

---

# SGD with Low-Dimensional Gradients with Applications to Private and Distributed Learning

---

Shiva Prasad Kasiviswanathan<sup>1</sup>

<sup>1</sup>Amazon, Palo Alto, USA

## Abstract

In this paper, we consider constrained optimization problems subject to a convex set  $\mathcal{C}$ . Stochastic gradient descent (SGD) is a simple and popular stochastic optimization algorithm that has been the workhorse of machine learning for many years. We show a new and surprising fact about SGD, in that depending on the constraint set  $\mathcal{C}$ , one can operate SGD with much lower-dimensional stochastic gradients without affecting its performance. In particular, we design an optimization algorithm that operates with the lower-dimensional (compressed) stochastic gradients, and establish that with the right set of parameters it has the exact same dimension-free convergence guarantees as that of regular SGD for popular convex and nonconvex optimization settings. We also present two applications of these bounds, one for improving the empirical risk minimization bounds in differentially private nonconvex optimization, and other for reducing communication costs with distributed SGD. Additionally, we also show that this connection between constraint set structure and gradient compression also extends beyond SGD to the conditional gradient (Frank-Wolfe) method. The geometry of the constraint set, captured by its Gaussian width, plays an important role in all our results.

## 1 INTRODUCTION

In this paper, we consider the classic optimization problem of minimizing a function over a convex set:

$$\min_{\mathbf{w} \in \mathcal{C}} F(\mathbf{w}), \quad (1)$$

where  $\mathcal{C} \subseteq \mathbb{R}^d$  is a compact convex set and  $F : \mathbb{R}^d \rightarrow \mathbb{R}$ . (Projected) Stochastic Gradient Descent (SGD) is one of the

simplest and most popular stochastic first-order optimization algorithms for solving (1). Stochastic gradient descent is widely used in machine learning applications due to its favorable computational scalability properties. This is notably true in the deep learning setting, where gradients can be computed efficiently via backpropagation. In convex optimization, SGD can be used to optimize any convex function  $F$  over a convex domain  $\mathcal{C}$ , given access only to unbiased estimates of  $F$ 's gradients (or more generally, subgradients<sup>1</sup>) through an oracle. This feature makes it very useful for learning problems, where our goal is to minimize generalization error based only on a finite sampled training set.

While classical results have focused on analyzing the performance of SGD in convex optimization problems, the most notable recent successes in machine learning have involved nonconvex optimization. For example, in the unconstrained setting (i.e.,  $\mathcal{C} = \mathbb{R}^d$ ) for a differentiable function that has an  $\mu$ -Lipschitz continuous gradient (i.e.,  $\mu$ -smooth), it is well known that SGD finds an  $\alpha$ -first-order stationary point (a point  $\mathbf{w}$  with  $\|\nabla F(\mathbf{w})\| \leq \alpha$ ) with  $O(\mu(F(\mathbf{w}_1) - F^*)/\alpha^2)$  iterations [Nesterov, 1998], where  $\mathbf{w}_1$  is the initial point and  $F^*$  is the optimal value of  $F$ . For constrained nonconvex smooth functions, a variety of similar dimension-free convergence results are known under the appropriate notion of stationarity [Ghadimi et al., 2016, Mokhtari et al., 2018]. These dimension-free convergence guarantees make SGD extremely attractive when the parameter space is very high dimensional.

### 1.1 OUR RESULTS

Our main result is a novel connection between constraint set structure and compression of the gradients.<sup>2</sup> We design a new SGD-style optimization algorithm that operates with

---

<sup>1</sup>Following a common convention, we still refer to the algorithm in this case as “gradient descent”.

<sup>2</sup>The term *compressed gradient* has been used in a variety of contexts in prior literature, here we use the term to denote a lower-dimensional representation of the gradient.

just the lower-dimensional (compressed) stochastic gradients and has the exact same dimension-free convergence guarantees as that of regular SGD for common convex and nonconvex optimization settings. Formally, instead of the usual stochastic gradient oracle, we assume the existence of a *compressed stochastic gradient* oracle, that on inputs  $\Phi \in \mathbb{R}^{m \times d}$  and  $\mathbf{w}_t \in \mathcal{C}$ , returns  $\vartheta_t = \Phi \hat{\mathbf{g}}_t \in \mathbb{R}^m$  where  $\hat{\mathbf{g}}_t$  is a stochastic subgradient of  $F$  at  $\mathbf{w}_t$ , and  $m$  could be much smaller than  $d$ . Our new SGD procedure (Algorithm COMP-SGD), that has access only to this compressed stochastic gradient oracle, first projects the current iterate  $\mathbf{w}_t$  onto the lower-dimensional space using  $\Phi$ , then updates the iterate here using the compressed stochastic gradient oracle, and then “lifts” back the result to  $\mathcal{C}$  to get the new iterate. We use ideas from geometry and high-dimensional estimation to perform this lifting. An immediate advantage of using compressed gradients (over regular SGD) comes in a distributed setup for reducing the cost of transmitting gradients.

We next investigate the convergence guarantees of this compressed SGD algorithm for various classes of functions. Our interest will be on  $\Phi$ ’s that are popularly referred to as *random projection* matrices such as subgaussian random matrices or sparse random matrices. Our analysis is based on exploiting the geometric structure of  $\mathcal{C}$ . We assume that in each iteration  $t$  of the algorithm, compressed stochastic gradient oracle is invoked using an *independent* random projection matrix  $\Phi_t$ , and provide conditions on the learning rate and the projection dimension under which the compressed SGD has the same (up to constant factors) dimension-free convergence guarantees as regular (uncompressed) SGD.

More precisely, let  $\Phi_t \in \mathbb{R}^{m_t \times d}$  be the random projection matrix used in iteration  $t$  with the compressed stochastic gradient oracle. We show that with appropriate setting of  $m_t$  one could match the regular SGD guarantees. The geometric parameter, *Gaussian width*, defined as  $\omega(\mathcal{C}) = \mathbb{E}_{\mathbf{r} \in \mathcal{N}(0,1)^d} [\sup_{\mathbf{w} \in \mathcal{C}} \langle \mathbf{w}, \mathbf{r} \rangle]$  plays an important role in our analysis, and shows up repeatedly as a geometric measure of the size of the set  $\mathcal{C}$ . Gaussian width roughly captures the expected width of  $\mathcal{C}$  along a random direction. Many constraint sets have a small Gaussian width, for example, a common choice of  $\mathcal{C}$  for sparsity purposes, the  $\ell_1$ -ball in  $\mathbb{R}^d$ , has a width of  $O(\sqrt{\log d})$ . Setting,  $m_t = \Omega(\omega(\mathcal{C})^2/\beta_t^2)$ , for a parameter  $\beta_t > 0$ , we establish the following bounds for the convergence guarantees of our compressed SGD algorithm over  $T$  iterations. For simplicity, we ignore dependence on other parameters such as variance of the stochastic gradient, the diameter of the convex set  $\mathcal{C}$ , etc.

**(a) Nonconvex Smooth Functions:** For a nonconvex function  $F$  that is  $\mu$ -smooth, we investigate two measures of (non)stationarity. Firstly, we show that a minibatch version of compressed SGD converges to  $\alpha$ -first-order sta-

tionary point ( $\alpha$ -FOSP)<sup>3</sup> in  $T = \Omega(\mu(F(\mathbf{w}_1) - F^*)/\alpha^2)$  iterations with  $\eta_t = 1/\mu$  and  $\beta_t \approx \alpha^2/\mu$  for all  $t \in [T]$ , and appropriate minibatch size (Theorem 2.4). We also investigate another popular measure of (non)stationarity defined through *gradient mapping*, which is a natural generalization of gradient for constrained problems [Nesterov, 1998, Ghadimi et al., 2016] (see Definition 8), and reach to a similar conclusion (Theorem B.5).

**(b) Convex and Strongly Convex Functions:** In this case, the goal is to bound the expected optimization error, defined as  $\mathbb{E}[F(\mathbf{w}_T) - F(\mathbf{w}^*)]$ , where  $\mathbf{w}^* \in \mathcal{C}$  is some minimizer of  $F$ . For a general convex function  $F$  (without any smoothness assumption), setting  $\beta_t = 1/t$  and  $\eta_t = 1/\sqrt{t}$ , we get an error bound of  $O(\log(T)/\sqrt{T})$  (Theorem 2.3). For a  $\lambda$ -strongly convex functions  $F$ , setting  $\beta_t = 1/t$  and  $\eta_t = 1/(\lambda t)$ , we get a  $O(\mu/(\lambda^2 T))$  bound on the expected error if  $F$  is  $\mu$ -smooth and a bound of  $O(\log T/(\lambda T))$  without smoothness (Theorem 2.3). These match the known regular SGD bounds, see, e.g., [Shamir and Zhang, 2013, Rakhlin et al., 2011].

These results demonstrate that for many problems, the compression of the gradients comes for “free”. The intuition behind this connection between constraint set structure and compression is as follows. Given a Gaussian random matrix and a finite set  $S$  of vectors, it is well-known (e.g., Johnson-Lindenstrauss Lemma) that projection with this matrix preserves the norms of all these vectors in  $S$  when the projection dimension is roughly  $\log(|S|)$  [Johnson and Lindenstrauss, 1984]. More generally, by measuring the complexity of the constraint set through the notion of Gaussian width, then the above compression guarantee can be achieved with a projection dimension that is roughly square of the Gaussian width of  $S$  [Gordon, 1988]. In our case, we compress the gradients (and not the iterates), and we show that by setting the projection dimension based on the square of the Gaussian width of the constraint set, suffices to well-approximate various norms that matter in convergence results (for example, the potential function defined as  $\|\mathbf{w}_t - \mathbf{w}\|$  for  $\mathbf{w} \in \mathcal{C}$ ).

To achieve this, we carefully combine ideas from modern SGD analyses, with ideas in convex geometry and high-dimensional estimation. One of main challenge comes in ensuring that the noise introduced by the random projection is always bounded. The compressed SGD algorithm updates the iterates in a lower-dimensional projected space, but we track the SGD progress in the original higher-dimensional space. To the best of our knowledge, these are the *first* rigorous results in SGD based on strictly utilizing *lower-dimensional* gradients.

**Applications.** We mention two applications of these com-

<sup>3</sup> $\bar{\mathbf{w}} \in \mathcal{C}$  is an  $\alpha$ -first-order stationary point if for all  $\mathbf{w} \in \mathcal{C}$ ,  $\langle \nabla F(\bar{\mathbf{w}}), \mathbf{w} - \bar{\mathbf{w}} \rangle \geq -\alpha$ ,  $\nabla F(\bar{\mathbf{w}})$  denotes the gradient of  $F$  at  $\bar{\mathbf{w}}$ .

pressed SGD results. Both these problems benefit from the reduced dimensionality of the gradients.

**(i) Differentially Private Empirical Risk Minimization**

**(ERM).** Machine learning algorithms are frequently run on sensitive data, and this has motivated the study of learning algorithms that have good performance guarantees while providing strong (mathematically proven) privacy protections for the training data. Differential Privacy (DP) is a formal algorithmic guarantee provides provable protection against adversaries with arbitrary side information and computational power, allows clear quantification of privacy losses, and satisfies graceful composition over multiple access to the same data [Dwork et al., 2006b]. We provide the first results in differentially private non-convex optimization for achieving first-order stationarity where the sample size  $n$  needed for a non-trivial result grows as  $\omega(\mathcal{C})$  and not as  $\sqrt{d}$  (see Table 1). As an example, if  $\mathcal{C}$  is the  $\ell_1$ -ball, then  $\omega(\mathcal{C}) = O(\sqrt{\log d})$ , and there is roughly an exponential improvement from  $\sqrt{d}$  to  $\sqrt{\log d}$  in the sample size needed compared to [Wang et al., 2017a, Zhang et al., 2017].

**(ii) Reducing Communication Costs in Distributed Synchronous SGD.**

Distributed stochastic gradient descent plays a very important role in distributed learning. In this work, we consider the data-distributed model of distributed SGD where the data sets are partitioned across various compute nodes. In each iteration of synchronous SGD, the compute nodes send their computed local gradients to a parameter server that averages and updates the global parameter. Since clients could be constrained devices that are on slow/expensive connections, communication cost is a principal constraint, especially, the cost of uploading gradients back to the server.

Firstly, our approach of utilizing lower-dimensional gradients already provides a way of reducing communication costs in many settings, without any change in the convergence rate. In addition, we show that combining this idea with any other gradient communication cost reduction scheme that provides unbiased estimates of the gradient can lead to even further communication cost savings, however now with an increase in convergence rate that only depends on the other chosen scheme. Our results provide the *first* communication cost bounds for various distributed optimization problems that depend on the geometry of  $\mathcal{C}$  (see Table 2).

**Extension to Conditional Gradient Method.** We show that the above-mentioned connection between constraint set and gradient compression extends beyond SGD to conditional gradient (Frank-Wolfe) method [Frank and Wolfe, 1956]. This method is a natural candidate for constrained optimization because of its projection free property and its ability to handle structured constraints [Jaggi, 2013]. We present a Frank-Wolfe style optimization algorithm that

utilizes the compressed stochastic gradient oracle for its gradient evaluations, and requires a linear minimization oracle over the set  $\Phi\mathcal{C}$ . As in the case of SGD, we show that a gradient dimension that is based on the square of the Gaussian width of  $\mathcal{C}$  suffices to obtain convergence bounds that match those of regular stochastic Frank-Wolfe method in convex/nonconvex settings (Theorems E.1 and E.2, Appendix E). Again, these are the first rigorous results for the Frank-Wolfe method based on strictly utilizing *lower-dimensional* gradients.

Due to space limitation, we leave many details, formal statements, proofs, and experimental studies in the supplement.

## 1.2 RELATED WORK

There has been a growing interest in understanding the convergence properties of SGD where instead of the true gradients some quantized/sparsified/sketched version of gradients is used in the SGD update step. The most common application of these techniques are in distributed/federated learning with the goal of reducing the communication costs, e.g., [Seide et al., 2014, Strom, 2015, De Sa et al., 2015, Konecny et al., 2016, Wen et al., 2017, Alistarh et al., 2017, Agarwal et al., 2018, Lin et al., 2017, Khirirat et al., 2018, Bernstein et al., 2018, Wu et al., 2018, Wang et al., 2018, Karimireddy et al., 2019, Stich et al., 2018, Mishchenko et al., 2019, Acharya et al., 2019, Alistarh et al., 2018, Ivkin et al., 2019, Gandikota et al., 2019, Horváth et al., 2019]. The general idea here is that a client could shrink the gradient communication cost by applying some encoding on the gradient before transmitting it to the server. Generally, these schemes come with an *increase* in the gradient variance. In other words, these schemes are generally “lossy” and will result in a slower convergence that using the true gradients. This is one major difference compared to our results, where we show that with the right set of parameters based on the geometry of the constraint set, our dimensionality-reduction scheme is “lossless”, in that we get the same convergence rate as if using the true gradients.

We also investigate whether our idea of utilizing lower-dimensional gradients can be combined with these existing gradient encoding techniques. Existing gradient encoding techniques can be categorized as either *unbiased* or *biased* based on whether the gradient estimates are unbiased or not. A number of gradient quantization methods are crafted specifically to yield unbiased estimates [Alistarh et al., 2017, Wen et al., 2017, Wang et al., 2018, Gandikota et al., 2019]. We show that one can combine our dimensionality-reduction scheme with *any* unbiased gradient encoding technique to get a reduction in communication cost at the expense of increased variance in gradient estimate.

We note that a number of gradient encoding techniques also produce biased gradient estimates [Seide et al., 2014,

Strom, 2015, Bernstein et al., 2018]. Using the idea of error-feedback, recently [Stich et al., 2018, Karimireddy et al., 2019] gave convergence guarantees for this kind of biased compression algorithm, showing that accumulating compression error locally in the workers can overcome the bias in the weight updates as long as the compression algorithm obeys certain properties. On a related note, gradient sparsification techniques with provable convergence were studied in [Alistarh et al., 2018, Stich et al., 2018]. Another idea proposed here is to use the *count-sketch* from the sketching literature for reducing communication [Ivkin et al., 2019, Rothchild et al., 2020]. An important point is that, in general, none of these prior encoding techniques inherently change the dimensionality of the gradient and also suffer from inflated variance at high compression rates. Also, in general it is hard to compare biased vs. unbiased approaches because of the different guarantees they are shown to provide. We leave the question of combining our dimensionality reduction idea with these biased gradient encoding schemes and/or stochastic variance reduction ideas [Alistarh et al., 2017, Horváth et al., 2019] for future work.

Another orthogonal idea to reduce communication, especially in the context of federated learning, is the idea of *local SGD* [Stich, 2019, Basu et al., 2019], where the clients perform local updates on their local data, and the clients communicate with the server only after few rounds.

**Beyond Gradient Descent.** It is well-known that for certain problems such as minimizing on the simplex with subgradients bounded in  $\ell_\infty$ -norm mirror descent outperforms SGD. Our focus here is on techniques that apply to a broad class of problems, including variety of constraint sets and non-convex functions. Since SGD is arguably the most popular and general optimization approach, we focus on it here.

Comparison with some additional related work based on sketching the Hessian matrix, differentially private ERM is presented in Appendix A.1.

### 1.3 PRELIMINARIES

**Notation.** We denote  $[n] = \{1, \dots, n\}$ . Vectors are denoted by boldface letters. We use  $\mathbb{I}_p$  to denote a  $p \times p$  identity matrix. For a vector  $\mathbf{v}$ ,  $\|\mathbf{v}\|$  denotes its Euclidean-norm. Throughout, we assume that  $\mathcal{C} \subseteq \mathbb{R}^d$  is a compact convex set. We define the diameter of  $\mathcal{C}$  as,  $\|\mathcal{C}\| = \sup_{\mathbf{w}, \mathbf{w}' \in \mathcal{C}} \|\mathbf{w} - \mathbf{w}'\|$ . Given a  $\Phi \in \mathbb{R}^{m \times d}$ , we define  $\Phi\mathcal{C} := \{\Phi\mathbf{w} : \mathbf{w} \in \mathcal{C}\}$ . Define the projection of  $\mathbf{w} \in \mathbb{R}^d$  onto a closed convex set  $\mathcal{C} \subseteq \mathbb{R}^d$  as:  $\Pi_{\mathcal{C}}(\mathbf{w}) := \operatorname{argmin}_{\mathbf{w}' \in \mathcal{C}} \|\mathbf{w} - \mathbf{w}'\|$ . We review some optimization fundamentals in Appendix A.2.

**Gaussian width.** Our analysis is based on exploiting the geometric properties of the constraint set. We use the well-studied quantity of *Gaussian width* that captures the  $\ell_2$ -geometric complexity of a set of points. Given a set  $S \subset \mathbb{R}^d$ , its Gaussian width  $\omega(S)$  is defined as:  $\omega(S) :=$

$$\mathbb{E}_{\mathbf{r} \in \mathcal{N}(0, \mathbb{I}_d)} [\sup_{\mathbf{a} \in S} \langle \mathbf{a}, \mathbf{r} \rangle].$$

Many interesting examples of  $S$  have low Gaussian width. For example, the  $\ell_1$ -ball and simplex in  $\mathbb{R}^d$  have width of  $O(\sqrt{\log d})$ , whereas a convex hull of  $l$  vectors with bounded  $\ell_2$ -norm of  $c$  has width  $O(c\sqrt{\log l})$ . See Table 3 in Appendix A.2 for additional examples of commonly used sets with low Gaussian width.

Given a set  $S$ , the square of Gaussian width of  $S$  measures the dimension size one needs to take on a random projection to still approximately preserve the norms of all the points in  $S$ . For Gaussian random matrices  $\Phi$ , this was first shown in [Gordon, 1988], popularly referred to as Gordon’s theorem (Theorem A.3 in Appendix A.2). This was recently extended to matrices drawn from subgaussian distributions [Dirksen, 2014] or distributions over sparse matrices [Bourgain et al., 2015] (Theorem A.4 in Appendix A.2). We will use this norm-preservation property and its consequences throughout our analyses.

## 2 SGD WITH LOW-DIM. GRADIENTS

In a regular SGD setup, the assumption is that we do not know  $F$ , and the only information is through a stochastic gradient oracle, which given  $\mathbf{w} \in \mathcal{C}$ , produces a vector  $\hat{\mathbf{g}}$  such that  $\mathbb{E}[\hat{\mathbf{g}}] = \mathbf{g}$  is a subgradient of  $F$  at  $\mathbf{w}$ . In this section, we introduce a dimensionality reduced gradient setup, where the stochastic gradient oracle does not return  $\hat{\mathbf{g}}$  but only a lower-dimensional projection of  $\hat{\mathbf{g}}$ , say obtained by applying a dimensionality reducing random projection to the subgradient. For a formal analysis, we define a compressed stochastic gradient oracle as follows.

**Definition 1** (Compressed Stochastic Subgradient Oracle (CSFO)). *Upon receiving query  $\mathbf{w}$  and  $\Phi \in \mathbb{R}^{m \times d}$ , the compressed stochastic gradient oracle returns  $\vartheta$  where  $\vartheta = \Phi\hat{\mathbf{g}}$ . Here,  $\hat{\mathbf{g}}$  is an independent random vector whose expectation  $\mathbb{E}[\hat{\mathbf{g}}] = \mathbf{g}$  is a subgradient of  $F$  at  $\mathbf{w}$ . Borrowing from the stochastic optimization literature, we denote this oracle as the compressed stochastic first-order oracle (CSFO), and use the notation  $\vartheta = \text{CSFO}(\mathbf{w}, \Phi)$ .*

Note that with  $\Phi = \mathbb{I}_d$ , CSFO defaults to the standard subgradient (SFO) oracle. Implementing a CSFO oracle is easy, as it just computes a projection of the subgradient. Hence, it can be efficiently implemented for any problem for which subgradients can be computed efficiently. Our convergence results, provide bounds on the number of calls needed to a CSFO oracle to get accurate results. Transmitting  $\Phi$  to the oracle is unnecessary as long as both the oracle and algorithm agree upon it. For example, as discussed in Section 3, in a distributed setup one can avoid transmitting  $\Phi$  by using shared randomness between clients and server.

Note that when  $m < d$ , the  $\vartheta = \text{CSFO}(\mathbf{w}, \Phi) \in \mathbb{R}^m$  has

a dimensionality lower than  $d$ .<sup>4</sup> This creates an immediate issue with using CSFO as there is now a dimensionality mismatch for the SGD iterate update step. We first design an SGD-style algorithm that overcomes this problem and operates with only access to CSFO.

Our algorithm (Algorithm COMPSGD) is based on a simple modification to the SGD procedure. Algorithm COMPSGD is parameterized by two functions  $\eta_t$  and  $\beta_t$ , where  $\eta_t$  describes the learning rate and  $\beta_t \in (0, 1)$  relates to the dimension of the compressed gradient. In each iteration  $t$ , the algorithm picks an independent random projection matrix  $\Phi_t$  of dimension  $m_t \times d$ . Various choices of  $\Phi_t$  could work here. For example,  $\Phi_t$  could be a Gaussian matrix with where each entry is sampled i.i.d. from  $N(0, 1/m_t)$  or  $\Phi_t$  could be a sparse JL matrix as defined by [Kane and Nelson, 2014]. Our approach can be informally described as: i) take an SGD step in the projected domain  $\mathbb{R}^{m_t}$  using the lower-dimensional subgradient, ii) do a projection onto  $\Phi\mathcal{C}$ , and iii) lift back the solution to  $\mathcal{C}$  to form the new iterate.

#### Algorithm COMPSGD

**Objective:**  $\min_{\mathbf{w} \in \mathcal{C}} F(\mathbf{w})$

**Input:** Convex set  $\mathcal{C}$ , learning rate parameters  $\{\eta_t\}$ , and projection dimension parameters  $\{\beta_t\}$ .

