

612 A. Sadiev, G. Malinovsky, E. Gorbunov, I. Sokolov, A. Khaled, K. Burlachenko, and P. Richtárik.
613 Federated optimization algorithms with random reshuffling and gradient compression. preprint
614 arXiv:2206.07021, 2022b.

616 K. Scaman, F. Bach, S. Bubeck, Y. T. Lee, and L. Massoulié. Optimal convergence rates for convex
617 distributed optimization in networks. *Journal of Machine Learning Research*, 20:1–31, 2019.

618 S. Shalev-Shwartz and S. Ben-David. *Understanding machine learning: From theory to algorithms*.
619 Cambridge University Press, 2014.

621 S. Sra, S. Nowozin, and S. J. Wright. *Optimization for Machine Learning*. The MIT Press, 2011.

622 A. Tyurin and P. Richtárik. DASHA: Distributed nonconvex optimization with communication
623 compression, optimal oracle complexity, and no client synchronization. In *Proc. of International*
624 *Conference on Learning Representations (ICLR)*, 2023a.

626 A. Tyurin and P. Richtárik. 2Direction: Theoretically faster distributed training with bidirectional
627 communication compression. In *Proc. of Conf. Neural Information Processing Systems (NeurIPS)*,
628 2023b.

629 J. Wang et al. A field guide to federated optimization. preprint arXiv:2107.06917, 2021.

631 H. Yang, M. Fang, and J. Liu. Achieving linear speedup with partial worker participation in non-IID
632 federated learning. In *Proc. of International Conference on Learning Representations (ICLR)*,
633 2021.

634 H. Yang, P. Qiu, P. Khanduri, M. Fang, and J. Liu. Understanding server-assisted federated learning
635 in the presence of incomplete client participation. In *Proc. of Int. Conf. Machine Learning (ICML)*,
636 2024.

637 Y. Zhao, M. Li, L. Lai, N. Suda, D. Civin, and V. Chandra. Federated learning with non-iid data.
638 preprint arXiv:1806.00582, 2018.

639
640
641
642
643
644
645
646
647

648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

Appendix

CONTENTS

1	Introduction	1
1.1	Problem and Motivation	1
1.2	State of the Art	2
1.3	A General Class of Unbiased Random Compressors	4
1.4	Challenge and Contributions	4
2	Proposed Algorithm LoCoDL	5
2.1	Principle: Double Lifting of the Problem to a Consensus Problem	5
2.2	Description of LoCoDL	5
3	Convergence and Complexity of LoCoDL	6
3.1	The Case $g = 0$	8
4	Experiments	8
5	Conclusion	9
A	Proof of Theorem 3.1	13
B	Additional Experiments	17
B.1	Experiment with the Dirichlet Distribution	17

A PROOF OF THEOREM 3.1

We define the Euclidean space $\mathcal{X} := \mathbb{R}^d$ and the product space $\mathcal{X} := \mathcal{X}^{n+1}$ endowed with the weighted inner product

$$\langle \mathbf{x}, \mathbf{x}' \rangle_{\mathcal{X}} := \sum_{i=1}^n \langle x_i, x'_i \rangle + n \langle y, y' \rangle, \quad \forall \mathbf{x} = (x_1, \dots, x_n, y), \mathbf{x}' = (x'_1, \dots, x'_n, y'). \quad (12)$$

We define the copy operator $\mathbf{1} : x \in \mathcal{X} \mapsto (x, \dots, x, x) \in \mathcal{X}$ and the linear operator

$$S : \mathbf{x} \in \mathcal{X} \mapsto \mathbf{1}\bar{x}, \quad \text{with } \bar{x} = \frac{1}{2n} \left(\sum_{i=1}^n x_i + ny \right). \quad (13)$$

S is the orthogonal projector in \mathcal{X} onto the consensus line $\{\mathbf{x} \in \mathcal{X} : x_1 = \dots = x_n = y\}$. We also define the linear operator

$$W := \text{Id} - S : \mathbf{x} = (x_1, \dots, x_n, y) \in \mathcal{X} \mapsto (x_1 - \bar{x}, \dots, x_n - \bar{x}, y - \bar{x}), \quad \text{with } \bar{x} = \frac{1}{2n} \left(\sum_{i=1}^n x_i + ny \right), \quad (14)$$

where Id denotes the identity. W is the orthogonal projector in \mathcal{X} onto the hyperplane $\{\mathbf{x} \in \mathcal{X} : x_1 + \dots + x_n + ny = 0\}$, which is orthogonal to the consensus line. As such, it is self-adjoint, positive semidefinite, its eigenvalues are $(1, \dots, 1, 0)$, its kernel is the consensus line, and its spectral

norm is 1. Also, $W^2 = W$. Note that we can write W in terms of the differences $d_i = x_i - y$ and $\bar{d} = \frac{1}{2n} \sum_{i=1}^n d_i$:

$$W : \mathbf{x} = (x_1, \dots, x_n, y) \mapsto (d_1 - \bar{d}, \dots, d_n - \bar{d}, -\bar{d}). \quad (15)$$

Since for every $\mathbf{x} = (x_1, \dots, x_n, y)$, $W\mathbf{x} = \mathbf{0} := (0, \dots, 0, 0)$ if and only if $x_1 = \dots = x_n = y$, we can reformulate the problem (1) as

$$\min_{\mathbf{x}=(x_1, \dots, x_n, y) \in \mathcal{X}} \mathbf{f}(\mathbf{x}) \quad \text{s.t.} \quad W\mathbf{x} = \mathbf{0}, \quad (16)$$

where $\mathbf{f}(\mathbf{x}) := \sum_{i=1}^n f_i(x_i) + ng(y)$. Note that in \mathcal{X} , \mathbf{f} is L -smooth and μ -strongly convex, and $\nabla \mathbf{f}(\mathbf{x}) = (\nabla f_1(x_1), \dots, \nabla f_n(x_n), \nabla g(y))$.