1. Pick  $\mathbf{w}_1$  as any point in  $\mathcal{C}$
2. for  $t = 1$  to  $T$  do
  - a. Set  $m_t \leftarrow \Omega(\min\{d, \omega(\mathcal{C})^2/\beta_t^2\})$
  - b. Let  $\Phi_t \in \mathbb{R}^{m_t \times d}$  be an i.i.d. random projection matrix (e.g., subgaussian or sparse JL matrix)
  - c. Let  $\vartheta_t \leftarrow \text{CSFO}(\mathbf{w}_t, \Phi_t)$  (from compressed stochastic gradient oracle)
  - d. Let  $\theta_t \leftarrow \Pi_{\Phi_t \mathcal{C}}(\Phi_t \mathbf{w}_t - \eta_t \vartheta_t)$
  - e. Let  $\mathbf{w}_{t+1} \leftarrow$  pick any element from the set  $\mathcal{S}_t = \{\mathbf{w} \in \mathcal{C} : \Phi_t \mathbf{w} = \theta_t\}$  (e.g., by solving (2))

Note that for any  $\Phi \in \mathbb{R}^{m \times d}$ ,  $\Phi\mathcal{C}$  is a convex set, and since  $\mathcal{C}$  is compact and closed, implies  $\Phi\mathcal{C}$  is closed. This follows because if  $\mathcal{C}$  is compact then the recess cone  $R_{\mathcal{C}} = \{0\}$  (Definition 6), which implies that  $\Phi\mathcal{C}$  is closed [Bertsekas, 2009]. Therefore, by properties of the projection operator  $\Pi_{\Phi\mathcal{C}}$  is well-defined and satisfies the projection properties in (5) and (6). For common  $\mathcal{C}$ 's of interest here such as convex hulls (polytopes), the linear transformation (through  $\Phi$ ) is also a convex hull (polytope) in  $\mathbb{R}^m$ , so the standard techniques for projection onto these sets are applicable for projection onto  $\Phi\mathcal{C}$ .

<sup>4</sup>Given  $\Phi \hat{\mathbf{g}} \in \mathbb{R}^m$  (with  $m \ll d$ ) it is not possible to recover  $\hat{\mathbf{g}}$  without further structural assumptions (such as sparsity) on  $\hat{\mathbf{g}}$ , which are generally not true.

In Algorithm COMPSGD,  $\mathbf{w}_{t+1}$  exists because  $\theta_t \in \Phi_t \mathcal{C} = \{\Phi_t \mathbf{w} : \mathbf{w} \in \mathcal{C}\}$  so we can represent  $\theta_t = \Phi_t \tilde{\mathbf{w}}$  with  $\tilde{\mathbf{w}} \in \mathcal{C}$ , therefore the set  $\mathcal{S}_t = \{\mathbf{w} \in \mathcal{C} : \Phi_t \mathbf{w} = \theta_t\}$  has at least one entry  $\tilde{\mathbf{w}}$ . For our theoretical results, it suffices that we pick any  $\mathbf{w}_{t+1} \in \mathcal{S}_t$ . We now discuss an idea for constructing  $\mathbf{w}_{t+1}$ , based on estimating  $\tilde{\mathbf{w}}$  which lies in  $\mathcal{C}$ , given  $\theta_t = \Phi_t \tilde{\mathbf{w}}$ . For any vector  $\mathbf{w} \in \mathbb{R}^d$ , the Minkowski functional of  $\mathcal{C} \subseteq \mathbb{R}^d$  is the non-negative number  $\|\mathbf{w}\|_{\mathcal{C}}$  defined by the rule:  $\|\mathbf{w}\|_{\mathcal{C}} = \inf\{\tau \in \mathbb{R} : \mathbf{w} \in \tau\mathcal{C}\}$ . Minkowski functional of a convex set is a convex function. In particular, when  $\mathcal{C}$  is a symmetric convex body, then  $\|\cdot\|_{\mathcal{C}}$  defines a norm (called Minkowski norm). Now consider the following optimization problem:

$$\text{Inv}_{\mathcal{C}}(\theta_t, \Phi_t) \stackrel{\text{def}}{=} \text{argmin}_{\mathbf{w}' \in \mathbb{R}^d} \|\mathbf{w}'\|_{\mathcal{C}} \text{ s.t. } \Phi_t \mathbf{w}' = \theta_t. \quad (2)$$

We now show that we can pick  $\mathbf{w}_{t+1}$  by solving (2).

**Lemma 2.1.** For any  $\mathbf{w}_{t+1} \in \text{Inv}_{\mathcal{C}}(\theta_t, \Phi_t)$ ,  $\mathbf{w}_{t+1} \in \mathcal{S}_t$ .

For any convex set  $\mathcal{C}$ , the optimization problem in (2) is a convex program, and hence can be solved using convex programming solvers. In fact, if  $\mathcal{C}$  is a polytope then this is a linear program. For example, if  $\mathcal{C}$  is the popular  $\ell_1$ -sparsity constraint, then (2) reduces to the well-studied *basis-pursuit* problem

$$\text{argmin}_{\mathbf{w}' \in \mathbb{R}^d} \|\mathbf{w}'\|_1 \text{ subject to } \Phi_t \mathbf{w}' = \theta_t. \quad (3)$$

for which lots of efficient solution techniques are known [Hastie et al., 2015]. In our experiments, we noticed that using sparse random matrices and known tricks for solving (3) (see, e.g., [Lorenz et al., 2015]), the per-iteration cost of Algorithm COMPSGD is almost identical to that of regular SGD. On the other hand, the reduction in gradient dimension could be quite substantial. Another point to note here is that in a distributed setup (as we discuss in Section 3.2), the selection of  $\mathbf{w}_{t+1}$  will be performed on the server, which is assumed to be computationally powerful, therefore this step will not be a matter of concern.

## 2.1 CONVERGENCE ANALYSIS

We now analyze the convergence guarantees of Algorithm COMPSGD. In all our results, the final expectation ( $\mathbb{E}[\cdot]$ ) is over the stochastic gradient noise and other randomness in the SGD algorithm (as is standard), but not over the randomness introduced by the random projection. The following inequality will play an important role in keeping track of the algorithm's progress.

**Lemma 2.2.** In Algorithm COMPSGD, for any  $t \in [T]$ ,  $(1 - \beta_t) \|\mathbf{w}_{t+1} - \mathbf{w}\|^2 \leq \|(\mathbf{w}_t - \eta_t \hat{\mathbf{g}}_t) - \mathbf{w}\|^2$ .

The challenge in our analysis comes from the fact that we have access to only to lower-dimensional gradients,

the updates happen in the compressed space, and the random projection introduces noise in the gradients. Our analyses will establish the values of  $\eta_t, \beta_t$  under which Algorithm COMPSGD has same (up to constant factors) dimension-free convergence guarantees as regular SGD. The analyses for the strongly convex and convex cases (with or without smoothness) are deferred to Appendix B.2 (Theorems B.7, B.8 and B.9). We get the following result.

**Theorem 2.3.** 1. [Strongly Convex, Smooth] Let  $F$  be a  $\lambda$ -strongly convex and  $\mu$ -smooth function over a convex set  $\mathcal{C}$ , and  $\mathbb{E}[\|\hat{\mathbf{g}}_t\|^2] \leq G^2$  for all  $t \in [T]$  (where  $\vartheta_t = \Phi_t \hat{\mathbf{g}}_t$ ). Then with  $\eta_t = 2/(\lambda t)$  and  $\beta_t = 1/t$ , Algorithm COMPSGD satisfies:  $\mathbb{E}[F(\mathbf{w}_T) - F(\mathbf{w}^*)] = O(\mu G^2/(\lambda^2 T))$ .

2. [Strongly Convex] Let  $F$  be a  $\lambda$ -strongly convex function over a convex set  $\mathcal{C}$ , and  $\mathbb{E}[\|\hat{\mathbf{g}}_t\|^2] \leq G^2$  for all  $t \in [T]$  (where  $\vartheta_t = \Phi_t \hat{\mathbf{g}}_t$ ). Then with  $\eta_t = 2/(\lambda t)$  and  $\beta_t = 1/t$ , Algorithm COMPSGD satisfies:  $\mathbb{E}[F(\mathbf{w}_T) - F(\mathbf{w}^*)] = O(G^2(1 + \log(T))/(\lambda T))$ .

3. [Convex] Let  $F$  be convex function over a convex set  $\mathcal{C}$ , and  $\mathbb{E}[\|\hat{\mathbf{g}}_t\|^2] \leq G^2$  for all  $t \in [T]$  (where  $\vartheta_t = \Phi_t \hat{\mathbf{g}}_t$ ). Then with  $\eta_t = \varphi/\sqrt{t}$  and  $\beta_t = 1/t$ , Algorithm COMPSGD satisfies:  $\mathbb{E}[F(\mathbf{w}_T) - F(\mathbf{w}^*)] = O((\|\mathcal{C}\|^2/\varphi + \varphi G^2)(\log T/\sqrt{T}))$ . In particular with  $\varphi = \|\mathcal{C}\|/G$ ,  $\mathbb{E}[F(\mathbf{w}_T) - F(\mathbf{w}^*)] = O(\|\mathcal{C}\|G \log T/\sqrt{T})$ .

**Convergence Analysis for Nonconvex Functions.** We now focus on the analysis for the nonconvex case. For nonconvex functions, in the constrained setting, several measures of (non)stationarity have been considered for projected gradient methods [Ghadimi et al., 2016, Mokhtari et al., 2018, Nouiehed et al., 2018]. We work with two popular stationarity notions here: a)  $\alpha$ -FOSP (Definition 7)<sup>3</sup> and b) norm of the gradient mapping (deferred to Appendix B.1.2). Let  $F^* = \min_{\mathbf{w} \in \mathcal{C}} F(\mathbf{w})$ . Missing details from this section are collected in Appendix B.1.

In this case, we assume a stronger stochastic approximation to the subgradient in that we use a minibatch version of Algorithm COMPSGD, where  $b$  is the minibatch size (formally described in Appendix B.1). Theorem 2.4 shows that after  $\approx O(\alpha^{-2})$  iterations, even with compressed gradients, minibatch Algorithm COMPSGD converges to an  $\alpha$ -FOSP (with high probability). The proof has two parts, first we show that when the difference between two consecutive iterates is small then we have an  $\alpha$ -FOSP, and then we bound the number of iterations before which this iterate difference condition is achieved. We show that the algorithm fails with a low probability to produce an  $\alpha$ -FOSP (in which case it outputs  $\perp$ ). We assume the following.

**Assumptions:**  $\mathbb{E}[\|\hat{\mathbf{g}}_t^{(i)}\|^2] \leq G^2$  and

$$\mathbb{E}[\|\nabla F(\mathbf{w}_t) - \hat{\mathbf{g}}_t^{(i)}\|^2] \leq \zeta^2, \quad \forall i \in [b], t \in [T]. \quad (4)$$

**Theorem 2.4.** Let  $F$  be  $\mu$ -smooth and continuously differentiable function over a convex set  $\mathcal{C}$ . Let the assumptions in (4) hold. Let  $\rho \in (0, 1)$  and  $\alpha > 0$ . Consider the minibatch version of Algorithm COMPSGD with output of first  $\mathbf{w}_\tau$  (if it exists) in  $(\mathbf{w}_1, \dots, \mathbf{w}_T)$  such that  $\|\mathbf{w}_\tau - \mathbf{w}_{\tau+1}\| \leq \alpha/(G + \mu\|\mathcal{C}\|)$ , and  $\perp$  (fail) otherwise. Set  $\eta_t = \eta = 1/\mu$ ,  $\beta_t = \beta = \min\{1/4, (\mu\alpha^2)/(64G\|\mathcal{C}\|(G + \mu\|\mathcal{C}\|)^2)\}$  for all  $t \in [T]$ , and batchsize  $b = \Omega((1/\rho) \cdot \max\{1, (\|\mathcal{C}\|\zeta(G + \mu\|\mathcal{C}\|)^2/(\alpha^2\mu))^2\})$ . Then, if  $T = \Omega((F(\mathbf{w}_1) - F^*)(G + \mu\|\mathcal{C}\|)^2/(\mu\alpha^2\rho))$ , with probability at least  $(1 - \rho)^2$ , this procedure outputs a  $\mathbf{w}_\tau$  that is an  $\alpha$ -FOSP for  $F$ .

One may notice that in the above theorem  $\beta_t$ 's are fixed across all iterations and scales as  $\approx \alpha^2$ , whereas in the convex settings we set  $\beta_t = 1/t$ . This distinction comes because the nature of guarantees differ in these two cases. In the convex case, as we get closer to the global optimum, we need a more precise estimate of the stochastic gradient (as compression adds noise), which is achieved by reducing  $\beta_t$ . In the nonconvex case, we can set  $\beta$  to a fixed value based on the required guarantee and there seems to be no advantage in adjusting it per iteration. We conjecture that for the convex cases, this dependence on  $\beta_t$  on  $t$  is unavoidable, if we want the corresponding convergence rates of Algorithm COMPSGD to match that of regular SGD.

The batchsize in Theorem 2.4 is carefully selected for our analysis to go through, and even with full stochastic gradients a similar batchsize would probably be needed, see for e.g., [Mokhtari et al., 2018] (note that the batchsize is independent of  $\beta_t$  and  $\omega(\mathcal{C})$ ). In practice, small batchsizes suffices as observed in our experiments.

## 2.2 EXPERIMENTAL RESULTS

We now experimentally compare Algorithm COMPSGD to SGD with the goal of validating our theoretical findings. Missing details about the datasets and experimental setup are provided in Appendix B.3. There we also have additional experimental results under other constraint sets such as subspace and (positive) simplex (Figure 5).

**Experimental Results on Convex Functions.** In Figure 1, we compare the performance of SGD and Algorithm COMPSGD on sparse linear regression [Hastie et al., 2015] and logistic regression with  $\ell_1$ -constraint [Lee et al., 2006]. We use a sparse random projection matrix in our experiments. We use the linear program in (3) to pick  $w_{t+1}$ , and use heuristic stopping ideas from [Lorenz et al., 2015] on an  $\ell_1$ -homotopy solver for efficiency. Based on the Gaussian width of the  $\ell_1$ -ball, we set  $m_t = \max\{d, t^2 \log(d)\}$

where  $d$  here is the dimensionality of the input. In Appendix B.3, we presents some results with other types of constraint sets.

With sparse linear regression, given a matrix  $A \in \mathbb{R}^{n \times d}$  and  $\mathbf{y} \in \mathbb{R}^n$  where  $\mathbf{y} = A\mathbf{w}^*$  with sparse  $\mathbf{w}^*$ , we solve:  $\min_{\mathbf{w}} \|\mathbf{y} - A\mathbf{w}\|^2$  subject to an  $\ell_1$ -constraint on  $\mathbf{w}$ . In Figure 1a, we plot the objective value (error)  $\|\mathbf{y} - A\mathbf{w}_t\|^2$  over SGD and COMPSGD iterations. We use synthetic random data for  $A$  with  $n = 1000, d = 10000$ . With logistic regression, we minimize the standard logistic regression objective but with an additional  $\ell_1$ -constraint [Lee et al., 2006].

The performance COMPSGD is almost identical (sometimes marginally better) than SGD. Computationally too Algorithm COMPSGD matches that of SGD. For example, on the MNIST, Reuters, and IMDB datasets, SGD took an average of 0.58, 6.36, 7.04 seconds per epoch respectively, whereas Algorithm COMPSGD took an average of 0.75, 6.76, 7.49 seconds per epoch respectively. So, for example on the IMDB dataset, CompSGD is only about 6.3% slower than SGD, however, as Figure 3d (Appendix B.3) shows the dimensionality of utilized gradients is significantly smaller, by a factor of  $\approx 19$  over the entire run. So, overall, Algorithm COMPSGD matches (or marginally improves) SGD performance, achieves significant gradient compression, with a minor increase in the computational cost. Training accuracy plots are in Figure 3 (Appendix B.3). The variance across runs for Algorithm COMPSGD is also low (see, e.g., Figure 4, Appendix B.3). In theory, this stems because of the tight concentration results we have with these random projections [Oymak et al., 2018, Theorem 1.3].

**Experimental Results on Nonconvex Functions.** We use an MLP with three hidden fully-connected layers (input dimension  $\times 50, 50 \times 50$ , and  $50 \times$  number of output classes) with ReLU activations, followed by softmax classifier. We optimize the network under an  $\ell_1$ -constraint on weights. Figure 2 presents the test accuracy plots for this network. For COMPSGD, we set  $m_t = d_i/10$  for each layer  $i$ , where  $d_i$  is the original number of parameters in layer  $i$ . The results show that convergence of SGD and COMPSGD are near identical, even though COMPSGD uses a *factor* of 10 lower-dimensional gradients. The training accuracy plots presented in Figure 6 (Appendix B.3) also exhibit a similar behavior.

### 3 APPLICATIONS

We now consider two different problems in which reducing the dimensionality of the gradient proves helpful.

#### 3.1 DIFF. PRIV. ERM WITH NONCONVEXITY

We consider the standard Empirical Risk Minimization (ERM) framework. Given a dataset  $D = (\mathbf{z}_1, \dots, \mathbf{z}_n)$  from a data universe  $\mathcal{D}^n$ , the goal in ERM is to:

$\min_{\mathbf{w} \in \mathcal{C}} F(\mathbf{w}; D) := \frac{1}{n} \sum_{i=1}^n f(\mathbf{w}; \mathbf{z}_i)$ , where  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is the loss of model  $\mathbf{w} \in \mathbb{R}^d$  for data record  $\mathbf{z}_i$ . Let us start with the definition of differential privacy. Let  $\mathcal{D}$  represent some domain. We say two datasets  $D \in \mathcal{D}^n$  and  $D' \in \mathcal{D}^n$  with  $n$  elements each are neighbors if they differ in one entry.

**Definition 2** ( $(\epsilon, \delta)$ -DP [Dwork et al., 2006b,a]). *A randomized algorithm  $\mathcal{A}$  is  $(\epsilon, \delta)$ -differentially private if for all neighboring datasets  $D, D'$  and for all outcomes  $\Gamma$  in the output space of  $\mathcal{A}$ , we have  $\Pr[\mathcal{A}(D) \in \Gamma] \leq \exp(\epsilon) \cdot \Pr[\mathcal{A}(D') \in \Gamma] + \delta$ , where the probability is taken over the randomness of the algorithm.*

Missing details and formal statements from this section are provided in Appendix C.

One of the basic ideas for achieving DP for ERM problems is add noise to the gradients using the (DP-SGD) algorithm [Song et al., 2013, Bassily et al., 2014]. The DP-SGD iteration is of the form,  $\theta_{t+1} \leftarrow \theta_t - \eta_t (\nabla F(\mathbf{w}; D) + \mathbf{e}_t)$ , where  $\mathbf{e}_t \sim \mathcal{N}(0, \sigma^2 \mathbb{I}_d)$  is the calibrated noise to achieve  $(\epsilon, \delta)$ -differential privacy. The fact that  $\nabla F(\mathbf{w}; D) \in \mathbb{R}^d$ , means that the popular *Gaussian mechanism* (Theorem A.1) idea in DP requires  $\mathbb{E}[\|\mathbf{e}_t\|^2] = \sigma^2 d$  that brings about the dependence on the dimension  $d$  in the utility analysis.

An approach to reduce the dependence on  $d$  arising from noise addition in DP-SGD iteration is to reduce the dimensionality of the gradient vector. For example, if we take  $\Phi \in \mathbb{R}^{m \times d}$  with entries drawn i.i.d from  $\mathcal{N}(0, 1/m)$ , then  $\Phi \nabla F(\mathbf{w}; D) \in \mathbb{R}^m$ . Therefore, by using the Gaussian mechanism now, one could add noise as:  $\Phi \nabla F(\mathbf{w}; D) + \mathbf{e}_t$  where  $\mathbf{e}_t \sim \mathcal{N}(0, \sigma^2 \mathbb{I}_m)$ . This roughly changes the dependence on  $d$  to  $m$  in the convergence analysis. This idea is formalized in Algorithm PRIVSGD (Appendix C). Table 1 summarizes the improved sample size bounds. These are the first results in differentially private nonconvex optimization which provides meaningful guarantees when  $n \gg \omega(\mathcal{C})$  rather than requiring  $n \gg \sqrt{d}$ . So for common  $\mathcal{C}$ 's, such as the  $\ell_1$ -ball, this reduction in sample size would be exponential from roughly  $d$  to  $\log d$ .

#### 3.2 REDUCING COMMUNICATION IN DIST. SGD

We consider the data-distributed model of distributed SGD. Let us assume that there are  $M$  clients, numbered  $1, \dots, M$ . Let  $F(\mathbf{w}) : \mathbb{R}^d \rightarrow \mathbb{R}$ . We investigate the following optimization problem:  $\min_{\mathbf{w} \in \mathcal{C}} F(\mathbf{w}) := \frac{1}{M} \sum_{i=1}^M f_i(\mathbf{w})$ , where each  $f_i$  resides at the  $i$ th client. As an illustration of the above setup, consider a machine learning problem, with data  $\{\mathbf{z}_1, \dots, \mathbf{z}_n\}$  partitioned on the clients, with  $P_i$  the set of indexes of datapoints on client  $i$ , then with  $f_i(\mathbf{w}) = (M/n) \sum_{j \in P_i} f(\mathbf{w}; \mathbf{z}_j)$ , we get  $F(\mathbf{w}) = (1/n) \sum_{i=1}^n f(\mathbf{w}; \mathbf{z}_i)$ . Missing details are provided in Appendix D.

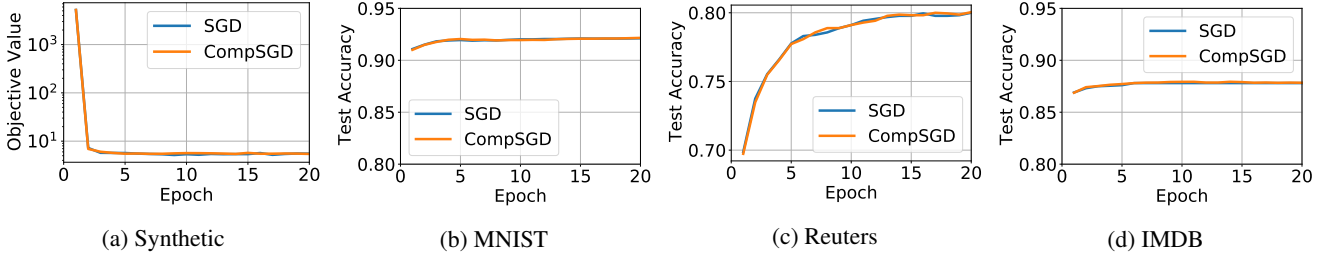


Figure 1: (a) Performance comparison for sparse linear regression on a synthetic dataset. (b), (c), and (d) Performance comparison for logistic regression with  $\ell_1$ -constraint on the MNIST dataset, Reuters dataset, and IMDB datasets, respectively. We use a constant learning rate at 0.1 based on its good performance for SGD and use minibatch size of 32.

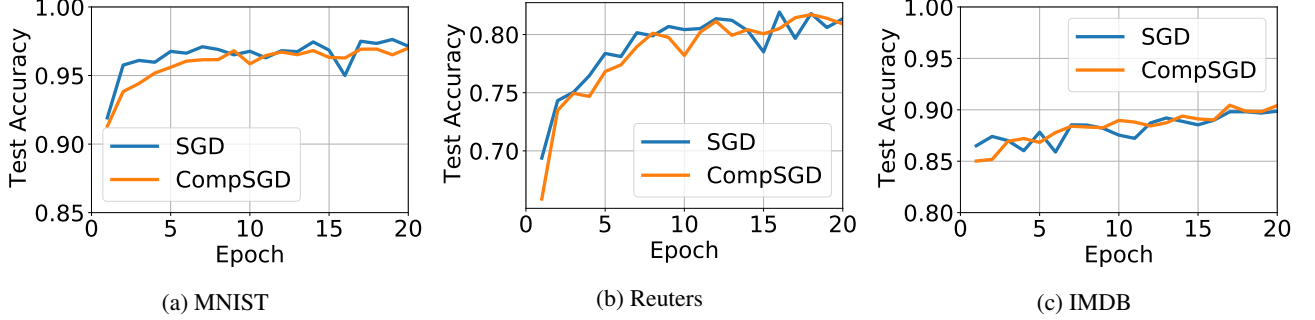


Figure 2: Performance comparison of training a MLP network with SGD vs. COMP-SGD on different datasets. Again, we use a constant learning rate at 0.1 based on its good performance for SGD and use minibatch size of 32.

Guarantee	DP-SGD Sample Size Req'd.	Algorithm PRIVSGD Sample Size Req'd.
$\alpha$ -FOSP	$\Omega\left(\sqrt{d} \cdot \frac{\ C\ (L+\mu\ C\ )^2 L\sqrt{T\log(1/\delta)}}{\sqrt{\rho\alpha^2\mu\epsilon}}\right)$	$\Omega\left(\min\left\{\frac{\omega(C)}{\beta}, \sqrt{d}\right\} \cdot \frac{\ C\ (L+\mu\ C\ )^2 L\sqrt{T\log(1/\delta)}}{\sqrt{\rho\alpha^2\mu\epsilon}}\right)$

Table 1: Comparison of sample size ( $n$ ) needed for achieving  $\alpha$ -FOSP for a nonconvex  $\mu$ -smooth function under  $(\epsilon, \delta)$ -DP. Assume  $\|\nabla f(\mathbf{w}; \cdot)\|$  is uniformly bounded by  $L$  for all  $\mathbf{w} \in \mathcal{C}$ . The settings of  $\beta$  and the upper bound on  $T$  are from Corollary C.2. Additional results are presented in Table 4, Appendix C.

Assumptions on $F$ and Guarantee	SGD (Total comm. cost)	SGD with Contraction $C$ (Total comm. cost)	Algorithm COMPDISTSGD (Total comm. cost)
Convex	$\tilde{O}\left(\frac{G^2\ C\ ^2 d}{\alpha^2}\right)$	$\tilde{O}\left(\frac{\chi_{C_d} G^2\ C\ ^2}{\alpha^2} \gamma_{C_d}\right)$	$\tilde{O}\left(\frac{\chi_{C_d} G^2\ C\ ^2}{\alpha^2} \min\{\gamma_{C_d}, \gamma_{C_r}\}\right)$
$\mathbb{E}[F(\mathbf{w}_T) - F(\mathbf{w}^*)] \leq \alpha$	Shamir and Zhang [2013]		(From Corollary D.2, Part 3)
$\mu$ -smooth and diff. nonconvex $\alpha$ -FOSP	$O(\kappa_1 T d)$	$O(\kappa T \gamma_{C_d})$	$O(\kappa T \min\{\gamma_{C_d}, \gamma_{C_r}\})$
	Mokhtari et al. [2018]		(From Corollary D.2, Part 4)

Table 2: Comparison of the total communication costs of distributed synchronous SGD, distributed synchronous SGD with contraction operator  $C$ , and our proposed Algorithm COMPDISTSGD. For convex case, we set  $\kappa = 1$  and  $r = (\omega(C)\chi_{C_d}G^2\|C\|^2/\alpha^2)^2$ . The  $\tilde{O}(\cdot)$  notation hides some logarithmic terms. The settings of  $\kappa, \beta, T$  for the nonconvex case is from Corollary D.2, Part 4 and  $r = \omega(C)^2/\beta^2$ . Here,  $\kappa_1$  is the value of  $\kappa$  obtained by setting  $\chi_{C_d} = 1$ . Additional results are presented in Table 5, Appendix D.

In each iteration  $t$  of synchronous SGD, the server randomly picks a set  $R_t$  of  $\kappa \geq 1$  clients and sends them the current model parameter  $\mathbf{w}_t$ . Each of these selected client  $i$  computes  $\hat{\mathbf{g}}_t^{(i)}$  an independent stochastic (sub)gradient of  $f_i$  at

$\mathbf{w}_t$ , and communicates  $\hat{\mathbf{g}}_t^{(i)}$  back to the server. The central server then aggregates these gradients and applies the update  $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \frac{\eta_t}{\kappa} \sum_{i \in R_t} \hat{\mathbf{g}}_t^{(i)}$ . Naively for the above protocol, the resulting per round communication is roughly  $\kappa \cdot (32d)$  bits (assuming 32-bit floating point numbers). Our



results (Theorems B.7, B.8, B.9, 2.4, and B.5) already show that, by transmitting lower-dimensional gradients, we can reduce communication costs, while preserving the convergence guarantees to within constant factor of the regular distributed SGD. In this section, we build on these results and show that in fact one can use any unbiased gradient compression scheme on top of our idea of utilizing lower-dimensional gradients.

For a formal analysis, we define a *contraction operator* to capture a general class of existing unbiased gradient encoding schemes. Following [Stich et al., 2018], we define a contraction operator as a (possibly randomized) function  $C_a$  defined as mapping from  $\mathbb{R}^a \rightarrow \mathbb{R}^a$  for  $a \in \mathbb{N}$  that satisfies these following common assumptions: (i)  $\forall \mathbf{v} \in \mathbb{R}^a, \mathbb{E}[C_a(\mathbf{v})] = \mathbf{v}$  (unbiasedness) and (ii)  $\mathbb{E}[\|C_a(\mathbf{v}) - \mathbf{v}\|^2] \leq \chi_{C_a} \|\mathbf{v}\|^2$  (variance bound). Let  $\gamma_{C_a}$  denote the bits needed to communicate  $C_a(\mathbf{v})$  from a client to server. This general notion captures many common unbiased quantization techniques such as *stochastic rounding* [Alistarh et al., 2017, Wen et al., 2017], vector quantization techniques such as *vqSGD* [Gandikota et al., 2019], and unbiased sparsification techniques such as *random sparsification* [Stich et al., 2018]. As an example, the quantization technique of [Alistarh et al., 2017] satisfies  $\chi_{C_a} = 1$  and  $\gamma_{C_a} \approx 2.8a + 32$  for all  $a \in \mathbb{N}$ .

Our procedure is presented in Algorithm COMPDISTSGD (Appendix D). The overall idea is simple, take a random projection of a gradient and then apply the contraction operator  $C_{m_t}$ . More formally, client  $i$  at iteration  $t$  will transmit  $C_{m_t}(\Phi_t \hat{\mathbf{g}}_t^{(i)})$ . The total communication cost in iteration  $t$  equals  $\kappa \cdot \gamma_{C_{m_t}}$  bits. This saves at least a factor  $\gamma_{C_d}/\gamma_{C_{m_t}}$  over just applying the contraction operator over the true gradients in each iteration  $t$ , which could be significant when  $m_t \ll d$ <sup>5</sup>. For example, if  $C$  is the quantization technique of [Alistarh et al., 2017], then  $\gamma_{C_d}/\gamma_{C_{m_t}} \approx d/m_t$  where  $m_t \approx \omega(\mathcal{C})^2/\beta_t^2$ , and for sets  $\mathcal{C}$  such as the unit  $l_1$ -ball,  $\omega(\mathcal{C})^2 = O(\log d)$ .

The cost of communicating  $\Phi_t$  is not significant as it can be achieved using various techniques such as one-to-all broadcasting. In practice,  $\Phi_t$  will be generated by a pseudorandom generator initialized by some seed, so by just communicating the seed we can regenerate  $\Phi_t$  at each device. In Table 2, we summarize the total communication cost (summed over all rounds) for a convex and nonconvex setting considered. Additional results and formal statements are provided in Appendix D.

In Figure 7, we provide experiments supporting our theoretical results. For illustration, we used the simple stochastic gradient quantization technique of [Alistarh et al., 2017] (described in Appendix D.1) for the contraction operator and assume just one client. Otherwise, we use the same

experimental setup as described in Section 2.2. Again, the main takeaway is that the performance of COMPDISTSGD matches that of SGD under the same contraction operator. The performance drop in both these schemes (compared to Figure 2 where we do not have any quantization) comes due to the applied quantization.

## 4 EXTENSIONS AND CONCLUSIONS

**Frank-Wolfe with Low-Dim. Gradients.** The Frank-Wolfe (FW) optimization algorithm requires access to a linear optimization oracle over  $\mathcal{C}$ . In Appendix E, we show how one could recover the standard convergence guarantees of Frank-Wolfe algorithm (for both convex and nonconvex functions) with only access to compressed gradients provided through the CSFO oracle (see Algorithm COMPFW). The main difference, compared to a traditional Frank-Wolfe algorithm, is how we invoke the linear optimization oracle. A traditional (stochastic) FW update computes a stochastic approximation to the gradient at the current iterate  $\mathbf{w}_t$  and invokes the LOO oracle with it. This gives the element in  $\mathcal{C}$  that correlates the most with the steepest descent (the negative stochastic gradient). We use a similar idea but instead solve a linear minimization problem with the compressed gradients over the set  $\Phi_t \mathcal{C}$ . We then utilize the lifting idea from (2) to compute a direction to take the step.

In Theorem E.2 (Appendix E), we establish that for convex functions the convergence guarantees of Algorithm COMPFW matches the known results with stochastic Frank-Wolfe algorithm [Hazan and Luo, 2016, Theorem 3]. In Theorem E.1 (Appendix E), we show that for smooth nonconvex functions Algorithm COMPFW, after  $O(\alpha^{-2})$  iterations, converges to an  $\alpha$ -FOSP (with high probability).

**Concluding Remarks.** We introduce the setting of SGD with compressed gradients, a fundamental question that studies how much bits of gradient information are truly needed for SGD to continue providing its guarantees, and also captures practical applications in private nonconvex ERM and distributed SGD. This new setup requires a rethinking of the SGD algorithm, and a subsequent careful analyses of the SGD guarantees. We also show that these ideas extend beyond SGD to the conditional gradient method. While we focused on getting bounds which hold in expectation it is possible that under assumptions on the tail of the noise distribution (such as [Harvey et al., 2019]) one could obtain high probability bounds. A natural question that arises from this work is whether the connection between constraint set structure and gradient compression that we observe here is inherent for *any first-order* optimization scheme. Extending the presented techniques to compress higher-order derivatives is an interesting open research direction.

<sup>5</sup> $\gamma_{C_a}$  is a non-decreasing function in  $a$ , i.e., cost of communicating a longer vector can't be less than one for a shorter vector.

## References

- Jayadev Acharya, Chris De Sa, Dylan Foster, and Karthik Sridharan. Distributed learning with sublinear communication. In *International Conference on Machine Learning*, pages 40–50, 2019.
- Naman Agarwal, Ananda Theertha Suresh, Felix Xinnan X Yu, Sanjiv Kumar, and Brendan McMahan. cpSGD: Communication-efficient and differentially-private distributed SGD. In *Advances in Neural Information Processing Systems*, pages 7564–7575, 2018.
- Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. QSGD: Communication-efficient SGD via gradient quantization and encoding. In *Advances in Neural Information Processing Systems*, pages 1709–1720, 2017.
- Dan Alistarh, Torsten Hoefer, Mikael Johansson, Nikola Konstantinov, Sarit Khirirat, and Cédric Renggli. The convergence of sparsified gradient methods. In *Advances in Neural Information Processing Systems*, pages 5973–5983, 2018.
- Raef Bassily, Adam Smith, and Abhradeep Thakurta. Differentially private empirical risk minimization: Efficient algorithms and tight error bounds. In *FOCS*. IEEE, 2014.
- Debraj Basu, Deepesh Data, Can Karakus, and Suhas Dig-gavi. Qsparse-local-sgd: Distributed sgd with quantization, sparsification, and local computations. *arXiv preprint arXiv:1906.02367 (Also in NeruIPS 2019)*, 2019.
- Jeremy Bernstein, Yu-Xiang Wang, Kamyar Azizzadenehsheli, and Animashree Anandkumar. SIGNSGD: Compressed optimisation for non-convex problems. In *International Conference on Machine Learning*, pages 559–568, 2018.
- Dimitri P Bertsekas. Nonlinear programming. *Journal of the Operational Research Society*, 48(3):334–334, 1997.
- Dimitri P Bertsekas. *Convex optimization theory*. Athena Scientific Belmont, 2009.
- Jean Bourgain, Dirksen Sjoerd, and Jelani Nelson. Toward a unified theory of sparse dimensionality reduction in euclidean space. In *Proceedings of the 47th ACM Symposium on Theory of Computing*. Association for Computing Machinery, 2015.
- Kamalika Chaudhuri and Claire Monteleoni. Privacy-preserving logistic regression. In *NIPS*, pages 289–296, 2009.
- Kamalika Chaudhuri, Claire Monteleoni, and Anand D Sarwate. Differentially private empirical risk minimization. *JMLR*, 12:1069–1109, 2011.
- Christopher M De Sa, Ce Zhang, Kunle Olukotun, and Christopher Ré. Taming the wild: A unified analysis of Hogwild!-style algorithms. In *Advances in neural information processing systems*, pages 2674–2682, 2015.
- Sjoerd Dirksen. Dimensionality reduction with sub-gaussian matrices: a unified theory. *arXiv preprint arXiv:1402.3973*, 2014.
- Cynthia Dwork, Krishnaram Kenthapadi, Frank McSherry, Ilya Mironov, and Moni Naor. Our data, ourselves: Privacy via distributed noise generation. In *EUROCRYPT*, LNCS, pages 486–503. Springer, 2006a.
- Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In *TCC*, volume 3876 of LNCS, pages 265–284. Springer, 2006b.
- Cynthia Dwork, Guy N Rothblum, and Salil Vadhan. Boosting and differential privacy. In *Foundations of Computer Science (FOCS), 2010 51st Annual IEEE Symposium on*, pages 51–60. IEEE, 2010.
- Cynthia Dwork, Aaron Roth, et al. The algorithmic foundations of differential privacy. *Foundations and Trends in Theoretical Computer Science*, 9(3-4):211–407, 2014.
- Marguerite Frank and Philip Wolfe. An algorithm for quadratic programming. *Naval research logistics quarterly*, 3(1-2):95–110, 1956.
- Venkata Gandikota, Raj Kumar Maity, and Arya Mazumdar. vqsgd: Vector quantized stochastic gradient descent. *arXiv preprint arXiv:1911.07971*, 2019.
- Saeed Ghadimi, Guanghui Lan, and Hongchao Zhang. Mini-batch stochastic approximation methods for nonconvex stochastic composite optimization. *Mathematical Programming*, 155(1-2):267–305, 2016.
- Yehoram Gordon. *On Milman’s inequality and random subspaces which escape through a mesh in  $\mathbb{R}^n$* . Springer, 1988.
- Nicholas JA Harvey, Christopher Liaw, Yaniv Plan, and Sikander Randhawa. Tight analyses for non-smooth stochastic gradient descent. In *ICML*, 2019.
- Trevor Hastie, Robert Tibshirani, and Martin Wainwright. *Statistical learning with sparsity: the lasso and generalizations*. CRC Press, 2015.
- Elad Hazan and Haipeng Luo. Variance-reduced and projection-free stochastic optimization. In *International Conference on Machine Learning*, pages 1263–1271, 2016.

- Samuel Horváth, Dmitry Kovalev, Konstantin Mishchenko, Sebastian Stich, and Peter Richtárik. Stochastic distributed learning with gradient quantization and variance reduction. *arXiv preprint arXiv:1904.05115*, 2019.
- Nikita Ivkin, Daniel Rothchild, Enayat Ullah, Ion Stoica, Raman Arora, et al. Communication-efficient distributed sgd with sketching. In *Advances in Neural Information Processing Systems*, pages 13144–13154, 2019.
- Martin Jaggi. Revisiting frank-wolfe: Projection-free sparse convex optimization. In *Proceedings of the 30th International Conference on Machine Learning (ICML-13)*, pages 427–435, 2013.
- Prateek Jain, Dheeraj Nagaraj, and Praneeth Netrapalli. Making the last iterate of sgd information theoretically optimal. In *Conference on Learning Theory*, pages 1752–1755, 2019.
- Peng Jiang and Gagan Agrawal. A linear speedup analysis of distributed deep learning with sparse and quantized communication. In *Advances in Neural Information Processing Systems*, pages 2525–2536, 2018.
- William B Johnson and Joram Lindenstrauss. Extensions of lipschitz mappings into a hilbert space. *Contemporary mathematics*, 26(189-206):1, 1984.
- Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Keith Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. Advances and open problems in federated learning. *arXiv preprint arXiv:1912.04977*, 2019.
- Daniel M Kane and Jelani Nelson. Sparser johnson-lindenstrauss transforms. *Journal of the ACM (JACM)*, 61(1):4, 2014.
- Sai Praneeth Karimireddy, Quentin Rebjock, Sebastian Stich, and Martin Jaggi. Error feedback fixes signsgd and other gradient compression schemes. In *International Conference on Machine Learning*, pages 3252–3261, 2019.
- Shiva Prasad Kasiviswanathan and Hongxia Jin. Efficient private empirical risk minimization for high-dimensional learning. In *International Conference on Machine Learning*, pages 488–497, 2016.
- Keras. Keras dataset, 2017. URL <https://keras.io/datasets/>.
- Sarit Khirirat, Hamid Reza Feyzmahdavian, and Mikael Johansson. Distributed learning with compressed gradients. *arXiv preprint arXiv:1806.06573*, 2018.
- Jakub Konecný, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh, and Dave Bacon. Federated learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492*, 2016.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Su-In Lee, Honglak Lee, Pieter Abbeel, and Andrew Y Ng. Efficient  $\ell_1$ -regularized logistic regression. In *AAAI*, volume 6, pages 401–408, 2006.
- Xiangru Lian, Wei Zhang, Ce Zhang, and Ji Liu. Asynchronous decentralized parallel stochastic gradient descent. In *International Conference on Machine Learning*, pages 3043–3052. PMLR, 2018.
- Yujun Lin, Song Han, Huizi Mao, Yu Wang, and William J Dally. Deep gradient compression: Reducing the communication bandwidth for distributed training. *arXiv preprint arXiv:1712.01887*, 2017.
- Dirk A Lorenz, Marc E Pfetsch, and Andreas M Tillmann. Solving basis pursuit: Heuristic optimality check and solver comparison. *ACM Transactions on Mathematical Software (TOMS)*, 41(2):1–29, 2015.
- Andrew L Maas, Raymond E Daly, Peter T Pham, Dan Huang, Andrew Y Ng, and Christopher Potts. Learning word vectors for sentiment analysis. In *Proceedings of the 49th annual meeting of the association for computational linguistics: Human language technologies-volume 1*, pages 142–150. Association for Computational Linguistics, 2011.
- Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial Intelligence and Statistics*, pages 1273–1282, 2017.
- Konstantin Mishchenko, Eduard Gorbunov, Martin Takáč, and Peter Richtárik. Distributed learning with compressed gradient differences. *arXiv preprint arXiv:1901.09269*, 2019.
- Aryan Mokhtari, Asuman Ozdaglar, and Ali Jadbabaie. Escaping saddle points in constrained optimization. In *Advances in Neural Information Processing Systems*, pages 3629–3639, 2018.
- Yurii Nesterov. Introductory lectures on convex programming volume i: Basic course. *Lecture notes*, 3(4):5, 1998.
- Maher Nouiehed, Jason D Lee, and Meisam Razaviyayn. Convergence to second-order stationarity for constrained non-convex optimization. *arXiv preprint arXiv:1810.02024*, 2018.
- Samet Oymak, Benjamin Recht, and Mahdi Soltanolkotabi. Isometric sketching of any set via the restricted isometry property. *Information and Inference: A Journal of the IMA*, 7(4):707–726, 2018.

- Mert Pilanci and Martin J Wainwright. Iterative hessian sketch: Fast and accurate solution approximation for constrained least-squares. *The Journal of Machine Learning Research*, 17(1):1842–1879, 2016.
- Mert Pilanci and Martin J Wainwright. Newton sketch: A near linear-time optimization algorithm with linear-quadratic convergence. *SIAM Journal on Optimization*, 27(1):205–245, 2017.
- Alexander Rakhlin, Ohad Shamir, and Karthik Sridharan. Making gradient descent optimal for strongly convex stochastic optimization. *arXiv preprint arXiv:1109.5647*, 2011.
- Daniel Rothchild, Ashwinee Panda, Enayat Ullah, Nikita Ivkin, Ion Stoica, Vladimir Braverman, Joseph Gonzalez, and Raman Arora. Fetchsgd: Communication-efficient federated learning with sketching. In *International Conference on Machine Learning*, pages 8253–8265. PMLR, 2020.
- Frank Seide, Hao Fu, Jasha Droppo, Gang Li, and Dong Yu. 1-bit stochastic gradient descent and its application to data-parallel distributed training of speech dnns. In *Fifteenth Annual Conference of the International Speech Communication Association*, 2014.
- Ohad Shamir and Tong Zhang. Stochastic gradient descent for non-smooth optimization: Convergence results and optimal averaging schemes. In *International Conference on Machine Learning*, pages 71–79, 2013.
- Shuang Song, Kamalika Chaudhuri, and Anand D Sarwate. Stochastic gradient descent with differentially private updates. In *Global Conference on Signal and Information Processing (GlobalSIP), 2013 IEEE*, pages 245–248. IEEE, 2013.
- Sebastian U Stich, Jean-Baptiste Cordonnier, and Martin Jaggi. Sparsified sgd with memory. In *Advances in Neural Information Processing Systems*, pages 4447–4458, 2018.
- Sebastian Urban Stich. Local sgd converges fast and communicates little. In *ICLR 2019 ICLR 2019 International Conference on Learning Representations*, 2019.
- Nikko Strom. Scalable distributed dnn training using commodity gpu cloud computing. In *Sixteenth Annual Conference of the International Speech Communication Association*, 2015.
- Kunal Talwar, Abhradeep Thakurta, and Li Zhang. Private empirical risk minimization beyond the worst case: The effect of the constraint set geometry. *arXiv preprint arXiv:1411.5417*, To appear in *NIPS*, 2015a.
- Kunal Talwar, Abhradeep Thakurta, and Li Zhang. Nearly optimal private lasso. In *NIPS*, 2015b.
- Di Wang, Minwei Ye, and Jinhui Xu. Differentially private empirical risk minimization revisited: Faster and more general. In *Advances in Neural Information Processing Systems*, pages 2722–2731, 2017a.
- Hongyi Wang, Scott Sievert, Shengchao Liu, Zachary Charles, Dimitris Papailiopoulos, and Stephen Wright. Atomo: Communication-efficient learning via atomic sparsification. In *Advances in Neural Information Processing Systems*, pages 9850–9861, 2018.
- Jialei Wang, Jason D Lee, Mehrdad Mahdavi, Mladen Kolar, Nathan Srebro, et al. Sketching meets random projection in the dual: A provable recovery algorithm for big and high-dimensional data. *Electronic Journal of Statistics*, 11(2):4896–4944, 2017b.
- 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 *Advances in neural information processing systems*, pages 1509–1519, 2017.
- Jiaxiang Wu, Weidong Huang, Junzhou Huang, and Tong Zhang. Error compensated quantized sgd and its applications to large-scale distributed optimization. In *International Conference on Machine Learning*, pages 5321–5329, 2018.
- Jiaqi Zhang, Kai Zheng, Wenlong Mou, and Liwei Wang. Efficient private erm for smooth objectives. *arXiv preprint arXiv:1703.09947*, 2017.
- Yuchen Zhang, John Duchi, Michael I Jordan, and Martin J Wainwright. Information-theoretic lower bounds for distributed statistical estimation with communication constraints. In *Advances in Neural Information Processing Systems*, pages 2328–2336, 2013.

---

# Appendix for “SGD with Low-Dimensional Gradients with Applications to Private and Distributed Learning”

---

Shiva Prasad Kasiviswanathan<sup>1</sup>

<sup>1</sup>Amazon, Palo Alto, USA

## A MISSING DETAILS FROM SECTION 1

### A.1 OTHER RELATED WORK

**Sketching Hessian.** Pilanci and Wainwright [2017], extending their previous work [Pilanci and Wainwright, 2016], proposed a technique for sketching the Hessian matrix square root by pre-multiplying it with a random matrix. As the column dimension is still  $d$ , the sketch matrix can be directly be used in Newton’s technique. The main technical difference compared to our setting is that our algorithms only sees a lower-dimensional summary (sketch) of the gradient which now due to dimension mismatch can’t be directed used in the SGD iteration. This requires a rethinking of the SGD algorithm. Also note that in practice first-order methods are preferred over second-order ones due to their simplicity and low per-iteration costs.

Wang et al. [2017b] presented sketching techniques for reducing the dimensionality of the data matrix for  $\ell_2$ -regularized least squares problem. Our focus is on general (non)convex optimization problems.