Let $t \geq 0$. We also introduce vector notations for the variables of the algorithm: $\mathbf{x}^t := (x_1^t, \dots, x_n^t, y^t)$, $\hat{\mathbf{x}}^t := (\hat{x}_1^t, \dots, \hat{x}_n^t, \hat{y}^t)$, $\mathbf{u}^t := (u_1^t, \dots, u_n^t, v^t)$, $\mathbf{u}^* := (u_1^*, \dots, u_n^*, v^*)$, $\mathbf{w}^t := \mathbf{x}^t - \gamma \nabla \mathbf{f}(\mathbf{x}^t)$, $\mathbf{w}^* := \mathbf{x}^* - \gamma \nabla \mathbf{f}(\mathbf{x}^*)$, where $\mathbf{x}^* := \mathbf{1}x^*$ is the unique solution to (16). We also define $\bar{x}^t := \frac{1}{2n} (\sum_{i=1}^n \hat{x}_i^t + n\hat{y}^t)$ and $\lambda := \frac{pX}{\gamma(1+2\omega)}$.

Then we can write the iteration of **LoCoDL** as

$$\left\{ \begin{array}{l} \hat{\mathbf{x}}^t := \mathbf{x}^t - \gamma \nabla \mathbf{f}(\mathbf{x}^t) + \gamma \mathbf{u}^t = \mathbf{w}^t + \gamma \mathbf{u}^t \\ \text{flip a coin } \theta^t \in \{0, 1\} \text{ with } \text{Prob}(\theta^t = 1) = p \\ \text{if } \theta^t = 1 \\ \quad \mathbf{d}^t := (\mathcal{C}_1^t(\hat{x}_1^t - \hat{y}^t), \dots, \mathcal{C}_n^t(\hat{x}_n^t - \hat{y}^t), 0) \\ \quad \bar{d}^t := \frac{1}{2n} \sum_{j=1}^n d_j^t \\ \quad \mathbf{x}^{t+1} := (1 - \rho)\hat{\mathbf{x}}^t + \rho \mathbf{1}(\hat{y}^t + \bar{d}^t) \\ \quad \mathbf{u}^{t+1} := \mathbf{u}^t + \lambda(\mathbf{1}\bar{d}^t - \mathbf{d}^t) = \mathbf{u}^t - \lambda W \mathbf{d}^t \\ \text{else} \\ \quad \mathbf{x}^{t+1} := \hat{\mathbf{x}}^t \\ \quad \mathbf{u}^{t+1} := \mathbf{u}^t \\ \text{end if} \end{array} \right. \quad (17)$$

We denote by \mathcal{F}^t the σ -algebra generated by the collection of \mathcal{X} -valued random variables $\mathbf{x}^0, \mathbf{u}^0, \dots, \mathbf{x}^t, \mathbf{u}^t$.

Since we suppose that $S\mathbf{u}^0 = \mathbf{0}$ and we have $S\bar{W}\mathbf{d}^{t'} = \mathbf{0}$ in the update of \mathbf{u} , we have $S\mathbf{u}^{t'} = \mathbf{0}$ for every $t' \geq 0$.

If $\theta^t = 1$, we have

$$\begin{aligned} \|\mathbf{u}^{t+1} - \mathbf{u}^*\|_{\mathcal{X}}^2 &= \|\mathbf{u}^t - \mathbf{u}^*\|_{\mathcal{X}}^2 + \lambda^2 \|W\mathbf{d}^t\|_{\mathcal{X}}^2 - 2\lambda \langle \mathbf{u}^t - \mathbf{u}^*, W\mathbf{d}^t \rangle_{\mathcal{X}} \\ &= \|\mathbf{u}^t - \mathbf{u}^*\|_{\mathcal{X}}^2 + \lambda^2 \|\mathbf{d}^t\|_{\mathcal{X}}^2 - \lambda^2 \|S\mathbf{d}^t\|_{\mathcal{X}}^2 - 2\lambda \langle \mathbf{u}^t - \mathbf{u}^*, \mathbf{d}^t \rangle_{\mathcal{X}}, \end{aligned}$$

because $S\mathbf{u}^t = S\mathbf{u}^* = \mathbf{0}$, so that $\langle \mathbf{u}^t - \mathbf{u}^*, S\mathbf{d}^t \rangle_{\mathcal{X}} = 0$.