**Differentially Private ERM.** Starting with the results of [Chaudhuri and Monteleoni, 2009, Chaudhuri et al., 2011], differentially private convex ERM mechanisms have been investigated extensively in the literature under both  $\epsilon$ - and  $(\epsilon, \delta)$ -DP notions. Using carefully calibrated noisy gradients in first-order optimization algorithms is a common technique for achieving differential privacy. This framework, with the use of a SGD algorithm for optimization, was first introduced by [Song et al., 2013] and later [Bassily et al., 2014] analyzed the excess empirical risk bounds for various classes of convex functions. In particular, under  $(\epsilon, \delta)$ -DP, for general convex 1-Lipschitz functions [Bassily et al., 2014] showed the expected excess empirical risk is at most  $O(\sqrt{d}/n)$  on a dataset of size  $n$ .<sup>1</sup> They also showed that this bound cannot be improved in general, even for squared loss functions. Subsequently there has been a line of work that go beyond these worst-case bounds by exploiting properties of the constraint set  $\mathcal{C}$  [Talwar et al., 2015a,b, Kasiviswanathan and Jin, 2016, Wang et al., 2017a]. The general picture that emerges is that we can replace  $\sqrt{d}$  term in the excess risk by the Gaussian width of  $\mathcal{C}$ . These above results are based on convex optimization ideas. To the best of our knowledge, there are no known similar results that exploit the geometry of  $\mathcal{C}$ , when the function is nonconvex.

### A.2 MISSING PRELIMINARIES

**Convex Optimization Fundamentals.** We now review some basic concepts from convex optimization.

**Definition 3.** (*Lipschitz Functions*) A function  $F : \mathbb{R}^d \rightarrow \mathbb{R}$  is  $L$ -Lipschitz with respect to  $\mathbf{w}$  over the domain  $\mathcal{C}$ , if for all  $\mathbf{w}_a, \mathbf{w}_b \in \mathcal{C}$ , we have  $|F(\mathbf{w}_a) - F(\mathbf{w}_b)| \leq L\|\mathbf{w}_a - \mathbf{w}_b\|$ . If  $F$  is a convex function, then  $F$  is  $L$ -Lipschitz iff for all  $\mathbf{w} \in \mathcal{C}$  and subgradients  $\mathbf{g}$  of  $F$  at  $\mathbf{w}$  we have  $\|\mathbf{g}\| \leq L$ .

**Definition 4.** (*Strongly Convex Functions*) A function  $F : \mathbb{R}^d \rightarrow \mathbb{R}$  is  $\lambda$ -strongly convex if for all  $\mathbf{w}_a, \mathbf{w}_b \in \mathcal{C}$ , all

---

<sup>1</sup>Ignoring privacy parameters and some polylog factors. It is easy to see that these results are non-trivial, i.e., better than a trivial bound, only when  $n \gg \sqrt{d}$ .

subgradients  $\mathbf{g}$  of  $F(\mathbf{w}_a)$ , we have  $F(\mathbf{w}_b) \geq F(\mathbf{w}_a) + \langle \mathbf{g}, \mathbf{w}_b - \mathbf{w}_a \rangle + (\lambda/2)\|\mathbf{w}_b - \mathbf{w}_a\|^2$  (i.e.,  $F$  is bounded below by a quadratic function tangent at  $\mathbf{w}_a$ ).

Such functions arise, for instance, in Support Vector Machines and other regularized learning algorithms.

**Definition 5.** (Smooth Functions) A function  $F : \mathbb{R}^d \rightarrow \mathbb{R}$  is  $\mu$ -smooth if for all  $\mathbf{w}_a, \mathbf{w}_b \in \mathcal{C}$ , all subgradients  $\mathbf{g}$  of  $F(\mathbf{w}_a)$ , we have  $F(\mathbf{w}_b) \leq F(\mathbf{w}_a) + \langle \mathbf{g}, \mathbf{w}_b - \mathbf{w}_a \rangle + (\mu/2)\|\mathbf{w}_b - \mathbf{w}_a\|^2$  (i.e.,  $F$  is bounded above by a quadratic function tangent at  $\mathbf{w}_a$ ).

Note that sometime smoothness is stated slightly differently as a condition on Lipschitz continuity on the gradient. A differentiable function  $F$  with  $\mu$ -Lipschitz continuous gradient, i.e., which means for all  $\mathbf{w}_a, \mathbf{w}_b \in \mathcal{C}$ ,  $\|\nabla F(\mathbf{w}_a) - \nabla F(\mathbf{w}_b)\| \leq \mu\|\mathbf{w}_a - \mathbf{w}_b\|$ , is also  $\mu$ -smooth under the above definition.  $\mu$ -smooth functions arise, for instance, in logistic and least-squares regression.

**Definition 6** (Recess Cone). Given a nonempty convex set  $\mathcal{C}$ , a vector  $\mathbf{z}$  is a direction of recession if starting at any  $\mathbf{x}$  in  $\mathcal{C}$  and going indefinitely along  $\mathbf{z}$ , we never cross the relative boundary of  $\mathcal{C}$  to points outside  $\mathcal{C}$ :  $\mathbf{x} + \tau\mathbf{z} \in \mathcal{C}$ ,  $\forall \mathbf{x} \in \mathcal{C}$ ,  $\forall \tau \geq 0$ . The recess cone ( $R_{\mathcal{C}}$ ) of  $\mathcal{C}$  is the set of all directions of recession for  $\mathcal{C}$ .

**Unit Ball and Sphere.** The  $d$ -dimensional unit ball in  $\ell_p$ -norm centered at the origin is denoted by  $B_p^d$  and the unit sphere in  $\mathbb{R}^d$  centered at origin is denoted by  $S^{d-1}$ .

**Properties of Projection.** By the properties of the projection operator, for any  $\mathbf{w} \in \mathbb{R}^d$ ,  $\mathbf{w}_p = \Pi_{\mathcal{C}}(\mathbf{w})$  iff

$$\langle \mathbf{w} - \mathbf{w}_p, \tilde{\mathbf{w}} - \mathbf{w}_p \rangle \leq 0, \quad \forall \tilde{\mathbf{w}} \in \mathcal{C}. \quad (5)$$

We also use the standard *contractive* property of the projection operator that states,

$$\|\mathbf{w}_p - \tilde{\mathbf{w}}\| \leq \|\mathbf{w} - \tilde{\mathbf{w}}\|, \quad \forall \tilde{\mathbf{w}} \in \mathcal{C}. \quad (6)$$

**Differential Privacy Background.** Differential privacy is a rigorous notion of privacy that emerged from a long line of work in theoretical computer science [Dwork et al., 2006b]. Differential Privacy (DP) is a formal algorithmic guarantee that provides provable protection against adversaries with arbitrary side information and computational power, allows clear quantification of privacy losses, and satisfies graceful composition over multiple access to the same data. The literature on differential privacy is now rich with tools for constructing differentially private analyses, and we refer the reader to a survey by Dwork et al. [2014] for a comprehensive review of developments there.

We use the popular idea of adding Gaussian noise to achieve  $(\epsilon, \delta)$ -DP (commonly referred to as the *Gaussian mechanism*).

**Theorem A.1** (Gaussian Mechanism [Dwork et al., 2006b]). Let  $F : \mathcal{D}^n \rightarrow \mathbb{R}^d$ . The algorithm that on an input  $D$  outputs  $F(D) + \mathbf{e}$  where  $\mathbf{e} \sim \mathcal{N}(0, (2\Delta^2 \ln(2/\delta))/\epsilon^2)^d$  and  $\Delta$  denotes  $\ell_2$ -global sensitivity of the function  $F$  defined as  $\sup_{D, D'} \text{neighbors} \|F(D) - F(D')\|$ , is  $(\epsilon, \delta)$ -differentially private.

Composition theorems for differential privacy allow a modular design of privacy preserving algorithms based on algorithms for simpler subtasks.

**Theorem A.2** (Strong Composition of DP [Dwork et al., 2010]). Let  $\epsilon, \delta, \delta' > 0$  and  $\epsilon \leq 1$ . A mechanism that permits  $k$  adaptive interactions with mechanisms that preserves  $(\epsilon, \delta)$ -differential privacy ensures  $(\epsilon\sqrt{2k \ln(1/\delta')} + 2k\epsilon^2, k\delta + \delta')$ -differential privacy.

**Gordon's Theorem and Extensions.** Gordon's theorem [Gordon, 1988] can be viewed as a generalization of the Johnson-Lindenstrauss (JL) lemma for the case of Gaussian random matrices. The following is a simple restatement of the original theorem, better suited for this paper.

**Theorem A.3** (Gordon's Theorem [Gordon, 1988] Restated). Let  $\Phi \in \mathbb{R}^{m \times d}$  be a random matrix with independent  $\mathcal{N}(0, 1/m)$  entries. Let  $S \subset S^{d-1}$  be a subset of the unit sphere in  $d$  dimensions. Then if  $m = \Theta(\omega(S)^2/\beta^2)$ ,

$$\mathbb{E}_{\Phi} \left[ \sup_{\mathbf{x} \in S} \left| \|\Phi\mathbf{x}\|^2 - 1 \right| \right] \leq \beta, \quad (7)$$

where  $\omega(S)$  is the Gaussian width of  $S$  and the expectation  $\mathbb{E}_{\Phi}[\cdot]$  is over the randomness in  $\Phi$ .

Set	Gaussian Width
Unit $\ell_1$ -ball in $\mathbb{R}^d$	$O(\sqrt{\log d})$
Probability simplex in $\mathbb{R}^d$	$O(\sqrt{\log d})$
$\ell_p$ -ball in $\mathbb{R}^d$ for $1 < p \leq \infty$	$O(d^{1-1/p})$
Convex hull of $l$ vectors with bounded $\ell_2$ -norm of $c$	$O(c\sqrt{\log l})$
$m$ -dimensional subspace of $\mathbb{R}^d$	$O(\sqrt{m})$
$d_1 \times d_2$ matrices of rank at most $r$ and unit Frobenius norm	$O(\sqrt{r(d_1 + d_2)})$
$d \times d$ matrices with unit nuclear norm	$O(\sqrt{d})$

Table 3: Gaussian width of some popular constraint sets.

Gordon’s theorem was extended to i.i.d. subgaussian random matrices by [Dirksen, 2014]. Recently, Bourgain et al. [2015] extended Gordon’s theorem to distributions over sparse matrices having at most  $s$  non-zeroes per column. There are multiple constructions of such distributions discussed in [Kane and Nelson, 2014] and a matrix  $\Phi$  drawn from these distributions are referred to as a sparse Johnson-Lindenstrauss matrix.

**Theorem A.4** (Gordon’s Theorem with Sparse JL Matrices [Bourgain et al., 2015]). *Let  $\Phi \in \mathbb{R}^{m \times d}$  be a sparse Johnson-Lindenstrauss transform  $\Phi$  with column sparsity  $s$ . Then the guarantees of Theorem A.3 stated in (7) holds if  $m = \Omega(\text{polylog}(n)\omega(S)^2/\beta^2)$ ,  $s = \Omega(\text{polylog}(n)/\beta^2)$ , and  $(\log m)^2(\log n)^{5/2}v(S) \leq \beta$ , where*

$$v(S) = \max_{q \leq m/s \log s} \left\{ \frac{1}{\sqrt{qs}} \left( \mathbb{E}_\eta \left( \mathbb{E}_g \sup_{x \in S} \left| \sum_{j=1}^d \eta_j g_j x_j \right| \right)^q \right)^{1/q} \right\},$$

where  $(g_j)_{j \in [d]}$  are i.i.d. standard Gaussian and  $(\eta_j)_{j \in [d]}$  are i.i.d. Bernoulli with mean  $qs/(m \log s)$ , and  $\mathbf{x} = (x_1, \dots, x_d)$ .

Bourgain et al. [2015] also show how the complexity parameter  $v(S)$  can be controlled quite easily for all common  $S$  arising in applications. A simple consequence of these above theorems is the following corollary about inner-products.

**Corollary A.5.** *Consider the settings as in Theorem A.3 or Theorem A.4. Then for any  $\mathbf{x}, \mathbf{x}' \in S$ , we have  $\langle \mathbf{x}, \mathbf{x}' \rangle - \beta \leq \mathbb{E}_\Phi[\langle \Phi \mathbf{x}, \Phi \mathbf{x}' \rangle] \leq \langle \mathbf{x}, \mathbf{x}' \rangle + \beta$ .*

## B MISSING DETAILS FROM SECTION 2

The following lemma shows that any feasible solution to the linear estimation problem in (2) lies in  $\mathcal{S}_t$ .

**Lemma B.1** (Lemma 2.1 Restated). *For any  $\mathbf{w}_{t+1} \in \text{Inv}_{\mathcal{C}}(\theta_t, \Phi_t)$ ,  $\mathbf{w}_{t+1} \in \mathcal{S}_t$ .*

*Proof.* Let  $\theta_t = \Phi_t \tilde{\mathbf{w}}$  with  $\tilde{\mathbf{w}} \in \mathcal{C}$ , therefore there exists a feasible solution to (2). By definition of Minkowski functional,  $\mathcal{C} = \{\mathbf{w} \in \mathbb{R}^d : \|\mathbf{w}\|_{\mathcal{C}} \leq 1\}$ . Since  $\tilde{\mathbf{w}} \in \mathcal{C}$ , ensures that  $\|\tilde{\mathbf{w}}\|_{\mathcal{C}} \leq 1$ . By construction,  $\|\mathbf{w}_{t+1}\|_{\mathcal{C}} \leq \|\tilde{\mathbf{w}}\|_{\mathcal{C}} \leq 1$ , therefore, by  $\mathbf{w}_{t+1} \in \mathcal{C}$ . This along with the fact that by construction  $\Phi_t \mathbf{w}_{t+1} = \theta_t$  implies that  $\mathbf{w}_{t+1} \in \mathcal{S}_t$ .  $\square$

In all our results, the final expectation ( $\mathbb{E}[\cdot]$ ) is over the stochastic gradient noise and other randomness in the SGD algorithm (as is standard), but not over the randomness introduced by the random projection.  $E_{\Phi_t}[\cdot]$  denotes expectation just over  $\Phi_t$ .

We start by establishing an inequality that helps us keep track of the Algorithm COMPSGD’s progress.

**Lemma B.2** (Lemma 2.2 Restated). *In Algorithm COMPSGD, for any  $t \in [T]$ ,*

$$(1 - \beta_t) \|\mathbf{w}_{t+1} - \mathbf{w}\|^2 \leq \|(\mathbf{w}_t - \eta_t \hat{\mathbf{g}}_t) - \mathbf{w}\|^2.$$

*Proof.* Let  $\vartheta_t = \Phi_t \hat{\mathbf{g}}_t$  and  $\theta_t = \Pi_{\Phi_t \mathcal{C}}(\Phi_t \mathbf{w}_t - \eta_t \Phi_t \hat{\mathbf{g}}_t) = \Phi_t \hat{\mathbf{w}}_t$  (as  $\Pi_{\Phi_t \mathcal{C}}(\cdot)$  is a projection on  $\Phi \mathcal{C}$ ). By construction,  $\Phi_t \mathbf{w}_{t+1} = \Phi_t \hat{\mathbf{w}}_t$ . We use the following inequality based on Theorem A.3 (or Theorem A.4). Consider any  $\mathbf{w} \in \mathcal{C}$  (picked independent of  $\Phi_t$ ). We have  $\mathbf{w}_{t+1} - \mathbf{w} \in \mathcal{C} - \mathcal{C} \in 2\mathcal{C}$  (Minkowski difference). Let  $\hat{\mathcal{C}}$  denote the convex set  $2\mathcal{C}$ . Let us fix a

$\Phi_t$ . Define a function,  $u : \mathbb{R}^d \rightarrow \mathbb{R}^d$  as  $u(\mathbf{v}) = \mathbf{v}/\|\mathbf{v}\|$ , i.e.,  $u(\mathbf{v})$  is just the unit vector along  $\mathbf{v}$ . Let  $\hat{S} = \{u(\mathbf{w}) \mid \mathbf{w} \in \hat{\mathcal{C}}\}$ . Note that  $\hat{S} \subseteq S^{d-1}$  (unit sphere). Let  $u(\mathbf{w}_{t+1} - \mathbf{w})$  be unit vector along  $\mathbf{w}_{t+1} - \mathbf{w}$ . Since  $u(\mathbf{w}_{t+1} - \mathbf{w}) \in \hat{S}$ , we have

$$\|\Phi_t u(\mathbf{w}_{t+1} - \mathbf{w})\|^2 - 1 \leq \sup_{\mathbf{w} \in \hat{S}} \|\Phi_t \mathbf{w}\|^2 - 1.$$

Since this argument holds for every  $\Phi_t$ , we have,

$$\mathbb{E}_{\Phi_t} [\|\Phi_t u(\mathbf{w}_{t+1} - \mathbf{w})\|^2 - 1] \leq \mathbb{E}_{\Phi_t} [\sup_{\mathbf{w} \in \hat{S}} \|\Phi_t \mathbf{w}\|^2 - 1].$$

From (7) (Theorem A.3 or Theorem A.4),  $\mathbb{E}_{\Phi_t} [\sup_{\mathbf{w} \in \hat{S}} \|\Phi_t \mathbf{w}\|^2 - 1] \leq \beta_t$ , we get

$$\mathbb{E}_{\Phi_t} [\|\Phi_t u(\mathbf{w}_{t+1} - \mathbf{w})\|^2 - 1] \leq \beta_t \Rightarrow (1 - \beta_t) \leq \mathbb{E}_{\Phi_t} [\|\Phi_t u(\mathbf{w}_{t+1} - \mathbf{w})\|^2] \leq (1 + \beta_t).$$

Multiplying by  $\|\mathbf{w}_{t+1} - \mathbf{w}\|^2$ , we get,

$$\begin{aligned} (1 - \beta_t) \|\mathbf{w}_{t+1} - \mathbf{w}\|^2 &\leq \mathbb{E}_{\Phi_t} [\|\Phi_t(\mathbf{w}_{t+1} - \mathbf{w})\|^2] = \mathbb{E}_{\Phi_t} [\|\Phi_t \hat{\mathbf{w}}_t - \Phi_t \mathbf{w}\|^2] \\ &= \mathbb{E}_{\Phi_t} [\|\Pi_{\Phi_t \mathcal{C}}(\Phi_t \mathbf{w}_t - \eta_t \Phi_t \hat{\mathbf{g}}_t) - \Phi_t \mathbf{w}\|^2] \\ &= \mathbb{E}_{\Phi_t} [\|\Pi_{\Phi_t \mathcal{C}}(\Phi_t \mathbf{w}_t - \eta_t \Phi_t \hat{\mathbf{g}}_t) - \Pi_{\Phi_t \mathcal{C}}(\Phi_t \mathbf{w})\|^2] \\ &\leq \mathbb{E}_{\Phi_t} [\|(\Phi_t \mathbf{w}_t - \eta_t \Phi_t \hat{\mathbf{g}}_t) - (\Phi_t \mathbf{w})\|^2] \\ &= \|(\mathbf{w}_t - \eta_t \hat{\mathbf{g}}_t) - \mathbf{w}\|^2. \end{aligned}$$

The second last equality follows because the projection operation is contractive. The last equality follows because all  $\mathbf{w}_t, \eta_t, \hat{\mathbf{g}}_t, \mathbf{w}$  are all independent of  $\Phi_t$ .  $\square$

## B.1 CONVERGENCE ANALYSIS FOR NONCONVEX FUNCTIONS

In this section, we consider the optimization problem from  $\min_{\mathbf{w} \in \mathcal{C}} F(\mathbf{w})$  where  $F : \mathbb{R}^d \rightarrow \mathbb{R}$  is a smooth, continuously differentiable function over a compact convex  $\mathcal{C}$ , but  $F$  is potentially nonconvex. In this case, it is known that finding a global minimum (if it even exists) is hard. Given the well-known hardness results in finding stationary points, recent focus has shifted in characterizing approximate stationary points. In the constrained setting, several measures of (non)stationarity have been considered for projected gradient methods [Ghadimi et al., 2016, Mokhtari et al., 2018, Nouiehed et al., 2018]. We consider two popular stationarity notions here: a) first-order stationarity and b) stationarity defined through the norm of the gradient mapping. Let  $F^* = \min_{\mathbf{w} \in \mathcal{C}} F(\mathbf{w})$ .

### B.1.1 First-Order Stationary Point

We start with definition of first-order stationary point (FOSP). If the problem is unconstrained (i.e.,  $\mathcal{C} = \mathbb{R}^d$ ), then a point  $\bar{\mathbf{w}}$  is an  $\alpha$ -first-order stationary point (FOSP) if it satisfies  $\|\nabla F(\bar{\mathbf{w}})\| \leq \alpha$ . For the constrained setting, we use the following generalization.

**Definition 7** ( $\alpha$ -FOSP [Bertsekas, 1997]). *A point  $\bar{\mathbf{w}} \in \mathcal{C}$  is an  $\alpha$ -first-order stationary point ( $\alpha$ -FOSP) for a function  $F$  over a convex set  $\mathcal{C}$  if,  $\langle \nabla F(\bar{\mathbf{w}}), \mathbf{w} - \bar{\mathbf{w}} \rangle \geq -\alpha$  for all  $\mathbf{w} \in \mathcal{C}$ .*

For constrained optimization, Mokhtari et al. [2018] established an  $\approx O(\alpha^{-2})$  convergence rate to  $\alpha$ -FOSP for the techniques of projected gradient descent with true gradients and conditional descent with stochastic gradients.

In the following, we assume that  $F$  is smooth function over  $\mathcal{C}$  (see Definition 5). We now analyze the number of steps needed to recover an  $\alpha$ -FOSP with only access to CSFO. In this case, we assume a stronger stochastic approximation to the subgradient in that we use a minibatch version of Algorithm COMPSGD, where  $b$  is the minibatch size.



### Minibatch version of Algorithm COMPSGD

**Objective:**  $\min_{\mathbf{w} \in \mathcal{C}} F(\mathbf{w})$

**Input:** Convex set  $\mathcal{C}$ , learning rate parameters  $\{\eta_t\}$ , projection dimension parameters  $\{\beta_t\}$

1. Pick  $\mathbf{w}_1$  as any point in  $\mathcal{C}$
2. for  $t = 1$  to  $T$  do
  - a. Set  $m_t \leftarrow \Omega(\omega(\mathcal{C})^2/\beta_t^2)$
  - b. Let  $\Phi_t \in \mathbb{R}^{m_t \times d} \sim_{i.i.d} \mathcal{N}(0, 1/m_t)$
  - c. Let  $\bar{\vartheta}_t \leftarrow \frac{1}{b} \sum_{i=1}^b \vartheta_t^{(i)}$  (each  $\vartheta_t^{(i)}$  is an independent call to CSFO( $\mathbf{w}_t, \Phi_t$ ))
  - d. Let  $\theta_t \leftarrow \Pi_{\Phi_t \mathcal{C}}(\Phi_t \mathbf{w}_t - \eta_t \bar{\vartheta}_t)$
  - e. Let  $\mathbf{w}_{t+1} \leftarrow$  pick any element from the set  $\mathcal{S}_t = \{\mathbf{w} \in \mathcal{C} : \Phi_t \mathbf{w} = \theta_t\}$  (e.g., by solving (2))

**Remark B.3.** Since,  $\vartheta_t = (1/b) \sum_{i=1}^b \vartheta_t^{(i)}$ , we can re-express it as  $\vartheta_t = (1/b) \sum_{i=1}^b \Phi_t \hat{\mathbf{g}}_t^{(i)}$ , where  $\hat{\mathbf{g}}_t^{(i)}$  are independent stochastic gradients of  $F(\mathbf{w}_t)$ . Let  $\bar{\mathbf{g}}_t = (1/b) \sum_{i=1}^b \hat{\mathbf{g}}_t^{(i)}$ . Under this notation,  $\bar{\vartheta}_t = \Phi_t \bar{\mathbf{g}}_t$ .

We now consider the minibatch version of Algorithm COMPSGD executed till we encounter the first iteration  $T$  where  $\|\mathbf{w}_T - \mathbf{w}_{T+1}\|$  satisfies a certain bound. Theorem 2.4 shows that after  $\approx O(\alpha^{-2})$  iterations, even with compressed gradients, Algorithm COMPSGD converges to an  $\alpha$ -FOSP (with high probability). The proof has two parts, first we show that when the difference between two consecutive iterates is small then we have an  $\alpha$ -FOSP, and then we bound the number of iterations before which this iterate difference condition is achieved. We show that the algorithm fails with low probability to produce an  $\alpha$ -FOSP.

**Theorem B.4** (Theorem 2.4 Restated). *Let  $F$  be  $\mu$ -smooth and continuously differentiable function over a convex set  $\mathcal{C}$ . Let the assumptions in (4) hold. Let  $\rho \in (0, 1)$  and  $\alpha > 0$ . Consider the minibatch version of Algorithm COMPSGD with output of first  $\mathbf{w}_\tau$  (if it exists) in  $(\mathbf{w}_1, \dots, \mathbf{w}_T)$  such that  $\|\mathbf{w}_\tau - \mathbf{w}_{\tau+1}\| \leq \alpha/(G + \mu\|\mathcal{C}\|)$ , and  $\perp$  (fail) otherwise. Set  $\eta_t = \eta = 1/\mu$ ,*

$$\beta_t = \beta = \min \left\{ 1/4, \frac{\mu\alpha^2}{64G\|\mathcal{C}\|(G + \mu\|\mathcal{C}\|)^2} \right\}$$

for all  $t \in [T]$ , and batchsize  $b = \Omega((1/\rho) \cdot \max\{1, (\frac{\|\mathcal{C}\|\zeta(G + \mu\|\mathcal{C}\|)^2}{\alpha^2\mu})^2\})$ . Then, if  $T = \Omega(\frac{(F(\mathbf{w}_1) - F^*)(G + \mu\|\mathcal{C}\|)^2}{\mu\alpha^2\rho})$ , with probability at least  $(1 - \rho)^2$ , this procedure outputs a  $\mathbf{w}_\tau$  that is an  $\alpha$ -FOSP for  $F$ .

*Proof.* We establish the proof over two parts.

1. If  $\mathbf{w}_\tau$  exists, then with probability at least  $1 - \rho$ ,  $\mathbf{w}_\tau$  is an  $\alpha$ -FOSP.

2. If  $T = \Omega(\frac{(F(\mathbf{w}_1) - F^*)(G + \mu\|\mathcal{C}\|)^2}{\mu\alpha^2\rho})$ , then with probability at least  $1 - \rho$ ,  $\mathbf{w}_\tau$  exists in  $(\mathbf{w}_1, \dots, \mathbf{w}_T)$ .

Note that under our assumptions,  $\|\nabla F(\mathbf{w}_t)\| \leq G$  for all  $t \in [T]$ .

Let us start with the first part. Let  $\mathbf{g}_t = \nabla F(\mathbf{w}_t)$  and  $\bar{\mathbf{g}}_t = \frac{1}{b} \sum_{i=1}^b \hat{\mathbf{g}}_t^{(i)}$  (see Remark B.3) for any  $t \in [T]$ . Note that  $\Phi_t \bar{\mathbf{g}}_t = \bar{\vartheta}_t$ . We use the following observation that for any  $\mathbf{w} \in \mathcal{C}$ ,

$$\begin{aligned} \langle \Phi_\tau \mathbf{g}_\tau, \Phi_\tau(\mathbf{w} - \mathbf{w}_\tau) \rangle &= \langle \Phi_\tau \mathbf{g}_\tau, \Phi_\tau(\mathbf{w}_{\tau+1} - \mathbf{w}_\tau) \rangle + \langle \Phi_\tau \mathbf{g}_\tau, \Phi_\tau(\mathbf{w} - \mathbf{w}_{\tau+1}) \rangle \\ &\geq -\|\Phi_\tau \mathbf{g}_\tau\| \|\Phi_\tau(\mathbf{w}_{\tau+1} - \mathbf{w}_\tau)\| + \langle \Phi_\tau \mathbf{g}_\tau, \Phi_\tau(\mathbf{w} - \mathbf{w}_{\tau+1}) \rangle. \end{aligned}$$

Now,

$$\begin{aligned} \langle \Phi_\tau \mathbf{g}_\tau, \Phi_\tau(\mathbf{w} - \mathbf{w}_{\tau+1}) \rangle &= \langle \Phi_\tau \bar{\mathbf{g}}_\tau, \Phi_\tau(\mathbf{w} - \mathbf{w}_{\tau+1}) \rangle + \langle \Phi_\tau \mathbf{g}_\tau - \Phi_\tau \bar{\mathbf{g}}_\tau, \Phi_\tau(\mathbf{w} - \mathbf{w}_{\tau+1}) \rangle \\ &\geq \langle \Phi_\tau \bar{\mathbf{g}}_\tau, \Phi_\tau(\mathbf{w} - \mathbf{w}_{\tau+1}) \rangle - \|\Phi_\tau \mathbf{g}_\tau - \Phi_\tau \bar{\mathbf{g}}_\tau\| \|\Phi_\tau(\mathbf{w} - \mathbf{w}_{\tau+1})\|. \end{aligned}$$

From properties of the projection operator, we also have

$$\begin{aligned} \langle \Phi_\tau \mathbf{w}_\tau - \eta_\tau \Phi_\tau \bar{\mathbf{g}}_\tau - \Phi_\tau \mathbf{w}_{\tau+1}, \Phi_\tau \mathbf{w} - \Phi_\tau \mathbf{w}_{\tau+1} \rangle &\leq 0 \\ \implies \langle \Phi_\tau \bar{\mathbf{g}}_\tau, \Phi_\tau \mathbf{w} - \Phi_\tau \mathbf{w}_{\tau+1} \rangle &\geq \frac{1}{\eta_\tau} \langle \Phi_\tau \mathbf{w}_\tau - \Phi_\tau \mathbf{w}_{\tau+1}, \Phi_\tau \mathbf{w} - \Phi_\tau \mathbf{w}_{\tau+1} \rangle. \end{aligned}$$

Putting the above expressions together,

$$\begin{aligned}
& \langle \Phi_\tau \mathbf{g}_\tau, \Phi_\tau(\mathbf{w} - \mathbf{w}_\tau) \rangle \\
& \geq \frac{1}{\eta_\tau} \langle \Phi_\tau \mathbf{w}_\tau - \Phi_\tau \mathbf{w}_{\tau+1}, \Phi_\tau \mathbf{w} - \Phi_\tau \mathbf{w}_{\tau+1} \rangle - \|\Phi_\tau \mathbf{g}_\tau\| \|\Phi_\tau(\mathbf{w}_{\tau+1} - \mathbf{w}_\tau)\| - \|\Phi_\tau \mathbf{g}_\tau - \Phi_\tau \bar{\mathbf{g}}_\tau\| \|\Phi_\tau(\mathbf{w} - \mathbf{w}_{\tau+1})\| \\
& \geq -\frac{1}{\eta_\tau} \|\Phi_\tau \mathbf{w}_\tau - \Phi_\tau \mathbf{w}_{\tau+1}\| \|\Phi_\tau \mathbf{w} - \Phi_\tau \mathbf{w}_{\tau+1}\| - \|\Phi_\tau \mathbf{g}_\tau\| \|\Phi_\tau(\mathbf{w}_{\tau+1} - \mathbf{w}_\tau)\| - \|\Phi_\tau \mathbf{g}_\tau - \Phi_\tau \bar{\mathbf{g}}_\tau\| \|\Phi_\tau(\mathbf{w} - \mathbf{w}_{\tau+1})\|.
\end{aligned}$$

Taking expectation with respect to  $\Phi_\tau$  on both sides.

$$\langle \mathbf{g}_\tau, \mathbf{w} - \mathbf{w}_\tau \rangle = \mathbb{E}_{\Phi_\tau}[\langle \Phi_\tau \mathbf{g}_\tau, \Phi_\tau(\mathbf{w} - \mathbf{w}_\tau) \rangle] \geq -(1 + \beta)(\mu \|\mathcal{C}\| \|\mathbf{w}_\tau - \mathbf{w}_{\tau+1}\| + G \|\mathbf{w}_\tau - \mathbf{w}_{\tau+1}\| + \|\mathcal{C}\| \|\mathbf{g}_\tau - \bar{\mathbf{g}}_\tau\|).$$

This implies that

$$\min_{\mathbf{w} \in \mathcal{C}} \langle \mathbf{g}_\tau, \mathbf{w} - \mathbf{w}_\tau \rangle \geq -(1 + \beta)(\mu \|\mathcal{C}\| \|\mathbf{w}_\tau - \mathbf{w}_{\tau+1}\| + G \|\mathbf{w}_\tau - \mathbf{w}_{\tau+1}\| + \|\mathcal{C}\| \|\mathbf{g}_\tau - \bar{\mathbf{g}}_\tau\|). \quad (8)$$

Since  $\mathbb{E}[\|\mathbf{g}_\tau - \bar{\mathbf{g}}_\tau\|^2] \leq \zeta^2/b$ , we obtain from Markov's inequality<sup>2</sup>

$$\Pr[\|\mathbf{g}_\tau - \bar{\mathbf{g}}_\tau\| \leq \alpha/\|\mathcal{C}\|] \geq 1 - \frac{\zeta^2 \|\mathcal{C}\|^2}{b\alpha^2}.$$

Using the assumption  $\|\mathbf{w}_\tau - \mathbf{w}_{\tau+1}\| \leq \frac{\alpha}{G + \mu \|\mathcal{C}\|}$  in (8) and the above probability bound, we get that with probability at least  $1 - \frac{\zeta^2 \|\mathcal{C}\|^2}{b\alpha^2}$ ,

$$\min_{\mathbf{w} \in \mathcal{C}} \langle \mathbf{g}_\tau, \mathbf{w} - \mathbf{w}_\tau \rangle \geq -c''\alpha,$$

for some constant  $c''$ . By rescaling  $\alpha$  to  $\alpha/c''$  and using our lower bound on batchsize  $b$ , we establish the  $\alpha$ -FOSP property under the condition that  $\|\mathbf{w}_\tau - \mathbf{w}_{\tau+1}\| \leq \alpha/(G + \mu \|\mathcal{C}\|)$ .

Let us look at the second part which bounds the expected number of iterations  $\mathcal{T}$ , starting from any  $\mathbf{w}_1 \in \mathcal{C}$ , before reaching the first  $\mathbf{w}_\mathcal{T}$  such that  $\|\mathbf{w}_\mathcal{T} - \mathbf{w}_{\mathcal{T}+1}\| \leq \alpha/(G + \mu \|\mathcal{C}\|)$ . Here we use a proof idea based on Theorem 2 of [Mokhtari et al., 2018]. Consider  $t \leq \mathcal{T}$ . By our assumption,  $\|\mathbf{w}_{t+1} - \mathbf{w}_t\| > \alpha/(G + \mu \|\mathcal{C}\|)$ . We start with the  $\mu$ -smooth condition which implies that,

$$F(\mathbf{w}_{t+1}) - F(\mathbf{w}_t) - \langle \mathbf{g}_t, \mathbf{w}_{t+1} - \mathbf{w}_t \rangle \leq \frac{\mu}{2} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|^2.$$

We have

$$F(\mathbf{w}_{t+1}) - F(\mathbf{w}_t) \leq \langle \bar{\mathbf{g}}_t, \mathbf{w}_{t+1} - \mathbf{w}_t \rangle - \langle \bar{\mathbf{g}}_t - \mathbf{g}_t, \mathbf{w}_{t+1} - \mathbf{w}_t \rangle + \frac{\mu}{2} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|^2.$$

Considering the  $\langle \bar{\mathbf{g}}_t, \mathbf{w}_{t+1} - \mathbf{w}_t \rangle$  term. By Corollary A.5,

$$\langle \bar{\mathbf{g}}_t, \mathbf{w}_{t+1} - \mathbf{w}_t \rangle \leq \mathbb{E}_{\Phi_t}[\langle \Phi_t \bar{\mathbf{g}}_t, \Phi_t(\mathbf{w}_{t+1} - \mathbf{w}_t) \rangle] + \beta \|\bar{\mathbf{g}}_t\| \|\mathbf{w}_{t+1} - \mathbf{w}_t\|. \quad (9)$$

Let  $\Phi_t \tilde{\mathbf{w}}_t = \Phi_t(\mathbf{w}_t - \eta \bar{\mathbf{g}}_t)$ . As before, let  $\Phi_t \hat{\mathbf{w}}_t = \vartheta_t = \Pi_{\Phi_t \mathcal{C}}(\Phi_t \mathbf{w}_t - \eta \Phi_t \bar{\mathbf{g}}_t) = \Pi_{\Phi_t \mathcal{C}}(\Phi_t \tilde{\mathbf{w}}_t)$ . By definition,  $\Phi_t \mathbf{w}_{t+1} = \Phi_t \hat{\mathbf{w}}_t = \Pi_{\Phi_t \mathcal{C}}(\Phi_t \tilde{\mathbf{w}}_t)$ . Therefore, by properties of the projection operator,

$$\begin{aligned}
\langle \Phi_t \bar{\mathbf{g}}_t, \Phi_t(\mathbf{w}_{t+1} - \mathbf{w}_t) \rangle &= \frac{1}{\eta} \langle \Phi_t \mathbf{w}_t - \Phi_t \tilde{\mathbf{w}}_t, \Phi_t \mathbf{w}_{t+1} - \Phi_t \mathbf{w}_t \rangle \\
&\leq \frac{1}{\eta} \langle \Phi_t \mathbf{w}_t - \Phi_t \mathbf{w}_{t+1}, \Phi_t \mathbf{w}_{t+1} - \Phi_t \mathbf{w}_t \rangle = -\frac{1}{\eta} \|\Phi_t \mathbf{w}_{t+1} - \Phi_t \mathbf{w}_t\|^2.
\end{aligned}$$

Using this along with Lemma 2.2 (with  $\mathbf{w} = \mathbf{w}_t$ ), we have

$$\mathbb{E}_{\Phi_t}[\langle \Phi_t \bar{\mathbf{g}}_t, \Phi_t(\mathbf{w}_{t+1} - \mathbf{w}_t) \rangle] \leq -\frac{(1 - \beta)}{\eta} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|^2.$$

Therefore, we have

$$\langle \bar{\mathbf{g}}_t, \mathbf{w}_{t+1} - \mathbf{w}_t \rangle \leq -\frac{(1 - \beta)}{\eta} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|^2 + \beta \|\bar{\mathbf{g}}_t\| \|\mathbf{w}_{t+1} - \mathbf{w}_t\|.$$

<sup>2</sup>We will use the Markov's higher moment inequality: for a random variable  $x$  and  $a \geq 1$ ,  $\Pr[|x| > t] \leq \mathbb{E}[|x|^a]/t^a$ .

Since  $\beta \leq 1/4$  and  $\eta = 1/\mu$ , we get

$$\begin{aligned} F(\mathbf{w}_{t+1}) - F(\mathbf{w}_t) &\leq -\frac{\mu}{4} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|^2 + \beta \|\bar{\mathbf{g}}_t\| \|\mathbf{w}_{t+1} - \mathbf{w}_t\| - \langle \bar{\mathbf{g}}_t - \mathbf{g}_t, \mathbf{w}_{t+1} - \mathbf{w}_t \rangle \\ &\leq -\frac{\mu}{4} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|^2 + \beta \|\bar{\mathbf{g}}_t\| \|\mathbf{w}_{t+1} - \mathbf{w}_t\| + \|\bar{\mathbf{g}}_t - \mathbf{g}_t\| \|\mathbf{w}_{t+1} - \mathbf{w}_t\|. \end{aligned}$$

Note that  $\|\mathbf{w}_{t+1} - \mathbf{w}_t\| \geq \alpha/(G + \mu\|\mathcal{C}\|)$ . Also note that  $\|\mathbf{w}_{t+1} - \mathbf{w}_t\| \leq \|\mathcal{C}\|$ .

$$F(\mathbf{w}_{t+1}) - F(\mathbf{w}_t) \leq -\frac{\mu\alpha^2}{4(G + \mu\|\mathcal{C}\|)^2} + \beta \|\bar{\mathbf{g}}_t\| \|\mathcal{C}\| + \|\bar{\mathbf{g}}_t - \mathbf{g}_t\| \|\mathcal{C}\|.$$

Consider  $\mathcal{F}_t$  as the sigma algebra that measures all sources of randomness up to iteration  $t$ . Then taking expectation on both sides,

$$\begin{aligned} \mathbb{E}[F(\mathbf{w}_{t+1}) | \mathcal{F}_t] &\leq F(\mathbf{w}_t) - \frac{\mu\alpha^2}{4(G + \mu\|\mathcal{C}\|)^2} + \beta \mathbb{E}[\|\bar{\mathbf{g}}_t\|] \|\mathcal{C}\| + \|\mathcal{C}\| \mathbb{E}[\|\bar{\mathbf{g}}_t - \mathbf{g}_t\|] \\ &\leq F(\mathbf{w}_t) - \frac{\mu\alpha^2}{4(G + \mu\|\mathcal{C}\|)^2} + \beta \|\mathcal{C}\| G + \frac{\zeta}{\sqrt{b}} \|\mathcal{C}\|. \end{aligned}$$

Under our setting of  $b$  and  $\beta$  we get,

$$\mathbb{E}[F(\mathbf{w}_{t+1}) | \mathcal{F}_t] \leq F(\mathbf{w}_t) - \frac{\mu\alpha^2}{8(G + \mu\|\mathcal{C}\|)^2}. \quad (10)$$

Let us take a look at the expected value of  $\mathbb{E}[F(\mathbf{w}_1) - F(\mathbf{w}_{\mathcal{T}+1})]$

$$\begin{aligned} \mathbb{E}[F(\mathbf{w}_1) - F(\mathbf{w}_{\mathcal{T}+1})] &= \mathbb{E} \left[ \sum_{t=1}^{\mathcal{T}} (F(\mathbf{w}_t) - F(\mathbf{w}_{t+1})) \right] \\ &= \mathbb{E} \left[ \mathbb{E} \left[ \sum_{t=1}^{\mathcal{T}} (F(\mathbf{w}_t) - F(\mathbf{w}_{t+1})) \mid \mathcal{T} = k \right] \right] = \mathbb{E} \left[ \mathbb{E} \left[ \sum_{t=1}^k (F(\mathbf{w}_t) - F(\mathbf{w}_{t+1})) \mid \mathcal{T} = k \right] \right] \\ &= \sum_{k=1}^{\infty} \mathbb{E} \left[ \sum_{t=1}^k (F(\mathbf{w}_t) - F(\mathbf{w}_{t+1})) \right] \Pr[\mathcal{T} = k] = \sum_{k=1}^{\infty} \sum_{t=1}^k \mathbb{E}[(F(\mathbf{w}_t) - F(\mathbf{w}_{t+1}))] \Pr[\mathcal{T} = k] \\ &\geq \sum_{k=1}^{\infty} \sum_{t=1}^k \frac{\mu\alpha^2}{8(G + \mu\|\mathcal{C}\|)^2} \Pr[\mathcal{T} = k] = \frac{\mu\alpha^2}{8(G + \mu\|\mathcal{C}\|)^2} \sum_{k=1}^{\infty} k \Pr[\mathcal{T} = k] \\ &= \frac{\mu\alpha^2}{8(G + \mu\|\mathcal{C}\|)^2} \mathbb{E}[\mathcal{T}]. \end{aligned}$$

This implies that  $\mathbb{E}[\mathcal{T}] \leq \mathbb{E}[F(\mathbf{w}_1) - F(\mathbf{w}_{\mathcal{T}+1})]/(8\mu\alpha^2/(G + \mu\|\mathcal{C}\|)^2)$ . Since  $F^* \leq F(\mathbf{w}_{\mathcal{T}+1})$ , we get that

$$\mathbb{E}[\mathcal{T}] \leq (F(\mathbf{w}_1) - F^*) \frac{8(G + \mu\|\mathcal{C}\|)^2}{\mu\alpha^2}.$$

Using Markov's inequality

$$\Pr \left[ \mathcal{T} \leq \frac{8(F(\mathbf{w}_1) - F^*)(G + \mu\|\mathcal{C}\|)^2}{\mu\alpha^2 \rho} \right] \geq 1 - \rho.$$

This completes both parts of the theorem. Putting them together gives the claimed result.  $\square$

### B.1.2 Stationarity based on Gradient Mapping

Nesterov [Nesterov, 1998] introduced the notion of gradient mapping in the context of constrained optimization, where the gradients required to be treated differently as compared with unconstrained optimization. In particular, in constrained minimization problems gradient mapping is known to preserve the most important properties of the gradient. We refer the reader to [Nesterov, 1998] for a thorough investigation of these connections. Following Ghadimi et al. [2016], for a learning

rate  $\eta$ , we define the mapping at  $\mathbf{w}$  as:  $\frac{1}{\eta}(\mathbf{w} - \hat{\mathbf{w}})$  where  $\hat{\mathbf{w}} := \Pi_{\mathcal{C}}(\mathbf{w} - \eta \nabla F(\mathbf{w}))$ . To get an intuition about gradient mapping, notice that If we let  $\eta \rightarrow 0$ , then the gradient mapping becomes simply the negative of the projection of  $-\nabla F(\mathbf{w})$  on the solid tangent cone to  $\mathcal{C}$  at  $\mathbf{w}$ . When  $\mathcal{C} = \mathbb{R}^d$  (unconstrained setting), then the gradient mapping becomes simply  $\nabla F(\mathbf{w})$ . In our setting, with only access to compressed (lower-dimensional) gradients, we never perform a projection on  $\mathcal{C}$  directly, so we modify the above definition of gradient mapping to that of compressed gradient mapping.

**Definition 8** (Compressed Gradient Mapping). *Given a continuously differentiable function  $F : \mathbb{R}^d \rightarrow \mathbb{R}$ , the compressed gradient mapping at  $\mathbf{w}$ , given compressed gradient  $\Phi \nabla F(\mathbf{w})$  for  $\Phi \in \mathbb{R}^{m \times d}$ , is defined as*

$$\frac{1}{\eta}(\mathbf{w} - \mathbf{w}^+) \text{ where } \Phi \mathbf{w}^+ := \Pi_{\Phi \mathcal{C}}(\Phi \mathbf{w} - \eta \Phi \nabla F(\mathbf{w})).$$

We say  $\mathbf{w}$  is an  $\alpha$ -compressed gradient mapping for  $\alpha > 0$ , if  $\|\frac{1}{\eta}(\mathbf{w} - \mathbf{w}^+)\| \leq \alpha$ .

When  $\Phi \sim_{i.i.d} \mathcal{N}(0, 1/m)$  and  $\mathcal{C} = \mathbb{R}^d$ , then with high probability  $\Phi \mathcal{C} = \mathbb{R}^m$ , and  $\frac{1}{\eta}(\Phi \mathbf{w} - \Phi \mathbf{w}^+) = \Phi \nabla F(\mathbf{w})$ , and  $\|\frac{1}{\eta}(\mathbf{w} - \mathbf{w}^+)\|^2 \approx \mathbb{E}_{\Phi}[\|\frac{1}{\eta}(\Phi \mathbf{w} - \Phi \mathbf{w}^+)\|^2] = \|\nabla F(\mathbf{w})\|^2$ . Therefore, compressed gradient mapping provides an interpretation that is similar to gradient mapping.

Again, we consider the minibatch version of Algorithm COMPSGD defined above. Our goal is to analyze the number of steps needed to achieve an  $\alpha$ -compressed gradient mapping in expectation. The following theorem establishes a rate of convergence for this quantity. In particular, we see that convergence rate is again  $\approx O(\alpha^{-2})$  as is typical with SGD stationarity analyses.

**Theorem B.5.** *Let  $F$  be  $\mu$ -smooth and continuously differentiable function over a convex set  $\mathcal{C}$ . Let the assumptions in (4) hold. Let  $\alpha > 0$ . Consider the minibatch version of Algorithm COMPSGD that outputs a uniformly at random iterate  $\mathbf{w}_k$  from  $\{\mathbf{w}_1, \dots, \mathbf{w}_T\}$ . Set  $\eta_t = \eta = 1/(2\mu)$ ,  $\beta_t = \beta = \min\{1/4, \frac{\alpha^2}{(\mu G \|\mathcal{C}\|)}\}$  for all  $t \in [T]$ , and batchsize  $b = \Omega(\max\{1, \frac{\zeta^2}{\alpha^2}, \frac{\zeta^2 \|\mathcal{C}\|^2 \mu^2}{\alpha^4}\})$ . Then, if  $T = \Omega(\frac{\mu(F(\mathbf{w}_1) - F^*)}{\alpha^2})$ ,  $\mathbf{w}_k$  satisfies*

$$\mathbb{E}_{\mathcal{R}} \left[ \left\| \frac{1}{\eta}(\mathbf{w}_k - \mathbf{w}_k^+) \right\| \right] \leq \alpha \text{ where } \Phi_k \mathbf{w}_k^+ := \Pi_{\Phi_k \mathcal{C}}(\Phi_k \mathbf{w}_k - \eta \Phi_k \nabla F(\mathbf{w}_k)).$$

where the expectation is over the choice of  $k$  and the stochastic gradient oracle noise (together denoted by  $\mathcal{R}$ ).