The variance inequality (2) satisfied by the compressors \mathcal{C}_i^t is equivalent to $\mathbb{E}[\|\mathcal{C}_i^t(x)\|^2] \leq (1 + \omega) \|x\|^2$, so that

$$\mathbb{E}[\|\mathbf{d}^t\|_{\mathcal{X}}^2 \mid \mathcal{F}^t, \theta^t = 1] \leq (1 + \omega) \|\hat{\mathbf{x}}^t - \mathbf{1}\hat{y}^t\|_{\mathcal{X}}^2.$$

Also,

$$\mathbb{E}[\mathbf{d}^t \mid \mathcal{F}^t, \theta^t = 1] = \hat{\mathbf{x}}^t - \mathbf{1}\hat{y}^t.$$

Thus,

$$\begin{aligned} \mathbb{E}[\|\mathbf{u}^{t+1} - \mathbf{u}^*\|_{\mathcal{X}}^2 \mid \mathcal{F}^t] &= (1 - p) \|\mathbf{u}^t - \mathbf{u}^*\|_{\mathcal{X}}^2 + p \mathbb{E}[\|\mathbf{u}^{t+1} - \mathbf{u}^*\|_{\mathcal{X}}^2 \mid \mathcal{F}^t, \theta^t = 1] \\ &\leq \|\mathbf{u}^t - \mathbf{u}^*\|_{\mathcal{X}}^2 + p\lambda^2(1 + \omega) \|\hat{\mathbf{x}}^t - \mathbf{1}\hat{y}^t\|_{\mathcal{X}}^2 - p\lambda^2 \mathbb{E}[\|S\mathbf{d}^t\|_{\mathcal{X}}^2 \mid \mathcal{F}^t, \theta^t = 1] \\ &\quad - 2p\lambda \langle \mathbf{u}^t - \mathbf{u}^*, \hat{\mathbf{x}}^t - \mathbf{1}\hat{y}^t \rangle_{\mathcal{X}} \\ &= \|\mathbf{u}^t - \mathbf{u}^*\|_{\mathcal{X}}^2 + p\lambda^2(1 + \omega) \|\hat{\mathbf{x}}^t - \mathbf{1}\hat{y}^t\|_{\mathcal{X}}^2 - p\lambda^2 \mathbb{E}[\|S\mathbf{d}^t\|_{\mathcal{X}}^2 \mid \mathcal{F}^t, \theta^t = 1] \\ &\quad - 2p\lambda \langle \mathbf{u}^t - \mathbf{u}^*, \hat{\mathbf{x}}^t \rangle_{\mathcal{X}}. \end{aligned}$$

Moreover, $\mathbb{E}\left[\|S\mathbf{d}^t\|_{\mathcal{X}}^2 \mid \mathcal{F}^t, \theta^t = 1\right] \geq \|\mathbb{E}[S\mathbf{d}^t \mid \mathcal{F}^t, \theta^t = 1]\|_{\mathcal{X}}^2 = \|S\hat{\mathbf{x}}^t - \mathbf{1}\hat{y}^t\|_{\mathcal{X}}^2$ and $\|\hat{\mathbf{x}}^t - \mathbf{1}\hat{y}^t\|_{\mathcal{X}}^2 = \|S\hat{\mathbf{x}}^t - \mathbf{1}\hat{y}^t\|_{\mathcal{X}}^2 + \|W\hat{\mathbf{x}}^t\|_{\mathcal{X}}^2$, so that

$$\begin{aligned} \mathbb{E}\left[\|\mathbf{u}^{t+1} - \mathbf{u}^*\|_{\mathcal{X}}^2 \mid \mathcal{F}^t\right] &\leq \|\mathbf{u}^t - \mathbf{u}^*\|_{\mathcal{X}}^2 + p\lambda^2(1 + \omega) \|\hat{\mathbf{x}}^t - \mathbf{1}\hat{y}^t\|_{\mathcal{X}}^2 - p\lambda^2 \|S\hat{\mathbf{x}}^t - \mathbf{1}\hat{y}^t\|_{\mathcal{X}}^2 \\ &\quad - 2p\lambda \langle \mathbf{u}^t - \mathbf{u}^*, \hat{\mathbf{x}}^t \rangle_{\mathcal{X}} \\ &= \|\mathbf{u}^t - \mathbf{u}^*\|_{\mathcal{X}}^2 + p\lambda^2\omega \|\hat{\mathbf{x}}^t - \mathbf{1}\hat{y}^t\|_{\mathcal{X}}^2 + p\lambda^2 \|W\hat{\mathbf{x}}^t\|_{\mathcal{X}}^2 - 2p\lambda \langle \mathbf{u}^t - \mathbf{u}^*, \hat{\mathbf{x}}^t \rangle_{\mathcal{X}}. \end{aligned}$$