*Proof.* Let  $\mathbf{g}_t = \nabla F(\mathbf{w}_t)$  and  $\delta_t = \bar{\mathbf{g}}_t - \mathbf{g}_t$ , where  $\bar{\mathbf{g}}_t = \frac{1}{b} \sum_{i=1}^b \hat{\mathbf{g}}_t^{(i)}$  (see Remark B.3). Define  $\mathbf{r}_t$  at iteration  $t$  as:

$$\mathbf{r}_t = \frac{1}{\eta}(\mathbf{w}_t - \mathbf{w}_{t+1}).$$

We have

$$\begin{aligned} F(\mathbf{w}_{t+1}) - F(\mathbf{w}_t) &\leq \langle \mathbf{g}_t, \mathbf{w}_{t+1} - \mathbf{w}_t \rangle + \frac{\mu}{2} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|^2 \\ &= \langle \bar{\mathbf{g}}_t, \mathbf{w}_{t+1} - \mathbf{w}_t \rangle - \langle \delta_t, \mathbf{w}_{t+1} - \mathbf{w}_t \rangle + \frac{\mu}{2} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|^2. \end{aligned}$$

As in the proof of Theorem 2.4, it is easy to see that

$$\langle \bar{\mathbf{g}}_t, \mathbf{w}_{t+1} - \mathbf{w}_t \rangle \leq -\frac{(1-\beta)}{\eta} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|^2 + \beta \|\bar{\mathbf{g}}_t\| \|\mathbf{w}_{t+1} - \mathbf{w}_t\|.$$

Hence, we have

$$\begin{aligned} F(\mathbf{w}_{t+1}) - F(\mathbf{w}_t) &\leq -\frac{(1-\beta)}{\eta} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|^2 + \beta \|\bar{\mathbf{g}}_t\| \|\mathbf{w}_{t+1} - \mathbf{w}_t\| - \langle \delta_t, \mathbf{w}_{t+1} - \mathbf{w}_t \rangle + \frac{\mu}{2} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|^2 \\ &\leq \left( \frac{-3\eta}{4} + \frac{\mu\eta^2}{2} \right) \|\mathbf{r}_t\|^2 + \eta \langle \delta_t, \mathbf{r}_t \rangle + \beta \|\bar{\mathbf{g}}_t\| \|\mathcal{C}\|. \end{aligned} \tag{11}$$

Note that

$$\Phi_t \mathbf{r}_t = \frac{1}{\eta}(\Phi_t \mathbf{w}_t - \Phi_t \mathbf{w}_{t+1}) = \frac{1}{\eta}(\Phi_t \mathbf{w}_t - \Pi_{\Phi_t \mathcal{C}}(\Phi_t \mathbf{w}_t - \eta \bar{\mathbf{g}}_t)).$$

Similarly define,

$$\Phi_t \hat{\mathbf{r}}_t = \frac{1}{\eta} (\Phi_t \mathbf{w}_t - \Pi_{\Phi_t \mathcal{C}} (\Phi_t \mathbf{w}_t - \eta \mathbf{g}_t)).$$

It is not hard to see that there exists a  $\hat{\mathbf{r}}_t \in \mathcal{C} - \mathcal{C}$ . Therefore, (11) can be re-expressed as

$$\begin{aligned} F(\mathbf{w}_{t+1}) - F(\mathbf{w}_t) &= \left( \frac{-3\eta}{4} + \frac{\mu\eta^2}{2} \right) \|\mathbf{r}_t\|^2 + \eta \langle \delta_t, \hat{\mathbf{r}}_t \rangle + \eta \langle \delta_t, \mathbf{r}_t - \hat{\mathbf{r}}_t \rangle + \beta \|\bar{\mathbf{g}}_t\| \|\mathcal{C}\| \\ &\leq \left( \frac{-3\eta}{4} + \frac{\mu\eta^2}{2} \right) \|\mathbf{r}_t\|^2 + \eta \langle \delta_t, \hat{\mathbf{r}}_t \rangle + \eta \|\delta_t\| \|\mathbf{r}_t - \hat{\mathbf{r}}_t\| + \beta \|\bar{\mathbf{g}}_t\| \|\mathcal{C}\|. \end{aligned}$$

Now  $\|\mathbf{r}_t - \hat{\mathbf{r}}_t\|$  can be bounded as,

$$\|\mathbf{r}_t - \hat{\mathbf{r}}_t\| \leq \frac{1}{\sqrt{1-\beta}} \mathbb{E}_{\Phi_t} [\|\Phi_t(\mathbf{r}_t - \hat{\mathbf{r}}_t)\|] \frac{1}{\sqrt{1-\beta}} \mathbb{E}_{\Phi_t} [\|\Phi_t \delta_t\|] \leq \frac{\sqrt{1+\beta}}{\sqrt{1-\beta}} \|\delta_t\| \leq c'' \|\delta_t\|.$$

for some constant  $c''$ . Substituting this bound on  $\|\mathbf{r}_t - \hat{\mathbf{r}}_t\|$ , we get

$$F(\mathbf{w}_{t+1}) - F(\mathbf{w}_t) \leq \left( \frac{-3\eta}{4} + \frac{\mu\eta^2}{2} \right) \|\mathbf{r}_t\|^2 + \eta \langle \delta_t, \mathbf{r}_t \rangle + c'' \eta \|\delta_t\|^2 + \beta \|\bar{\mathbf{g}}_t\| \|\mathcal{C}\|.$$

Following the idea in [Ghadimi et al., 2016, Theorem 2], we get starting from a similar to above condition, for a  $k$  drawn uniformly at random from  $[T]$

$$\mathbb{E}_{\mathcal{R}} \left[ \frac{1}{\eta^2} \|(\mathbf{w}_k - \mathbf{w}_{k+1})\|^2 \right] = O \left( \frac{\mu(F(\mathbf{w}_1) - F^*)}{T} + \frac{\mu\zeta\|\mathcal{C}\|}{\sqrt{b}} + \frac{\zeta^2}{b} + \mu\beta G\|\mathcal{C}\| \right). \quad (12)$$

Therefore,

$$\begin{aligned} &\mathbb{E}_{\Phi_k, \mathcal{R}} \left[ \left\| \frac{1}{\eta} (\Phi_k \mathbf{w}_k - \Phi_k \mathbf{w}_k^+) \right\|^2 \right] = \mathbb{E}_{\Phi_k, \mathcal{R}} \left[ \frac{1}{\eta^2} \|\Phi_k \mathbf{w}_k - \Pi_{\Phi_k \mathcal{C}} (\Phi_k \mathbf{w}_k - \eta \Phi_k \mathbf{g}_k)\|^2 \right] \\ &\leq \frac{1}{\eta^2} (2 \mathbb{E}_{\Phi_k, \mathcal{R}} [\|\Phi_k \mathbf{w}_k - \Pi_{\Phi_k \mathcal{C}} (\Phi_k \mathbf{w}_k - \eta \Phi_k \bar{\mathbf{g}}_k)\|^2] + 2 \mathbb{E}_{\Phi_k, \mathcal{R}} [\|\Pi_{\Phi_k \mathcal{C}} (\Phi_k \mathbf{w}_k - \eta \Phi_k \mathbf{g}_k) - \Pi_{\Phi_k \mathcal{C}} (\Phi_k \mathbf{w}_k - \eta \Phi_k \bar{\mathbf{g}}_k)\|^2]) \\ &\leq 2 \mathbb{E}_{\Phi_k, \mathcal{R}} \left[ \frac{1}{\eta^2} \|\Phi_k \mathbf{w}_k - \Phi_k \mathbf{w}_{k+1}\|^2 \right] + 2 \mathbb{E}_{\Phi_k, \mathcal{R}} [\|\Phi_k \bar{\mathbf{g}}_k - \Phi_k \mathbf{g}_k\|^2] \\ &= O \left( \mathbb{E}_{\mathcal{R}} \left[ \frac{1}{\eta^2} \|\mathbf{w}_k - \mathbf{w}_{k+1}\|^2 \right] + \mathbb{E}_{\mathcal{R}} [\|\bar{\mathbf{g}}_k - \mathbf{g}_k\|^2] \right) \\ &= O \left( \frac{\mu(F(\mathbf{w}_1) - F^*)}{T} + \frac{\mu\zeta\|\mathcal{C}\|}{\sqrt{b}} + \frac{\zeta^2}{b} + \mu\beta G\|\mathcal{C}\| \right) = O(\alpha^2), \end{aligned}$$

where in the second to last line we used (12). We now have

$$(1 - \beta) \mathbb{E}_{\mathcal{R}} \left[ \left\| \frac{1}{\eta} (\mathbf{w}_k - \mathbf{w}_k^+) \right\|^2 \right] \leq \mathbb{E}_{\Phi_k, \mathcal{R}} \left[ \left\| \frac{1}{\eta} (\Phi_k \mathbf{w}_k - \Phi_k \mathbf{w}_k^+) \right\|^2 \right] = O(\alpha^2).$$

Rescaling  $\alpha$  proves the claimed bound.  $\square$

## B.2 CONVERGENCE ANALYSIS FOR CONVEX FUNCTIONS

In this section, we analyze the expected error bounds of Algorithm COMPSGD for convex functions. We assume that  $F$  is minimized at some  $\mathbf{w}^* \in \mathcal{C}$ . We consider three cases: a) strongly convex and smooth functions, b) strongly convex functions (without smoothness), and c) convex function (without smoothness). We build on ideas from modern SGD analyses, e.g. [Rakhlin et al., 2011, Shamir and Zhang, 2013].

**Analysis for Strongly Convex and Smooth Functions.** We start with just utilizing the strong convexity assumption. The following lemma provides a bound on the expected error on individual iterate  $\mathbf{w}_t$  for strongly convex functions.

**Lemma B.6.** Let  $F$  be a  $\lambda$ -strongly convex function over a convex set  $\mathcal{C}$ , and  $\mathbb{E}[\|\hat{\mathbf{g}}_t\|^2] \leq G^2$  for  $t \in [T]$  (where  $\vartheta_t = \Phi_t \hat{\mathbf{g}}_t$ ). Then if we pick  $\eta_t = 2/(\lambda t)$  and  $\beta_t = 1/t$ , Algorithm COMPSGD satisfies for every  $t \in [T]$ :

$$\mathbb{E}[\|\mathbf{w}_t - \mathbf{w}^*\|^2] \leq \frac{4G^2}{\lambda^2 t}.$$

*Proof.* We start with Lemma 2.2 with  $\mathbf{w} = \mathbf{w}^*$

$$(1 - \beta_t)\|\mathbf{w}_{t+1} - \mathbf{w}^*\|^2 \leq \|\mathbf{w}_t - \eta_t \hat{\mathbf{g}}_t - \mathbf{w}^*\|^2.$$

Using strong convexity,

$$\mathbb{E}[\langle \hat{\mathbf{g}}_t, \mathbf{w}_t - \mathbf{w}^* \rangle] \geq \mathbb{E}[F(\mathbf{w}_t) - F(\mathbf{w}^*) + (\lambda/2)\|\mathbf{w}_t - \mathbf{w}^*\|^2].$$

The fact that  $\mathbf{w}^*$  is a minimizer of  $F$  in  $\mathcal{C}$  also implies that

$$F(\mathbf{w}_t) - F(\mathbf{w}^*) \geq (\lambda/2)\|\mathbf{w}_t - \mathbf{w}^*\|^2.$$

Using these inequalities,  $\mathbb{E}[\|(\mathbf{w}_t - \eta_t \hat{\mathbf{g}}_t) - \mathbf{w}^*\|^2]$  can be bounded as,

$$\begin{aligned} \mathbb{E}[\|(\mathbf{w}_t - \eta_t \hat{\mathbf{g}}_t) - \mathbf{w}^*\|^2] &= \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}^*\|^2] - 2\eta_t \mathbb{E}[\langle \hat{\mathbf{g}}_t, \mathbf{w}_t - \mathbf{w}^* \rangle] + \eta_t^2 G^2 \\ &\leq \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}^*\|^2] - 2\eta_t \mathbb{E}[F(\mathbf{w}_t) - F(\mathbf{w}^*) + (\lambda/2)\|\mathbf{w}_t - \mathbf{w}^*\|^2] + \eta_t^2 G^2 \\ &\leq \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}^*\|^2] - 2\eta_t \mathbb{E}[(\lambda/2)\|\mathbf{w}_t - \mathbf{w}^*\|^2 + (\lambda/2)\|\mathbf{w}_t - \mathbf{w}^*\|^2] + \eta_t^2 G^2 \\ &= (1 - 2\lambda\eta_t) \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}^*\|^2] + \eta_t^2 G^2. \end{aligned}$$

Plugging this into Lemma 2.2,

$$(1 - \beta_t) \mathbb{E}[\|\mathbf{w}_{t+1} - \mathbf{w}^*\|^2] \leq (1 - 2\eta_t \lambda) \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}^*\|^2] + \eta_t^2 G^2.$$

Setting  $\eta_t = 2/(\lambda t)$  gives,

$$(1 - \beta_t) \mathbb{E}[\|\mathbf{w}_{t+1} - \mathbf{w}^*\|^2] \leq \left(1 - \frac{4}{t}\right) \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}^*\|^2] + \frac{4G^2}{\lambda^2 t^2}.$$

Now,  $1/(1 - \beta_t) \leq (1 + 2\beta_t)$  if  $\beta_t \leq 1/2$ . Therefore,

$$\mathbb{E}[\|\mathbf{w}_{t+1} - \mathbf{w}^*\|^2] \leq (1 + 2\beta_t) \left( \left(1 - \frac{4}{t}\right) \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}^*\|^2] + \frac{4G^2}{\lambda^2 t^2} \right).$$

From Lemma 2 of [Rakhlin et al., 2011], for any  $\mathbf{w}_1 \in \mathcal{C}$ ,  $\mathbb{E}[\|\mathbf{w}_1 - \mathbf{w}^*\|^2] \leq 4G^2/\lambda^2$ . The rest of the proof proceeds by induction using  $t = 2$  as the base case. The base case can be verified by substituting  $t = 2$  in the above equation, and by noting that the left hand side is always non-negative. By inductive hypothesis,  $\mathbb{E}[\|\mathbf{w}_t - \mathbf{w}^*\|^2] \leq 4G^2/(\lambda^2 t)$ , we can complete the proof by noting (after some algebraic manipulation) that,

$$\left(1 + \frac{2}{t}\right) \left( \frac{4G^2}{\lambda^2 t} - \frac{12G^2}{\lambda^2 t^2} \right) \leq \frac{4G^2}{\lambda^2 (t+1)}.$$

□

The following theorem is a simple consequence of the above lemma using the definition of smoothness at the minimizer  $\mathbf{w}^*$ .

**Theorem B.7.** Let  $F$  be a  $\lambda$ -strongly convex and  $\mu$ -smooth function over a convex set  $\mathcal{C}$ , and  $\mathbb{E}[\|\hat{\mathbf{g}}_t\|^2] \leq G^2$  for all  $t \in [T]$  (where  $\vartheta_t = \Phi_t \hat{\mathbf{g}}_t$ ). Then with  $\eta_t = 2/(\lambda t)$  and  $\beta_t = 1/t$ , Algorithm COMPSGD satisfies:

$$\mathbb{E}[F(\mathbf{w}_T) - F(\mathbf{w}^*)] = O\left(\frac{\mu G^2}{\lambda^2 T}\right).$$

**Analysis for Strongly Convex Functions without Smoothness.** We now focus on non-smooth optimization that appear in many machine learning applications (such as in SVMs). For SGD, the  $\log T$  factor appearing in the analysis can be removed through careful stepsize selection [Jain et al., 2019].

**Theorem B.8.** Let  $F$  be a  $\lambda$ -strongly convex function over a convex set  $\mathcal{C}$ , and  $\mathbb{E}[\|\hat{\mathbf{g}}_t\|^2] \leq G^2$  for all  $t \in [T]$  (where  $\vartheta_t = \Phi_t \hat{\mathbf{g}}_t$ ). Then with  $\eta_t = 2/(\lambda t)$  and  $\beta_t = 1/t$ , Algorithm COMPSGD satisfies:

$$\mathbb{E}[F(\mathbf{w}_T) - F(\mathbf{w}^*)] = O\left(\frac{G^2(1 + \log(T))}{\lambda T}\right).$$

*Proof.* Let  $\mathbf{g}_t$  denote a subgradient of  $F$  at  $\mathbf{w}_t$ . Again, we start with Lemma 2.2,

$$(1 - \beta_t)\|\mathbf{w}_{t+1} - \mathbf{w}\|^2 \leq \|\mathbf{w}_t - \eta_t \hat{\mathbf{g}}_t - \mathbf{w}\|^2.$$

Taking expectation on both sides,

$$(1 - \beta_t) \mathbb{E}[\|\mathbf{w}_{t+1} - \mathbf{w}\|^2] \leq \mathbb{E}[\|\mathbf{w}_t - \eta_t \hat{\mathbf{g}}_t - \mathbf{w}\|^2] \leq \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}\|^2] - 2\eta_t \mathbb{E}[\langle \mathbf{g}_t, \mathbf{w}_t - \mathbf{w} \rangle] + \eta_t^2 G^2.$$

Let  $k$  be an arbitrary element in  $\{1, \dots, \lceil T/2 \rceil\}$ . Extracting the inner product, summing over all  $t = T - k, \dots, T$ , and rearranging, we get

$$\begin{aligned} & \sum_{t=T-k}^T \mathbb{E}[\langle \mathbf{g}_t, \mathbf{w}_t - \mathbf{w} \rangle] \\ & \leq \frac{1}{2\eta_{T-k}} \mathbb{E}[\|\mathbf{w}_{T-k} - \mathbf{w}\|^2] + \frac{1}{2} \sum_{t=T-k+1}^T \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}\|^2] \left( \frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} + \frac{\beta_t}{\eta_{t-1}} \right) + \frac{G^2}{2} \sum_{t=T-k}^T \eta_t. \end{aligned}$$

By convexity, we can lower bound  $\mathbb{E}[\langle \mathbf{g}_t, \mathbf{w}_t - \mathbf{w} \rangle] \geq \mathbb{E}[F(\mathbf{w}_t) - F(\mathbf{w})]$ . Using this and substituting for  $\eta_t$  and  $\beta_t$ ,

$$\begin{aligned} & \mathbb{E} \left[ \sum_{t=T-k}^T (F(\mathbf{w}_t) - F(\mathbf{w})) \right] \\ & \leq \frac{\lambda(T-k)}{4} \mathbb{E}[\|\mathbf{w}_{T-k} - \mathbf{w}\|^2] + \frac{\lambda}{4} \sum_{t=T-k+1}^T \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}\|^2] + \frac{\lambda}{4} \sum_{t=T-k+1}^T \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}\|^2] \frac{(t-1)}{2t} + \frac{G^2}{\lambda} \sum_{t=T-k}^T \frac{1}{t} \\ & \leq \frac{\lambda(T-k)}{4} \mathbb{E}[\|\mathbf{w}_{T-k} - \mathbf{w}\|^2] + \frac{3\lambda}{8} \sum_{t=T-k+1}^T \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}\|^2] + \frac{G^2}{\lambda} \sum_{t=T-k}^T \frac{1}{t}. \end{aligned} \quad (13)$$

Picking  $\mathbf{w} = \mathbf{w}_{T-k}$ , and using Lemma B.6,

$$\mathbb{E}[\|\mathbf{w}_t - \mathbf{w}_{T-k}\|^2] \leq 2 \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}^*\|^2 + \|\mathbf{w}_{T-k} - \mathbf{w}^*\|^2] \leq \frac{8G^2}{\lambda^2} \left( \frac{1}{t} + \frac{1}{T-k} \right) \leq \frac{32G^2}{\lambda^2 T},$$

where we have used that  $t \geq T - k$ . Plugging this into (13)

$$\mathbb{E} \left[ \sum_{t=T-k}^T (F(\mathbf{w}_t) - F(\mathbf{w})) \right] \leq \frac{12G^2 k}{\lambda T} + \frac{G^2}{\lambda} \sum_{t=T-k}^T \frac{1}{t}.$$

The remainder of the proof follows as in Theorem 1 of [Shamir and Zhang, 2013], where starting from the previous equation it can be shown that,

$$\mathbb{E}[F(\mathbf{w}_T) - F(\mathbf{w}^*)] = O\left(\frac{G^2(1 + \log(T))}{\lambda T}\right).$$

□

**Analysis for Convex Functions.** The analysis from Theorem B.8 can also be extended to a case of a general convex function without any assumption on strong convexity or smoothness.

**Theorem B.9.** Let  $F$  be convex function over a convex set  $\mathcal{C}$ , and  $\mathbb{E}[\|\hat{\mathbf{g}}_t\|^2] \leq G^2$  for all  $t \in [T]$  (where  $\vartheta_t = \Phi_t \hat{\mathbf{g}}_t$ ). Then with  $\eta_t = \varphi/\sqrt{t}$  (for some constant  $\varphi$ ) and  $\beta_t = 1/t$ , Algorithm COMPSGD satisfies:

$$\mathbb{E}[F(\mathbf{w}_T) - F(\mathbf{w}^*)] = O\left(\left(\frac{\|\mathcal{C}\|^2}{\varphi} + \varphi G^2\right) \frac{\log T}{\sqrt{T}}\right).$$

In particular with  $\varphi = \|\mathcal{C}\|/G$ ,  $\mathbb{E}[F(\mathbf{w}_T) - F(\mathbf{w}^*)] = O\left(\|\mathcal{C}\|G \log T/\sqrt{T}\right)$ .

*Proof.* The proof begins the same as in Theorem B.8, with  $k$  now being from the set  $[T - 1]$ .

$$\begin{aligned} & \sum_{t=T-k}^T \mathbb{E}[\langle \mathbf{g}_t, \mathbf{w}_t - \mathbf{w} \rangle] \\ & \leq \frac{1}{2\eta_{T-k}} \mathbb{E}[\|\mathbf{w}_{T-k} - \mathbf{w}\|^2] + \frac{1}{2} \sum_{t=T-k+1}^T \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}\|^2] \left( \frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} + \frac{\beta_t}{\eta_{t-1}} \right) + \frac{G^2}{2} \sum_{t=T-k}^T \eta_t. \end{aligned}$$

Setting  $\mathbf{w} = \mathbf{w}_{T-k}$ ,  $\eta_t = \varphi/\sqrt{t}$ ,  $\mathbb{E}[\|\mathbf{w}_t - \mathbf{w}\|^2] \leq 4\|\mathcal{C}\|^2$  (by assumption on the diameter of  $\mathcal{C}$ ).

$$\begin{aligned} \sum_{t=T-k}^T \mathbb{E}[\langle \mathbf{g}_t, \mathbf{w}_t - \mathbf{w} \rangle] & \leq \frac{2\|\mathcal{C}\|^2}{\varphi} (\sqrt{T} - \sqrt{T-k}) + \frac{2\|\mathcal{C}\|^2}{\varphi} \sum_{t=T-k+1}^T \frac{\sqrt{t-1}}{t} + \frac{G^2\varphi}{2} \sum_{t=T-k}^T \frac{1}{\sqrt{t}} \\ & \leq \frac{2\|\mathcal{C}\|^2}{\varphi} (\sqrt{T} - \sqrt{T-k}) + \frac{2\|\mathcal{C}\|^2}{\varphi} \sum_{t=T-k+1}^T \frac{1}{\sqrt{t}} + \frac{G^2\varphi}{2} \sum_{t=T-k}^T \frac{1}{\sqrt{t}} \\ & \leq \frac{2\|\mathcal{C}\|^2}{\varphi} (\sqrt{T} - \sqrt{T-k}) + \frac{2\|\mathcal{C}\|^2}{\varphi} 2(\sqrt{T} - \sqrt{T-k}) + \frac{G^2\varphi}{2} 2(\sqrt{T} - \sqrt{T-k-1}), \end{aligned}$$

where we have used the fact that  $\sum_{t=T-k}^T 1/\sqrt{t} \leq 2(\sqrt{T} - \sqrt{T-k-1})$  and  $\sum_{t=T-k+1}^T 1/\sqrt{t} \leq 2(\sqrt{T} - \sqrt{T-k})$ . Simplifying the above bound (using that  $(\sqrt{T} - \sqrt{T-k-1}) \geq (\sqrt{T} - \sqrt{T-k})$ ), we get

$$\sum_{t=T-k}^T \mathbb{E}[\langle \mathbf{g}_t, \mathbf{w}_t - \mathbf{w} \rangle] \leq \left( \frac{6\|\mathcal{C}\|^2}{\varphi} + G^2\varphi \right) (\sqrt{T} - \sqrt{T-k-1}).$$

Using the convexity bound  $\mathbb{E}[\langle \mathbf{g}_t, \mathbf{w}_t - \mathbf{w}_{T-k} \rangle] \geq \mathbb{E}[F(\mathbf{w}_t) - F(\mathbf{w}_{T-k})]$ ,

$$\mathbb{E} \left[ \sum_{t=T-k}^T F(\mathbf{w}_t) - F(\mathbf{w}_{T-k}) \right] \leq \left( \frac{6\|\mathcal{C}\|^2}{\varphi} + G^2\varphi \right) (\sqrt{T} - \sqrt{T-k-1}).$$

The remainder of the proof follows as in Theorem 2 of [Shamir and Zhang, 2013], whose proof can be used to show, starting from the previous equation,

$$\mathbb{E}[(F(\mathbf{w}_T) - F(\mathbf{w}^*))] = O \left( \left( \frac{\|\mathcal{C}\|^2}{\varphi} + \varphi G^2 \right) \frac{\log T}{\sqrt{T}} \right).$$

□

### B.3 MISSING DETAILS FROM SECTION 2.2

**Dataset Description.** Here we evaluate on three common datasets: a) MNIST (60,000  $28 \times 28$  grayscale images of the 10 digits, along with a test set of 10,000 images) [LeCun et al., 1998], b) Reuters newswire topics (11,228 newswires from Reuters split as 8982 train and 2246 test, labeled over 46 topics) [Keras, 2017], and c) IMDB movie reviews sentiment (50,000 movies reviews from IMDB split as 25000 train and 25000 test, labeled by sentiment positive/negative) [Maas et al., 2011]. For the Reuters and IMDB datasets, we use the top 10000 words from each dataset respectively for input representation, i.e.,  $d = 10000$ .

**Different Constraint Sets.** In Figure 5, we present results on some other popular constraint sets. In Figure 5a, we plot the results (objective value) for a least-squares problem:  $\min_{\mathbf{w}} \|\mathbf{y} - A\mathbf{w}\|^2$  subject to a constraint that  $\mathbf{w}$  lies in a 10-dimensional subspace  $U \subset \mathbb{R}^d$ . Here, matrix  $A \in \mathbb{R}^{n \times d}$  is random standard Gaussian matrix and  $\mathbf{y} \in \mathbb{R}^n$  where  $\mathbf{y} = A\mathbf{w}^*$  with a random  $\mathbf{w}^*$  picked from  $U$ , and  $n = 1000$ ,  $d = 10000$ . The Gaussian width of  $U$  is  $\sqrt{10}$ , and the dimensionality of utilized gradients by Algorithm COMPSGD in various epochs is presented in Figure 5b.

In Figure 5c, we use the same setup, except now the constraint on  $\mathbf{w}$  is that it lies in a positive simplex  $S$ , and  $\mathbf{w}^*$  is randomly picked from this positive simplex  $S$ . Again, and the dimensionality of utilized gradients by Algorithm COMPSGD



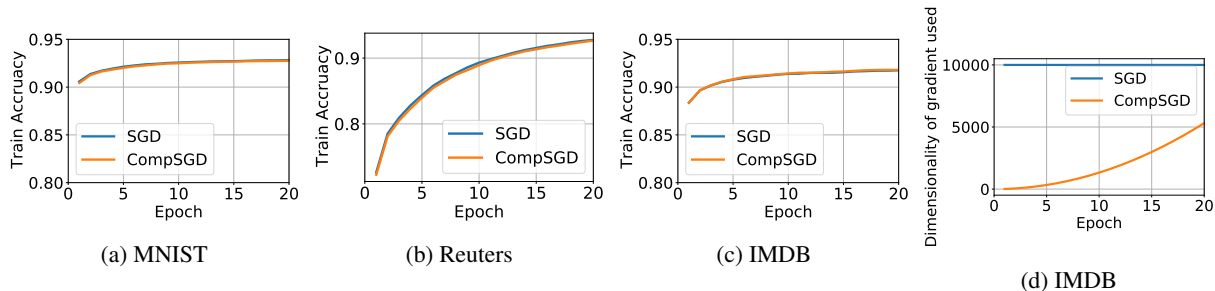


Figure 3: Training accuracy for: (a) logistic regression with  $\ell_1$ -constraint on MNIST dataset, (b) logistic regression with  $\ell_1$ -constraint on Reuters dataset, and (c) logistic regression with  $\ell_1$ -constraint on IMDB dataset. Note that the performance of SGD and COMPSGD are almost identical. (d) Comparison of dimensionality of the utilized gradients on IMDB dataset.

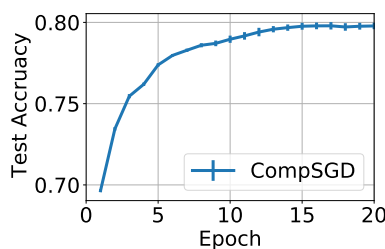


Figure 4: Performance of Algorithm COMPSGD for solving logistic regression with  $\ell_1$ -constraint on Reuters dataset. The error bars show one standard deviation computed across 10 runs.

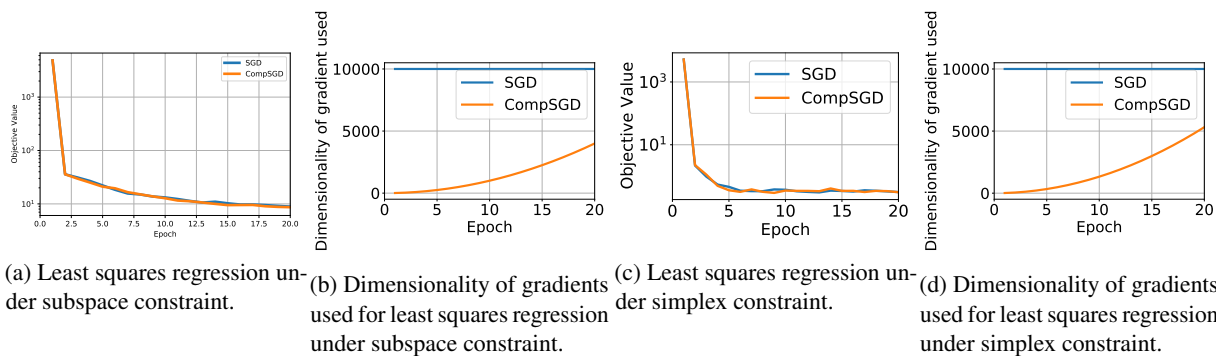


Figure 5: Least-squares objective with different constraint sets on a synthetic dataset.

in various epochs is presented in Figure 5d. The key takeaway is that the performance of COMPSGD is again quite identical to regular SGD while utilizing much lower-dimensional gradients.

**Experimental Results on Nonconvex Functions.** We optimize the network under an  $\ell_1$ -constraint on weights in each layer. Since computing exactly the  $\beta_t$  and batchsize provided in the theoretical bounds is hard in practice, we use a fixed batchsize of 32 and set  $m_t = d_i/10$  for each layer  $i$ , where  $d_i$  is the original number of parameters in layer  $i$ . The results show that one can match SGD (or even marginally improve) performance with lower-dimensional gradients using Algorithm COMPSGD. The training accuracy plots presented in Figure 6 also exhibit a similar behavior.

In terms of the running time, for the MNIST, Reuters, and IMDB datasets, SGD took an average of 24.77, 30.23, 72.48 seconds per epoch respectively, whereas Algorithm COMPSGD took an average of 28.83, 34.01, 78.51 seconds per epoch respectively. So, for example on the IMDB dataset, CompSGD is only about 8.3% slower than SGD, however the dimensionality of utilized gradients is reduced by a factor of 10.

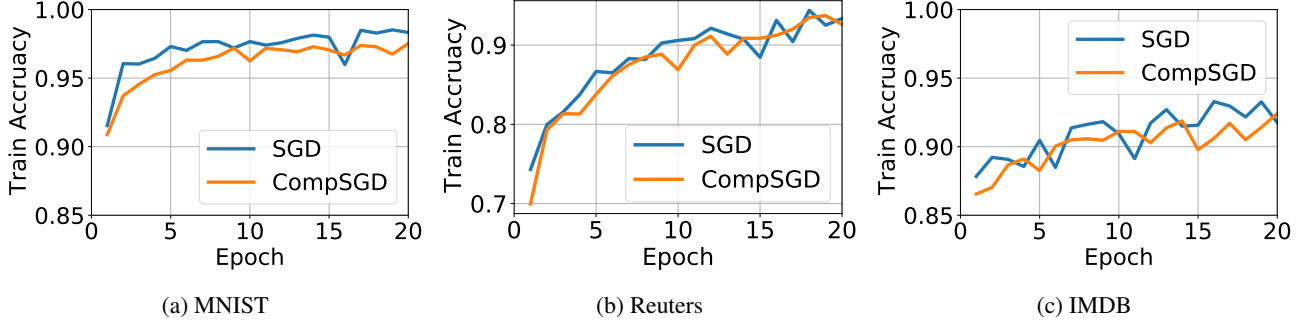


Figure 6: Training accuracy of a MLP network with SGD and COMPSGD on different datasets. Note that the performance of SGD and COMPSGD are almost identical.

## C DIFFERENTIALLY PRIVATE ERM WITH NONCONVEX FUNCTIONS

We consider the standard Empirical Risk Minimization (ERM) framework. Given a dataset  $D = \{\mathbf{z}_1, \dots, \mathbf{z}_n\}$  from a data universe  $\mathcal{D}^n$ , the goal in ERM is to

$$\min_{\mathbf{w} \in \mathcal{C}} F(\mathbf{w}; D) := \frac{1}{n} \sum_{i=1}^n f(\mathbf{w}; \mathbf{z}_i),$$

where  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is the loss of model  $\mathbf{w} \in \mathbb{R}^d$  for data record  $\mathbf{z}_i$ . This formulation also captures regularized ERM, in which an additional function  $r(\mathbf{w})$  is added to the loss function to penalize certain types of solutions. One can fold the regularizer  $r(\cdot)$  into the data-dependent functions by replacing  $f(\mathbf{w}; \mathbf{z}_i)$  by  $f(\mathbf{w}; \mathbf{z}_i) + r(\mathbf{w})/n$ .

**Our Approach.** We start with the simple differentially private gradient descent (DP-SGD) algorithm [Song et al., 2013, Bassily et al., 2014].<sup>3</sup> The basic idea is to perturb each intermediate update by adding noise to the gradients.

$$\text{DP-SGD iteration: } \theta_{t+1} \leftarrow \theta_t - \eta_t (\nabla F(\mathbf{w}; D) + \mathbf{e}), \quad (14)$$

where  $\eta_t$  is the learning rate and  $\mathbf{e} \sim \mathcal{N}(0, \sigma^2 \mathbb{I}_d)$  is the calibrated noise to achieve  $(\epsilon, \delta)$ -differential privacy. The fact that  $\nabla F(\mathbf{w}; D) \in \mathbb{R}^d$ , means that the popular *Gaussian mechanism* (Theorem A.1) idea in DP has  $\mathbb{E}[\|\mathbf{e}\|^2] = \sigma^2 d$  that brings about the dependence on the dimension  $d$  into the utility analysis.

### Algorithm PRIVSGD: Differentially Private SGD with Dimensionality-reduced Gradients

**Input:** Dataset  $D = (\mathbf{z}_1, \dots, \mathbf{z}_n)$ , privacy parameters  $(\epsilon, \delta)$ , learning rate parameters  $\{\eta_t\}$ , projection dimension parameters  $\{\beta_t\}$ .

1. Pick  $\mathbf{w}_1$  as any point in  $\mathcal{C}$
2. for  $t = 1$  to  $T$  do
  - a. Set  $m_t \leftarrow \min\{d, c \cdot \omega(\mathcal{C})^2 / \beta_t^2\}$  (for some constant  $c$ )
  - b. Let  $\Phi_t \in \mathbb{R}^{m_t \times d} \sim_{i.i.d} \mathcal{N}(0, 1/m_t)$
  - c. Set  $\sigma^2 \leftarrow \frac{32L^2 T \log(1/\delta)}{n^2 \epsilon^2}$
  - d. Let  $s_t \leftarrow \frac{\|\nabla F(\mathbf{w}_t; D)\|}{\|\Phi_t \nabla F(\mathbf{w}_t; D)\|}$
  - e. Let  $\theta_t \leftarrow \Pi_{\Phi_t \mathcal{C}}(\Phi_t \mathbf{w}_t - \eta_t (s_t \Phi_t \nabla F(\mathbf{w}_t; D) + \mathbf{e}))$  where  $\mathbf{e} \sim \mathcal{N}(0, \sigma^2 \mathbb{I}_m)$  (fresh noise)
  - f. Let  $\mathbf{w}_{t+1} \leftarrow$  pick any element from the set  $\mathcal{S}_t = \{\mathbf{w} \in \mathcal{C} : \Phi_t \mathbf{w} = \theta_t\}$  (e.g., by solving (2))

An approach to reduce the dependence on  $d$  arising from noise addition in (14) is to reduce the dimensionality of the gradient vector. For example, if we take  $\Phi \in \mathbb{R}^{m \times d}$  with entries drawn i.i.d from  $\mathcal{N}(0, 1/m)$ , then  $\Phi \nabla F(\mathbf{w}; D) \in \mathbb{R}^m$ . Therefore, by using the Gaussian mechanism now, one could add noise as:  $\Phi \nabla F(\mathbf{w}; D) + \mathbf{e}$  where  $\mathbf{e} \sim \mathcal{N}(0, \sigma^2 \mathbb{I}_m)$ . This roughly

<sup>3</sup>A background in differential privacy is provided in Appendix A.

changes the dependence on  $d$  to  $m$  in the convergence analysis. This idea is formalized in Algorithm PRIVSGD. We assume that  $\|\nabla f(\mathbf{w}; \cdot)\|$  is uniformly bounded by  $L$  for all  $\mathbf{w} \in \mathcal{C}$ , and also  $f$  is  $\mu$ -smooth and continuously differentiable in the first parameter. We use (full) gradient descent here for simplicity.

The proof of  $(\epsilon, \delta)$ -DP of Algorithm PRIVSGD is fairly straightforward based on the privacy guarantees of Gaussian mechanism (see Theorem A.1), and the strong composition (see Theorem A.2) and follows as in [Bassily et al., 2014, Theorem 2.1]. The normalizing factor  $s_t$  ensures that  $\|s_t \Phi_t \nabla F(\mathbf{w}_t; D)\| = \|\nabla F(\mathbf{w}_t; D)\| \leq L$  (by our assumption), thus bounding the  $\ell_2$ -global sensitivity (see Appendix A) defined as

$$\sup_{D, D' \text{ neighbors}} \|s_t \Phi_t \nabla F(\mathbf{w}_t; D) - s_t \Phi_t \nabla F(\mathbf{w}_t; D')\| \leq 2L.$$

Moreover, Algorithm COMPSGD can release all the intermediate  $\mathbf{w}_t$ 's without damaging the privacy guarantee as differential privacy is immune to post processing [Dwork et al., 2006b].

**Claim C.1.** *Algorithm PRIVSGD is  $(\epsilon, \delta)$ -DP.*

We focus on the utility guarantees of Algorithm PRIVSGD using the previously established convergence bounds on Algorithm COMPSGD. Table 4 summarizes our results. These are the first results in differentially private nonconvex optimization which provides meaningful guarantees when  $n \gg \omega(\mathcal{C})$  rather than requiring  $n \gg \sqrt{d}$ . Note that there will be an increase in the sample size with the number of rounds for any differentially private SGD algorithm due to the fact it will have to use some composition result (such as Theorem A.2) to account for accumulation of privacy parameters over multiple rounds.

The following corollaries follow from Theorems 2.4 and B.5, where the role of  $\zeta/\sqrt{b}$  is replaced by  $\sigma$  (the variance in the coordinates of  $\mathbf{e}$ ).

**Corollary C.2** (From Theorem 2.4). *Let  $f$  be  $\mu$ -smooth and continuously differentiable function over a convex set  $\mathcal{C}$ . Let  $\|\nabla f(\mathbf{w}; \cdot)\| \leq L$  for all  $\mathbf{w} \in \mathcal{C}$ . Let  $\rho \in (0, 1)$  and  $\alpha > 0$ . In Algorithm PRIVSGD, let  $T = \Theta((F(\mathbf{w}_1) - F^*)(L + \mu\|\mathcal{C}\|)^2 / (\mu\alpha^2\rho))$ ,  $\eta_t = \eta = 1/\mu$ , and  $\beta_t = \beta = \min\{1/4, (\mu\alpha^2)/(64L\|\mathcal{C}\|(L + \mu\|\mathcal{C}\|)^2)\}$ . Output the first  $\mathbf{w}_\tau$  (if it exists) in  $(\mathbf{w}_1, \dots, \mathbf{w}_T)$  such that  $\|\mathbf{w}_\tau - \mathbf{w}_{\tau+1}\| \leq \alpha/(L + \mu\|\mathcal{C}\|)$  and  $\perp$  (fail) otherwise. Then, with probability at least  $(1 - \rho)^2$ , this procedure outputs a  $\mathbf{w}_\tau$  that is an  $\alpha$ -FOSP for  $F$ , if the sample size  $n \geq 1$  satisfies:*

$$n = \Omega \left( \min \left\{ \frac{\omega(\mathcal{C})}{\beta}, \sqrt{d} \right\} \cdot \frac{\|\mathcal{C}\|(L + \mu\|\mathcal{C}\|)^2 L \sqrt{T \log(1/\delta)}}{\sqrt{\rho} \alpha^2 \mu \epsilon} \right).$$

*Proof Sketch.* The proof is identical to Theorem 2.4, so we just give the high-level idea. Again, we establish the proof over two parts.

1. If  $\mathbf{w}_\tau$  exists, then with probability at least  $1 - \rho$ ,  $\mathbf{w}_\tau$  is an  $\alpha$ -FOSP.
2. If  $T = \Omega \left( \frac{(F(\mathbf{w}_1) - F^*)(L + \mu\|\mathcal{C}\|)^2}{\mu\alpha^2\rho} \right)$ , then with probability at least  $1 - \rho$ ,  $\mathbf{w}_\tau$  exists in  $(\mathbf{w}_1, \dots, \mathbf{w}_T)$ .

Let  $\mathbf{g}_t = \nabla F(\mathbf{w}_t; D)$ . For the first part assume  $\mathbf{w}_\tau$  exists. We replace the role of  $\Phi_t \bar{\mathbf{g}}_t$  in Theorem 2.4 by  $\Phi_t \mathbf{g}_t + \mathbf{e}_t$ . By noting that  $\mathbb{E}_{\Phi_t}[\mathbf{s}_t] = 1$ , it is easy to derive an expression similar to (8),

$$\min_{\mathbf{w} \in \mathcal{C}} \langle \mathbf{g}_\tau, \mathbf{w} - \mathbf{w}_\tau \rangle \geq -(1 + \beta)(\mu\|\mathcal{C}\| \|\mathbf{w}_\tau - \mathbf{w}_{\tau+1}\| + L\|\mathbf{w}_\tau - \mathbf{w}_{\tau+1}\| + \|\mathcal{C}\| \|\mathbf{e}_\tau\|),$$

which suffices to establish the first part. For the second part, similar to (9), we get,

$$\langle \mathbf{g}_t, \mathbf{w}_{t+1} - \mathbf{w}_t \rangle \leq \mathbb{E}_{\Phi_t}[\langle s_t \Phi_t \mathbf{g}_t + \mathbf{e}_t, \Phi_t(\mathbf{w}_{t+1} - \mathbf{w}_t) \rangle] + \|\mathbf{e}_t\| \|\mathbf{w}_{t+1} - \mathbf{w}_t\| + \beta \|\mathbf{g}_t\| \|\mathbf{w}_{t+1} - \mathbf{w}_t\|.$$

The remainder of the proof proceeds identical to Theorem 2.4. □

Using Theorem B.5 one could also get a bound on Algorithm PRIVSGD for achieving stationarity based on compressed gradient mapping.

**Corollary C.3** (From Theorem B.5). *Let  $f$  be  $\mu$ -smooth and continuously differentiable function over a convex set  $\mathcal{C}$ . Let  $\|\nabla f(\mathbf{w}; \cdot)\| \leq L$  for all  $\mathbf{w} \in \mathcal{C}$ . Let  $\alpha > 0$ . In Algorithm PRIVSGD, let  $T = \Theta(\mu(F(\mathbf{w}_1) - F^*)/\alpha^2)$ ,  $\eta_t = \eta = 1/(2\mu)$ , and  $\beta_t = \beta = \min\{1/4, \alpha^2/(\mu L\|\mathcal{C}\|)\}$ . Output a uniformly at random iterate  $\mathbf{w}_k$  from  $\{\mathbf{w}_1, \dots, \mathbf{w}_T\}$ . Then,  $\mathbf{w}_k$  satisfies the  $\alpha$ -compressed gradient mapping condition (Definition 8) in expectation, if the sample size  $n \geq 1$  satisfies:*

$$n = \Omega \left( \max \left\{ \min \left\{ \frac{\omega(\mathcal{C})}{\beta}, \sqrt{d} \right\} \cdot \frac{\mu\|\mathcal{C}\| L \sqrt{T \log(1/\delta)}}{\epsilon \alpha^2}, \min \left\{ \frac{\omega(\mathcal{C})}{\beta}, \sqrt{d} \right\} \cdot \frac{L \sqrt{T \log(1/\delta)}}{\epsilon \alpha} \right\} \right).$$

Guarantee	DP-SGD Sample Size Req.	Algorithm PRIVSGD Sample Size Req.
$\alpha$ -FOSP	$\Omega \left( \sqrt{d} \cdot \frac{\ \mathcal{C}\ (L+\mu\ \mathcal{C}\ )^2 L \sqrt{T \log(1/\delta)}}{\sqrt{\rho} \alpha^2 \mu \epsilon} \right)$	$\Omega \left( \min \left\{ \frac{\omega(\mathcal{C})}{\beta}, \sqrt{d} \right\} \cdot \frac{\ \mathcal{C}\ (L+\mu\ \mathcal{C}\ )^2 L \sqrt{T \log(1/\delta)}}{\sqrt{\rho} \alpha^2 \mu \epsilon} \right)$
$\alpha$ -comp. gradient mapping	$\approx \Omega \left( \sqrt{d} \cdot \frac{\mu \ \mathcal{C}\  L \sqrt{T \log(1/\delta)}}{\epsilon \alpha^2} \right)$	$\approx \Omega \left( \min \left\{ \frac{\omega(\mathcal{C})}{\beta}, \sqrt{d} \right\} \cdot \frac{\mu \ \mathcal{C}\  L \sqrt{T \log(1/\delta)}}{\epsilon \alpha^2} \right)$

Table 4: Comparison of sample size ( $n$ ) needed for achieving a first-order stationarity condition for a nonconvex  $\mu$ -smooth function under  $(\epsilon, \delta)$ -DP. Assume  $\|\nabla f(\mathbf{w}; \cdot)\|$  is uniformly bounded by  $L$  for all  $\mathbf{w} \in \mathcal{C}$ . The settings of  $\beta$  and the upper bound on  $T$  are from Corollaries C.2 and C.3.