From the Peter–Paul inequality $\|a + b\|^2 \leq 2\|a\|^2 + 2\|b\|^2$ for any a and b , we have

$$\begin{aligned} \|\hat{\mathbf{x}}^t - \mathbf{1}\hat{y}^t\|_{\mathcal{X}}^2 &= \sum_{i=1}^n \|\hat{x}_i^t - \hat{y}^t\|^2 = \sum_{i=1}^n \|(\hat{x}_i^t - \bar{x}^t) - (\hat{y}^t - \bar{x}^t)\|^2 \\ &\leq \sum_{i=1}^n \left(2\|\hat{x}_i^t - \bar{x}^t\|^2 + 2\|\hat{y}^t - \bar{x}^t\|^2\right) \\ &= 2 \left(\sum_{i=1}^n \|\hat{x}_i^t - \bar{x}^t\|^2 + n\|\hat{y}^t - \bar{x}^t\|^2 \right) \\ &= 2 \|\hat{\mathbf{x}}^t - \mathbf{1}\bar{x}^t\|_{\mathcal{X}}^2 = 2 \|W\hat{\mathbf{x}}^t\|_{\mathcal{X}}^2. \end{aligned} \tag{18}$$

Hence,

$$\mathbb{E}\left[\|\mathbf{u}^{t+1} - \mathbf{u}^*\|_{\mathcal{X}}^2 \mid \mathcal{F}^t\right] \leq \|\mathbf{u}^t - \mathbf{u}^*\|_{\mathcal{X}}^2 + p\lambda^2(1 + 2\omega) \|W\hat{\mathbf{x}}^t\|_{\mathcal{X}}^2 - 2p\lambda \langle \mathbf{u}^t - \mathbf{u}^*, \hat{\mathbf{x}}^t \rangle_{\mathcal{X}}.$$

On the other hand,

$$\begin{aligned} \mathbb{E}\left[\|\mathbf{x}^{t+1} - \mathbf{x}^*\|_{\mathcal{X}}^2 \mid \mathcal{F}^t, \theta = 1\right] &= (1 - \rho)^2 \|\hat{\mathbf{x}}^t - \mathbf{x}^*\|_{\mathcal{X}}^2 + \rho^2 \mathbb{E}\left[\|\mathbf{1}(\hat{y}^t + \bar{d}^t) - \mathbf{x}^*\|_{\mathcal{X}}^2 \mid \mathcal{F}^t, \theta = 1\right] \\ &\quad + 2\rho(1 - \rho) \langle \hat{\mathbf{x}}^t - \mathbf{x}^*, \mathbf{1}(\hat{y}^t + \mathbb{E}[\bar{d}^t \mid \mathcal{F}^t, \theta = 1]) - \mathbf{x}^* \rangle_{\mathcal{X}}. \end{aligned}$$

We have $\mathbb{E}[\bar{d}^t \mid \mathcal{F}^t, \theta = 1] = \frac{1}{2n} \sum_{i=1}^n \hat{x}_i^t - \frac{1}{2}\hat{y}^t = \bar{x}^t - \hat{y}^t$, so that

$$\mathbf{1}(\hat{y}^t + \mathbb{E}[\bar{d}^t \mid \mathcal{F}^t, \theta = 1]) = \mathbf{1}\bar{x}^t = S\hat{\mathbf{x}}^t.$$

In addition,

$$\langle \hat{\mathbf{x}}^t - \mathbf{x}^*, S\hat{\mathbf{x}}^t - \mathbf{x}^* \rangle_{\mathcal{X}} = \langle \hat{\mathbf{x}}^t - \mathbf{x}^*, S(\hat{\mathbf{x}}^t - \mathbf{x}^*) \rangle_{\mathcal{X}} = \|S(\hat{\mathbf{x}}^t - \mathbf{x}^*)\|_{\mathcal{X}}^2.$$

Moreover,

$$\begin{aligned} \mathbb{E}\left[\|\mathbf{1}(\hat{y}^t + \bar{d}^t) - \mathbf{x}^*\|_{\mathcal{X}}^2 \mid \mathcal{F}^t, \theta = 1\right] &= \|\mathbf{1}(\hat{y}^t + \mathbb{E}[\bar{d}^t \mid \mathcal{F}^t, \theta = 1]) - \mathbf{x}^*\|_{\mathcal{X}}^2 \\ &\quad + \mathbb{E}\left[\|\mathbf{1}(\bar{d}^t - \mathbb{E}[\bar{d}^t \mid \mathcal{F}^t, \theta = 1])\|_{\mathcal{X}}^2 \mid \mathcal{F}^t, \theta = 1\right] \\ &= \|S\hat{\mathbf{x}}^t - \mathbf{x}^*\|_{\mathcal{X}}^2 \\ &\quad + 2n\mathbb{E}\left[\|\bar{d}^t - \mathbb{E}[\bar{d}^t \mid \mathcal{F}^t, \theta = 1]\|^2 \mid \mathcal{F}^t, \theta = 1\right] \end{aligned}$$

and, using (3),

$$\begin{aligned} \mathbb{E}\left[\|\bar{d}^t - \mathbb{E}[\bar{d}^t \mid \mathcal{F}^t, \theta = 1]\|^2 \mid \mathcal{F}^t, \theta = 1\right] &\leq \frac{\omega_{\text{av}}}{4n} \sum_{i=1}^n \|\hat{x}_i^t - \hat{y}^t\|^2 \\ &\leq \frac{\omega_{\text{av}}}{2n} \|W\hat{\mathbf{x}}^t\|_{\mathcal{X}}^2, \end{aligned}$$

where the second inequality follows from (18). Hence,

$$\begin{aligned}
\mathbb{E}\left[\|\mathbf{x}^{t+1} - \mathbf{x}^*\|_{\mathcal{X}}^2 \mid \mathcal{F}^t, \theta = 1\right] &\leq (1 - \rho)^2 \|\hat{\mathbf{x}}^t - \mathbf{x}^*\|_{\mathcal{X}}^2 + \rho^2 \|S\hat{\mathbf{x}}^t - \mathbf{x}^*\|_{\mathcal{X}}^2 + \rho^2 \omega_{\text{av}} \|W\hat{\mathbf{x}}^t\|_{\mathcal{X}}^2 \\
&\quad + 2\rho(1 - \rho) \|S(\hat{\mathbf{x}}^t - \mathbf{x}^*)\|_{\mathcal{X}}^2 \\
&= (1 - \rho)^2 \|\hat{\mathbf{x}}^t - \mathbf{x}^*\|_{\mathcal{X}}^2 + \rho^2 \omega_{\text{av}} \|W\hat{\mathbf{x}}^t\|_{\mathcal{X}}^2 \\
&\quad + (2\rho - \rho^2) \|S(\hat{\mathbf{x}}^t - \mathbf{x}^*)\|_{\mathcal{X}}^2 \\
&= (1 - \rho)^2 \|\hat{\mathbf{x}}^t - \mathbf{x}^*\|_{\mathcal{X}}^2 + \rho^2 \omega_{\text{av}} \|W\hat{\mathbf{x}}^t\|_{\mathcal{X}}^2 \\
&\quad + (2\rho - \rho^2) \left(\|\hat{\mathbf{x}}^t - \mathbf{x}^*\|_{\mathcal{X}}^2 - \|W\hat{\mathbf{x}}^t\|_{\mathcal{X}}^2 \right) \\
&= \|\hat{\mathbf{x}}^t - \mathbf{x}^*\|_{\mathcal{X}}^2 - (2\rho - \rho^2 - \rho^2 \omega_{\text{av}}) \|W\hat{\mathbf{x}}^t\|_{\mathcal{X}}^2
\end{aligned}$$

and

$$\begin{aligned}
\mathbb{E}\left[\|\mathbf{x}^{t+1} - \mathbf{x}^*\|_{\mathcal{X}}^2 \mid \mathcal{F}^t\right] &= (1 - p) \|\hat{\mathbf{x}}^t - \mathbf{x}^*\|_{\mathcal{X}}^2 + p\mathbb{E}\left[\|\mathbf{x}^{t+1} - \mathbf{x}^*\|_{\mathcal{X}}^2 \mid \mathcal{F}^t, \theta^t = 1\right] \\
&\leq \|\hat{\mathbf{x}}^t - \mathbf{x}^*\|_{\mathcal{X}}^2 - p(2\rho - \rho^2(1 + \omega_{\text{av}})) \|W\hat{\mathbf{x}}^t\|_{\mathcal{X}}^2.
\end{aligned}$$

Furthermore,

$$\begin{aligned}
\|\hat{\mathbf{x}}^t - \mathbf{x}^*\|_{\mathcal{X}}^2 &= \|\mathbf{w}^t - \mathbf{w}^*\|_{\mathcal{X}}^2 + \gamma^2 \|\mathbf{u}^t - \mathbf{u}^*\|_{\mathcal{X}}^2 + 2\gamma \langle \mathbf{w}^t - \mathbf{w}^*, \mathbf{u}^t - \mathbf{u}^* \rangle_{\mathcal{X}} \\
&= \|\mathbf{w}^t - \mathbf{w}^*\|_{\mathcal{X}}^2 - \gamma^2 \|\mathbf{u}^t - \mathbf{u}^*\|_{\mathcal{X}}^2 + 2\gamma \langle \hat{\mathbf{x}}^t - \mathbf{x}^*, \mathbf{u}^t - \mathbf{u}^* \rangle_{\mathcal{X}} \\
&= \|\mathbf{w}^t - \mathbf{w}^*\|_{\mathcal{X}}^2 - \gamma^2 \|\mathbf{u}^t - \mathbf{u}^*\|_{\mathcal{X}}^2 + 2\gamma \langle \hat{\mathbf{x}}^t, \mathbf{u}^t - \mathbf{u}^* \rangle_{\mathcal{X}},
\end{aligned}$$

which yields

$$\begin{aligned}
\mathbb{E}\left[\|\mathbf{x}^{t+1} - \mathbf{x}^*\|_{\mathcal{X}}^2 \mid \mathcal{F}^t\right] &\leq \|\mathbf{w}^t - \mathbf{w}^*\|_{\mathcal{X}}^2 - \gamma^2 \|\mathbf{u}^t - \mathbf{u}^*\|_{\mathcal{X}}^2 + 2\gamma \langle \hat{\mathbf{x}}^t, \mathbf{u}^t - \mathbf{u}^* \rangle_{\mathcal{X}} \\
&\quad - p(2\rho - \rho^2(1 + \omega_{\text{av}})) \|W\hat{\mathbf{x}}^t\|_{\mathcal{X}}^2.
\end{aligned}$$

Hence, with $\lambda = \frac{p\chi}{\gamma(1+2\omega)}$,

$$\begin{aligned}
&\frac{1}{\gamma} \mathbb{E}\left[\|\mathbf{x}^{t+1} - \mathbf{x}^*\|_{\mathcal{X}}^2 \mid \mathcal{F}^t\right] + \frac{\gamma(1+2\omega)}{p^2\chi} \mathbb{E}\left[\|\mathbf{u}^{t+1} - \mathbf{u}^*\|_{\mathcal{X}}^2 \mid \mathcal{F}^t\right] \\
&\leq \frac{1}{\gamma} \|\mathbf{w}^t - \mathbf{w}^*\|_{\mathcal{X}}^2 - \gamma \|\mathbf{u}^t - \mathbf{u}^*\|_{\mathcal{X}}^2 + 2\langle \hat{\mathbf{x}}^t, \mathbf{u}^t - \mathbf{u}^* \rangle_{\mathcal{X}} - \frac{p}{\gamma} (2\rho - \rho^2(1 + \omega_{\text{av}})) \|W\hat{\mathbf{x}}^t\|_{\mathcal{X}}^2 \\
&\quad + \frac{\gamma(1+2\omega)}{p^2\chi} \|\mathbf{u}^t - \mathbf{u}^*\|_{\mathcal{X}}^2 + \frac{p\chi}{\gamma} \|W\hat{\mathbf{x}}^t\|_{\mathcal{X}}^2 - 2\langle \mathbf{u}^t - \mathbf{u}^*, \hat{\mathbf{x}}^t \rangle_{\mathcal{X}} \\
&= \frac{1}{\gamma} \|\mathbf{w}^t - \mathbf{w}^*\|_{\mathcal{X}}^2 + \frac{\gamma(1+2\omega)}{p^2\chi} \left(1 - \frac{p^2\chi}{1+2\omega}\right) \|\mathbf{u}^t - \mathbf{u}^*\|_{\mathcal{X}}^2 \\
&\quad - \frac{p}{\gamma} (2\rho - \rho^2(1 + \omega_{\text{av}}) - \chi) \|W\hat{\mathbf{x}}^t\|_{\mathcal{X}}^2.
\end{aligned}$$

Therefore, assuming that $2\rho - \rho^2(1 + \omega_{\text{av}}) - \chi \geq 0$,

$$\mathbb{E}[\Psi^{t+1} \mid \mathcal{F}^t] \leq \frac{1}{\gamma} \|\mathbf{w}^t - \mathbf{w}^*\|_{\mathcal{X}}^2 + \left(1 - \frac{p^2\chi}{1+2\omega}\right) \frac{\gamma(1+2\omega)}{p^2\chi} \|\mathbf{u}^t - \mathbf{u}^*\|_{\mathcal{X}}^2.$$

According to Condat & Richtárik (2023, Lemma 1),

$$\begin{aligned}
\|\mathbf{w}^t - \mathbf{w}^*\|_{\mathcal{X}}^2 &= \|(\text{Id} - \gamma\nabla\mathbf{f})\mathbf{x}^t - (\text{Id} - \gamma\nabla\mathbf{f})\mathbf{x}^*\|_{\mathcal{X}}^2 \\
&\leq \max(1 - \gamma\mu, \gamma L - 1)^2 \|\mathbf{x}^t - \mathbf{x}^*\|_{\mathcal{X}}^2.
\end{aligned}$$

Hence,

$$\mathbb{E}[\Psi^{t+1} \mid \mathcal{F}^t] \leq \max\left((1 - \gamma\mu)^2, (1 - \gamma L)^2, 1 - \frac{p^2\chi}{1+2\omega}\right) \Psi^t. \quad (19)$$

Using the tower rule, we can unroll the recursion in (19) to obtain the unconditional expectation of Ψ^{t+1} .

Using classical results on supermartingale convergence (Bertsekas, 2015, Proposition A.4.5), it follows from (19) that $\Psi^t \rightarrow 0$ almost surely. Almost sure convergence of \mathbf{x}^t and \mathbf{u}^t follows.

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

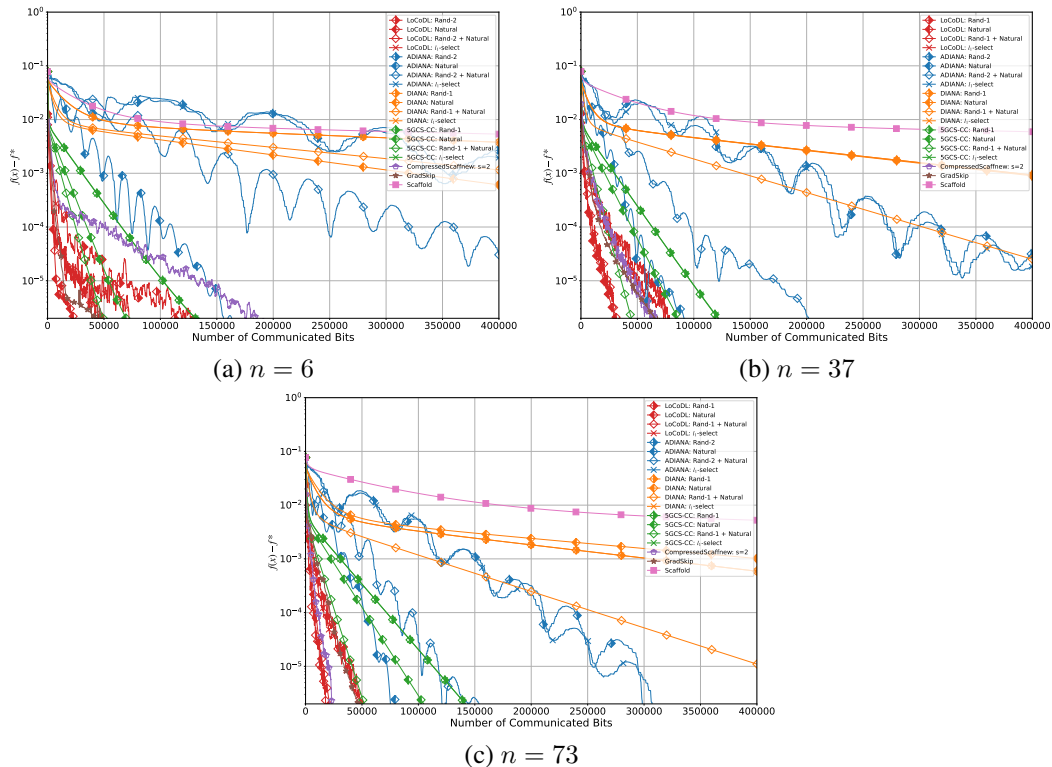


Figure 2: Comparison of several algorithms with several compressors on logistic regression with the ‘diabetes’ dataset from the LibSVM, which has $d = 8$ and 768 data points. We chose different values of n to illustrate the three regimes $n < d$, $n > d$, $n > d^2$, as discussed at the end of Section 3.

B ADDITIONAL EXPERIMENTS

The results for the experiments in Section 4 with the ‘diabetes’ dataset from the LibSVM library (Chang & Lin, 2011) are shown in Figure 2. The results with the ‘w1a’, ‘australian’ and ‘covtype.binary’ datasets, for the same logistic regression problem with $\kappa = 10^4$, are shown in Figures 3, 4 and 5. Finally, we also run experiments on MNIST dataset (LeCun et al., 1998) in Figure 6.

Consistent with our previous findings, **LoCoDL** outperforms the other algorithms in terms of communication efficiency.

B.1 EXPERIMENT WITH THE DIRICHLET DISTRIBUTION

In this section, we investigate how heterogeneity of the functions influences the convergence. We consider logistic regression as above, but with synthetic data sampled from the Dirichlet distribution of parameter α . If α is small, the Dirichlet distribution becomes similar to the uniform distribution over the simplex, which corresponds to the heterogeneous case where there is no similarity between the data. If α is large, the samples of the Dirichlet distribution tend to be similar to each other and concentrated around the middle point $(1/d, \dots, 1/d)$ of the simplex. We set $n = 100$ and $d = 10$, and a single random sample is assigned to each client. The results are shown in Figure 7, for $\alpha = 1$ and $\alpha = 10$. With these values of n and d , **ADIANA** has a better theoretical complexity than **LoCoDL**. However, in practice, we observe that **LoCoDL** again outperforms **ADIANA**. For both methods, joint sparsification and quantization with rand-1 and natural compression performs best. There is no significant difference depending on the value of α .

918
 919
 920
 921
 922
 923
 924
 925
 926
 927
 928
 929
 930
 931
 932
 933
 934
 935
 936
 937
 938
 939
 940
 941
 942
 943
 944
 945
 946
 947
 948
 949
 950
 951
 952
 953
 954
 955
 956
 957
 958
 959
 960
 961
 962
 963
 964
 965
 966
 967
 968
 969
 970
 971

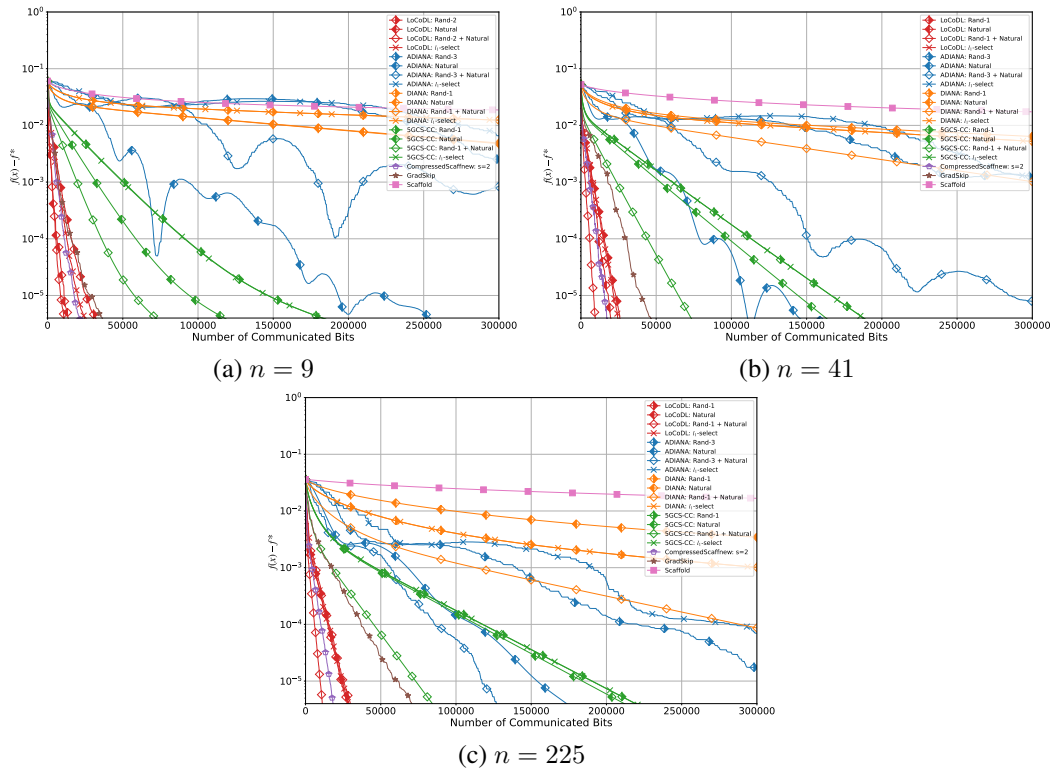


Figure 3: Comparison of several algorithms with various compressors on logistic regression with the ‘australian’ dataset from the LibSVM, which has $d = 14$ and 690 data points. We chose different values of n to illustrate the three regimes: $n < d$, $n > d$, $n > d^2$, as discussed at the end of Section 3.

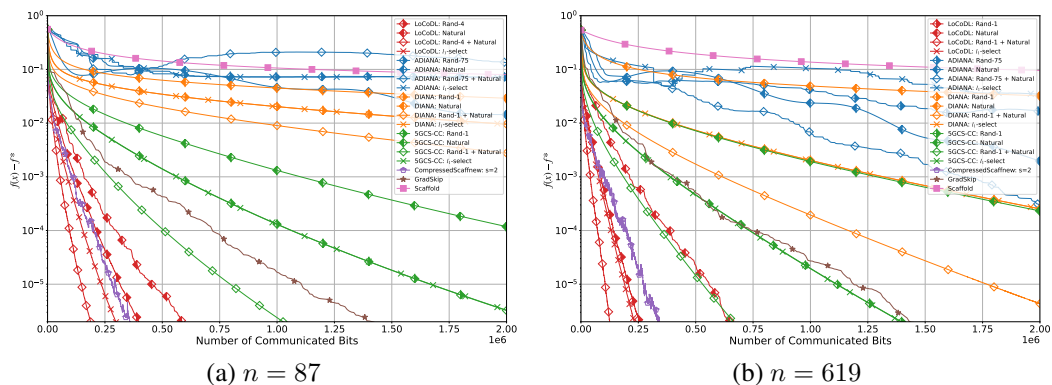


Figure 4: Comparison of several algorithms with various compressors on logistic regression with the ‘w1a’ dataset from the LibSVM, which has $d = 300$ and 2,477 data points. We chose different values of n to illustrate the two regimes, $n < d$ and $n > d$, as discussed at the end of Section 3.

1026
 1027
 1028
 1029
 1030
 1031
 1032
 1033
 1034
 1035
 1036
 1037
 1038
 1039
 1040
 1041
 1042
 1043
 1044
 1045
 1046
 1047
 1048
 1049
 1050
 1051
 1052
 1053
 1054
 1055
 1056
 1057
 1058
 1059
 1060
 1061
 1062
 1063
 1064
 1065
 1066
 1067
 1068
 1069
 1070
 1071
 1072
 1073
 1074
 1075
 1076
 1077
 1078
 1079

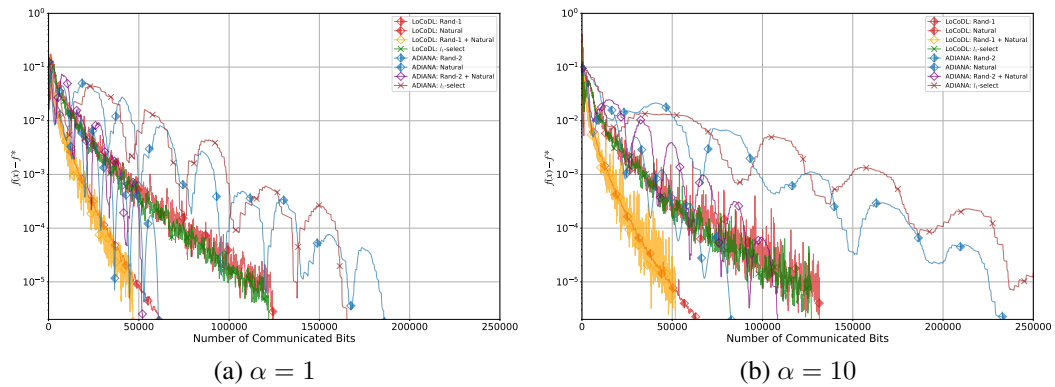


Figure 7: Comparison of **LoCoDL** and **ADIANA** with various compressors on logistic regression ($n = 100, d = 10$), with samples from the Dirichlet distribution of parameter α , with $\alpha = 1$ on the left and $\alpha = 10$ on the right.