## D REDUCING COMMUNICATION IN DISTRIBUTED SYNCHRONOUS SGD

Training very large machine learning models requires a distributed computing approach, with communication of the model updates often being the bottleneck. This is especially true in the federated learning setting where clients are mobile devices with expensive up-link communication cost [McMahan et al., 2017].

Here we consider a distributed synchronous model. Let us assume that there are  $M$  clients, numbered  $1, \dots, M$ . Let  $F(\mathbf{w}) : \mathbb{R}^d \rightarrow \mathbb{R}$ . We investigate the following optimization problem

$$\min_{\mathbf{w} \in \mathcal{C}} F(\mathbf{w}) := \frac{1}{M} \sum_{i=1}^M f_i(\mathbf{w}), \quad (15)$$

where each  $f_i$  resides at the  $i$ th client. As an illustration of the above setup, consider a machine learning problem, with data  $\{\mathbf{z}_1, \dots, \mathbf{z}_n\}$  partitioned on the clients, with  $P_i$  the set of indexes of datapoints on client  $i$ , then with  $f_i(\mathbf{w}) = \frac{M}{n} \sum_{j \in P_i} f(\mathbf{w}; \mathbf{z}_j)$ , we get  $F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n f(\mathbf{w}; \mathbf{z}_i)$ .

In each iteration  $t$  of synchronous SGD, the server randomly picks a set  $R_t$  of  $\kappa \geq 1$  clients and sends them the current model parameter  $\mathbf{w}_t$ . Each of these selected client  $i$  computes  $\hat{\mathbf{g}}_t^{(i)}$  an independent stochastic (sub)gradient of  $f_i$  at  $\mathbf{w}_t$ , and communicates  $\hat{\mathbf{g}}_t^{(i)}$  back to the server. The central server then aggregates these gradients and applies the update

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \frac{\eta t}{\kappa} \sum_{i \in R_t} \hat{\mathbf{g}}_t^{(i)}.$$

Naively for the above protocol, the resulting per round communication is roughly  $\kappa \cdot (32d)$  bits (assuming 32-bit floating point numbers).

Prior gradient encoding works focus on reducing the communication cost of  $\kappa \cdot (32d)$  to  $\kappa \cdot (\gamma d)$  where  $\gamma < 32$ . However, an important point is that, in general, none of these prior encoding techniques inherently change the dimensionality of the gradient and also suffer from inflated variance when  $\gamma$  is really small. With only constant variance blowup,  $\Omega(d)$  bits per iteration is also necessary if  $\mathcal{C} = \mathbb{R}^d$  (unconstrained setting), see, e.g., communication complexity lower bound of distributed mean estimation [Zhang et al., 2013]. We show that in-fact we can go beyond these worst-case bounds by exploiting properties of the constraint set  $\mathcal{C}$ . Our results (Theorems B.7, B.8, B.9, 2.4, and B.5) already show that, by transmitting lower-dimensional gradients, we can reduce communication costs, while preserving the convergence guarantees to within constant factor of the regular distributed SGD. The reduction of  $d$  is especially important as it is common to have  $d > 10^8$  in deep networks. In this section, we build on these results and show that in fact one can use any unbiased gradient compression scheme on top of our idea of utilizing lower-dimensional gradients.

**Our Approach.** For a formal analysis, we define a *contraction operator* to capture a general class of existing unbiased gradient encoding schemes. We define a contraction operator as a (possibly randomized) function  $C_a$  defined as mapping from  $\mathbb{R}^a \rightarrow \mathbb{R}^a$  for  $a \in \mathbb{N}$  that satisfies these following common assumptions:

$$\forall \mathbf{v} \in \mathbb{R}^a, \mathbb{E}[C_a(\mathbf{v})] = \mathbf{v} \quad (\text{unbiasedness}) \text{ and } \mathbb{E}[\|C_a(\mathbf{v}) - \mathbf{v}\|^2] \leq \chi_{C_a} \|\mathbf{v}\|^2 \quad (\text{variance bound}). \quad (16)$$

Let  $\gamma_{C_a}$  denote the bits needed to communicate  $C_a(\mathbf{v})$  from a client to server. This general notion captures many common unbiased quantization techniques such as *stochastic rounding* [Alistarh et al., 2017, Wen et al., 2017], vector quantization techniques such as *vqSGD* [Gandikota et al., 2019], and sparsification techniques such as *random sparsification* [Stich et al., 2018]. As an example, the quantization technique of [Alistarh et al., 2017] satisfies  $\chi_{C_a} = 1$  and  $\gamma_{C_a} \approx 2.8a + 32$  for all  $a \in \mathbb{N}$ .

Our procedure is presented in Algorithm COMPDISTSGD. The overall idea is simple, take a random projection of a gradient and then apply the contraction operator  $C_{m_t}$ . More formally, client  $i$  at iteration  $t$  will transmit  $C_{m_t}(\Phi_t \hat{\mathbf{g}}_t^{(i)})$ . The total communication cost in iteration  $t$  equals  $\kappa \cdot \gamma_{C_{m_t}}$  bits. This saves at least a factor  $\gamma_{C_d}/\gamma_{C_{m_t}}$  over just applying the contraction operator over the true gradients in each iteration  $t$ , which could be significant when  $m_t \ll d$ .<sup>4</sup> For example, if  $C$  is the quantization technique of [Alistarh et al., 2017], then  $\gamma_{C_d}/\gamma_{C_{m_t}} \approx d/m_t$  where  $m_t \approx \omega(\mathcal{C})^2/\beta_t^2$ , and for sets  $\mathcal{C}$  such as the unit  $l_1$ -ball,  $\omega(\mathcal{C})^2 = O(\log d)$ .

#### Algorithm COMPDISTSGD

**Setup:**  $M$  clients numbered  $1, \dots, M$ , each with function  $f_i$ . Contraction operator  $C$ . Batchsize  $b$ . Learning rate parameters are  $\{\eta_t\}$  and projection dimension parameters are  $\{\beta_t\}$ .

ServerUpdate: //Run on the server

1. Pick  $\mathbf{w}_1$  as any point in  $\mathcal{C}$
2. Set  $m_t \leftarrow \min\{d, c \cdot \omega(\mathcal{C})^2/\beta_t^2\}$  (for some constant  $c$ )
3. For  $t = 1$  to  $T$  do
  - a. Let  $\Phi_t \in \mathbb{R}^{m_t \times d}$  be an i.i.d. random projection matrix
  - b. Let  $R_t \leftarrow$  random set of  $\kappa$  clients chosen from  $[M]$  (with replacement)
  - c. For each client  $i \in R_t$  (in parallel do)
    - $\vartheta_t^{(i)} \leftarrow$  ClientUpdate( $i, \mathbf{w}_t, \Phi_t$ )
  - c. Let  $\theta_t \leftarrow \Pi_{\Phi_t \mathcal{C}} \left( \Phi_t \mathbf{w}_t - \frac{\eta_t}{\kappa} \sum_{i \in R_t} \vartheta_t^{(i)} \right)$
  - d. Let  $\mathbf{w}_{t+1} \leftarrow$  pick any element from the set  $\mathcal{S}_t = \{\mathbf{w} \in \mathcal{C} : \Phi_t \mathbf{w} = \theta_t\}$  (e.g., by solving (2))

ClientUpdate( $i, \mathbf{w}_t, \Phi_t$ ): //Run on Client  $i$

1. Compute  $\hat{\mathbf{g}}_t^{(i)}$  as independent stochastic (sub)gradient of  $f_i$  at  $\mathbf{w}_t$
2. Return  $C_{m_t}(\Phi_t \hat{\mathbf{g}}_t^{(i)})$  back to the server

The cost of communicating  $\Phi_t$  is not significant as it can be achieved using various techniques such as one-to-all broadcasting. In practice,  $\Phi_t$  will be generated by a pseudorandom generator initialized by some seed, so by just communicating the seed we can regenerate  $\Phi_t$  at each device.

For our analysis, we make the following assumptions. Let  $\nabla f_i(\mathbf{w}_t)$  denote the (sub)gradient of  $f_i$  at  $\mathbf{w}_t$ .

$$\begin{aligned} \textbf{Assumption A: } \mathbb{E}[\|\hat{\mathbf{g}}_t^{(i)}\|^2] &\leq G^2, \quad \textbf{Assumption B: } \mathbb{E}[\|\hat{\mathbf{g}}_t^{(i)} - \nabla f_i(\mathbf{w}_t)\|^2] \leq \zeta^2, \quad \forall i \in [M], t \in [T], \\ \textbf{Assumption C: } \mathbb{E}_i[\|\nabla f_i(\mathbf{w}) - \nabla F(\mathbf{w})\|^2] &\leq \vartheta^2 \quad \forall \mathbf{w} \in \mathcal{C}. \end{aligned} \quad (17)$$

Assumptions A and B are similar to (4) on the subgradients. Assumption C is also common in the distributed SGD literature as a measure of non-iidness among clients, and is referred to as *bounded inter-client gradient variance* [Kairouz et al., 2019, Lian et al., 2018, Jiang and Agrawal, 2018].

**Lemma D.1.** Let  $\bar{\mathbf{g}}_t = \frac{1}{\kappa} \sum_{i \in R_t} C_{m_t}(\Phi_t \hat{\mathbf{g}}_t^{(i)})$ . Then (over all randomness in Algorithm COMPDISTSGD) (a)  $\mathbb{E}[\bar{\mathbf{g}}_t] = \frac{1}{\kappa} \sum_{i \in R_t} \Phi_t \nabla f_i(\mathbf{w}_t)$  (unbiasedness), (b)  $\mathbb{E}[\|\bar{\mathbf{g}}_t\|^2] \leq 2(\frac{\chi_{C_{m_t}} G^2}{\kappa} + G^2)$  (second moment), and (c)  $\mathbb{E}[\|\bar{\mathbf{g}}_t - \Phi_t \nabla F(\mathbf{w}_t)\|^2] \leq 2(\chi_{C_{m_t}} G^2 + \zeta^2 + \vartheta^2)$  (variance).

*Proof.* From unbiasedness of contraction operator and  $\hat{\mathbf{g}}_t^{(i)}$ , the unbiasedness follows. For Part (b) (expectation over all

<sup>4</sup>As  $\gamma_{C_a}$  would be a non-decreasing function in  $a$ , i.e., cost of communicating a longer vector can't be less than one for a shorter vector.

randomness in Algorithm COMPDISTSGD),

$$\begin{aligned}
\mathbb{E}[\|\frac{1}{\kappa} \sum_{i \in R_t} C_{m_t}(\Phi_t \hat{\mathbf{g}}_t^{(i)})\|^2] &= \mathbb{E}[\|\frac{1}{\kappa} \sum_{i \in R_t} C_{m_t}(\Phi_t \hat{\mathbf{g}}_t^{(i)}) - \frac{1}{\kappa} \sum_{i \in R_t} \Phi_t \hat{\mathbf{g}}_t^{(i)} + \frac{1}{\kappa} \sum_{i \in R_t} \Phi_t \hat{\mathbf{g}}_t^{(i)}\|^2] \\
&= \mathbb{E}[\|\frac{1}{\kappa} \sum_{i \in R_t} C_{m_t}(\Phi_t \hat{\mathbf{g}}_t^{(i)}) - \frac{1}{\kappa} \sum_{i \in R_t} \Phi_t \hat{\mathbf{g}}_t^{(i)}\|^2] + \mathbb{E}[\|\frac{1}{\kappa} \sum_{i \in R_t} \Phi_t \hat{\mathbf{g}}_t^{(i)}\|^2] \\
&= \frac{1}{\kappa^2} \sum_{i \in R_t} \mathbb{E}[\|C_{m_t}(\Phi_t \hat{\mathbf{g}}_t^{(i)}) - \Phi_t \hat{\mathbf{g}}_t^{(i)}\|^2] + \mathbb{E}[\|\frac{1}{\kappa} \sum_{i \in R_t} \Phi_t \hat{\mathbf{g}}_t^{(i)}\|^2] \\
&\leq (1 + \beta_t) \left( \frac{\chi_{C_{m_t}} G^2}{\kappa} + G^2 \right).
\end{aligned}$$

Here, we used  $C_{m_t}(\Phi_t \hat{\mathbf{g}}_t^{(i)}) - \Phi_t \hat{\mathbf{g}}_t^{(i)}$  are mean 0 and independent random variables. The last inequality follows from assumption on the contraction operator  $C$ , Assumption A in (17). We also use the guarantees from Theorem A.3 or Theorem A.4 to remove  $\Phi_t$ . Using the fact that  $\beta_t \leq 1/4$ , completes this part.

For Part (c) (expectation over all randomness in Algorithm COMPDISTSGD),

$$\begin{aligned}
&\mathbb{E}[\|\frac{1}{\kappa} \sum_{i \in R_t} C_{m_t}(\Phi_t \hat{\mathbf{g}}_t^{(i)}) - \Phi_t \nabla F(\mathbf{w}_t)\|^2] \\
&= \mathbb{E}[\|\frac{1}{\kappa} \sum_{i \in R_t} C_{m_t}(\Phi_t \hat{\mathbf{g}}_t^{(i)}) - \frac{1}{\kappa} \sum_{i=1}^{\kappa} \Phi_t \nabla F(\mathbf{w}_t) - \frac{1}{\kappa} \sum_{i \in R_t} \Phi_t \nabla f_i(\mathbf{w}_t) + \frac{1}{\kappa} \sum_{i \in R_t} \Phi_t \nabla f_i(\mathbf{w}_t)\|^2] \\
&= \mathbb{E}[\|\frac{1}{\kappa} \sum_{i \in R_t} C_{m_t}(\Phi_t \hat{\mathbf{g}}_t^{(i)}) - \frac{1}{\kappa} \sum_{i \in R_t} \Phi_t \nabla f_i(\mathbf{w}_t)\|^2] + \mathbb{E}[\|\frac{1}{\kappa} \sum_{i \in R_t} \Phi_t \nabla f_i(\mathbf{w}_t) - \frac{1}{\kappa} \sum_{i=1}^{\kappa} \Phi_t \nabla F(\mathbf{w}_t)\|^2] \\
&= \frac{1}{\kappa^2} \mathbb{E}[\|\sum_{i \in R_t} C_{m_t}(\Phi_t \hat{\mathbf{g}}_t^{(i)}) - \sum_{i \in R_t} \Phi_t \nabla f_i(\mathbf{w}_t) - \sum_{i \in R_t} \Phi_t \hat{\mathbf{g}}_t^{(i)} + \sum_{i \in R_t} \Phi_t \hat{\mathbf{g}}_t^{(i)}\|^2] + \frac{1}{\kappa^2} \mathbb{E}[\|\sum_{i \in R_t} \Phi_t \nabla f_i(\mathbf{w}_t) - \sum_{i=1}^{\kappa} \Phi_t \nabla F(\mathbf{w}_t)\|^2] \\
&= \frac{1}{\kappa} \sum_{i \in R_t} \mathbb{E}[\|C_{m_t}(\Phi_t \hat{\mathbf{g}}_t^{(i)}) - \Phi_t \hat{\mathbf{g}}_t^{(i)}\|^2] + \frac{1}{\kappa} \sum_{i \in R_t} \mathbb{E}[\|\Phi_t \hat{\mathbf{g}}_t^{(i)} - \Phi_t \nabla f_i(\mathbf{w}_t)\|^2] + \frac{1}{\kappa} \sum_{i \in R_t} \mathbb{E}[\|\nabla f_i(\mathbf{w}_t) - \nabla F(\mathbf{w}_t)\|^2] \\
&\leq (1 + \beta_t) (\chi_{C_{m_t}} G^2 + \zeta^2 + \vartheta^2),
\end{aligned}$$

where the last inequality follows from assumption on the contraction operator  $C$ , Assumptions A, B, C in (17). We also use the guarantees from Theorem A.3 or Theorem A.4 to remove  $\Phi_t$ . Using the fact that  $\beta_t \leq 1/4$ , completes this part.  $\square$

The following results follow from Theorems B.7, B.8, B.9, 2.4, and B.5 using the bounds from Lemma D.1. For simplicity, we assume that  $\chi_{C_a}$  is a non-decreasing function in  $a$ , an assumption that is true for all known contraction operators.

**Corollary D.2.** Consider the optimization problem  $\min_{\mathbf{w} \in \mathcal{C}} F(\mathbf{w}) := \frac{1}{M} \sum_{i=1}^M f_i(\mathbf{w})$ . Let  $G'^2 = 2(\chi_{C_d} G^2 / \kappa + G^2)$ . Let  $\zeta'^2 = 2(\chi_{C_d} G^2 + \zeta^2 + \vartheta^2)$ .

1. Let Assumption A in (17) hold. Let  $F$  be  $\lambda$ -strongly convex and  $\mu$ -smooth over  $\mathcal{C}$ , then Algorithm COMPDISTSGD with  $\eta_t = 2/(\lambda t)$  and  $\beta_t = 1/t$  satisfies:  $\mathbb{E}[F(\mathbf{w}_T) - F(\mathbf{w}^*)] = O(\mu G'^2 / (\lambda^2 T))$ .
2. Let Assumption A in (17) hold. Let  $F$  be  $\lambda$ -strongly convex over  $\mathcal{C}$ , then Algorithm COMPDISTSGD with  $\eta_t = 2/(\lambda t)$  and  $\beta_t = 1/t$  satisfies:  $\mathbb{E}[F(\mathbf{w}_T) - F(\mathbf{w}^*)] = O(G'^2 (1 + \log(T)) / (\lambda T))$ .
3. Let Assumption A in (17) hold. Let  $F$  be a convex function over  $\mathcal{C}$ , then Algorithm COMPDISTSGD with  $\eta_t = \|\mathcal{C}\| / (G' \sqrt{t})$  and  $\beta_t = 1/t$  satisfies:  $\mathbb{E}[F(\mathbf{w}_T) - F(\mathbf{w}^*)] = O(\|\mathcal{C}\| G' \log T / \sqrt{T})$ .
4. Let Assumptions A, B, C in (17) hold. Let  $F$  be a  $\mu$ -smooth and continuously differentiable function over  $\mathcal{C}$ . Let  $\rho \in (0, 1)$  and  $\alpha > 0$ . Consider Algorithm COMPDISTSGD with  $\eta_t = 1/\mu$ ,  $\beta_t = \min\{1/4, (\mu \alpha^2) / (64 G' \|\mathcal{C}\| (G + \mu \|\mathcal{C}\|)^2)\}$ , and  $\kappa = \Omega((1/\rho) \cdot \max\{1, (\|\mathcal{C}\| \zeta' (G + \mu \|\mathcal{C}\|)^2 / (\alpha^2 \mu))^2\})$ . Output the first  $\mathbf{w}_\tau$  (if it exists) in  $(\mathbf{w}_1, \dots, \mathbf{w}_T)$  such that  $\|\mathbf{w}_\tau - \mathbf{w}_{\tau+1}\| \leq \alpha / (G + \mu \|\mathcal{C}\|)$  and  $\perp$  (fail) otherwise. Then, if  $T = \Omega((F(\mathbf{w}_1) - F^*) (G + \mu \|\mathcal{C}\|)^2 / (\mu \alpha^2 \rho))$ , with probability at least  $(1 - \rho)^2$ , this procedure outputs a  $\mathbf{w}_\tau$  that is an  $\alpha$ -FOSP for  $F$ .
5. Let Assumptions A, B, C in (17) hold. Let  $F$  be a  $\mu$ -smooth and continuously differentiable function over  $\mathcal{C}$ . Let  $\alpha > 0$ . Consider Algorithm COMPDISTSGD with  $\eta_t = 1/(2\mu)$ ,  $\beta_t = \min\{1/4, \alpha^2 / (\mu G' \|\mathcal{C}\|)\}$ , and

$\kappa = \Omega(\max\{1, \zeta'^2/\alpha^2, \zeta'^2\|\mathcal{C}\|^2\mu^2/\alpha^4\})$ . Output a random iterate  $\mathbf{w}_k$  from  $\{\mathbf{w}_1, \dots, \mathbf{w}_T\}$ . Then, if  $T = \Omega(\mu(F(\mathbf{w}_1) - F^*)/\alpha^2)$ ,  $\mathbf{w}_k$  satisfies the  $\alpha$ -compressed gradient mapping condition in expectation.

*Proof Sketch.* For Parts 1, 2, and 3, the following variant of Lemma 2.2 holds. For any  $t \in [T]$ ,

$$(1 - \beta_t) \mathbb{E}[\|\mathbf{w}_{t+1} - \mathbf{w}\|^2] \leq \mathbb{E}[\|\mathbf{w}_t - \eta_t \nabla F(\mathbf{w}_t) - \mathbf{w}\|^2] + \eta_t^2 \left( \frac{\chi C}{\kappa} + 1 \right) G^2.$$

The remainder of the proofs for Parts 1, 2, and 3 follow identical to Theorems B.7, B.8, B.9, respectively (after replacing Lemma 2.2 with the above inequality).

For Part 4, assuming  $\mathbf{w}_\tau$  exists, we can derive an expression similar to (8),

$$\begin{aligned} & \min_{\mathbf{w} \in \mathcal{C}} \langle \nabla F(\mathbf{w}_\tau), \mathbf{w} - \mathbf{w}_\tau \rangle \\ & \geq -(1 + \beta) \left( \mu \|\mathcal{C}\| \|\mathbf{w}_\tau - \mathbf{w}_{\tau+1}\| + G \|\mathbf{w}_\tau - \mathbf{w}_{\tau+1}\| + \|\mathcal{C}\| \left\| \Phi_\tau \nabla F(\mathbf{w}_\tau) - \frac{1}{\kappa} \sum_{i \in R_\tau} C_{m_\tau}(\Phi_\tau \hat{\mathbf{g}}_\tau^{(i)}) \right\| \right). \end{aligned}$$

We can use Lemma D.1 and follow proof of Theorem 2.4.

Part 5 follows from Theorem B.5 using similar ideas.  $\square$

To take the blow-up in variance into account, in Table 5, we summarize the total communication cost (summed over all rounds) for: i) distributed synchronous SGD, ii) distributed synchronous SGD with only the contraction operator, and iii) Algorithm COMPDISTSGD which utilizes both the lower-dimensional gradients and the contraction operator. Note that Algorithm COMPDISTSGD improves the communication costs for many constraint sets, and is never more than that the case of distributed synchronous SGD with only the contraction operator.

Assumptions on $F$ and Guarantee	SGD (Total comm. cost)	SGD with Contraction $C$ (Total comm. cost)	Algorithm COMPDISTSGD (Total comm. cost)
$\lambda$ -strongly convex, $\mu$ -smooth $\mathbb{E}[F(\mathbf{w}_T) - F(\mathbf{w}^*)] \leq \alpha$	$O\left(\frac{\mu G^2}{\lambda^2 \alpha} d\right)$ Rakhlín et al. [2011]	$O\left(\frac{\mu \chi C_d G^2}{\lambda^2 \alpha} \gamma_{C_d}\right)$	$O\left(\frac{\mu \chi C_d G^2}{\lambda^2 \alpha} \min\{\gamma_{C_d}, \gamma_{C_r}\}\right)$ (where $r = \left(\frac{\omega(C) \mu \chi C_d G^2}{\lambda^2 \alpha}\right)^2$ ) (From Corollary D.2, Part 1)
$\lambda$ -strongly convex $\mathbb{E}[F(\mathbf{w}_T) - F(\mathbf{w}^*)] \leq \alpha$	$\tilde{O}\left(\frac{G^2}{\lambda \alpha} d\right)$ Shamir and Zhang [2013]	$\tilde{O}\left(\chi C_d \frac{G^2}{\lambda \alpha} \gamma_{C_d}\right)$	$\tilde{O}\left(\frac{\chi C_d G^2}{\lambda \alpha} \min\{\gamma_{C_d}, \gamma_{C_r}\}\right)$ (where $r = \left(\frac{\omega(C) \chi C_d G^2}{\lambda \alpha}\right)^2$ ) (From Corollary D.2, Part 2)
Convex $\mathbb{E}[F(\mathbf{w}_T) - F(\mathbf{w}^*)] \leq \alpha$	$\tilde{O}\left(\frac{G^2 \ \mathcal{C}\ ^2}{\alpha^2} d\right)$ Shamir and Zhang [2013]	$\tilde{O}\left(\frac{\chi C_d G^2 \ \mathcal{C}\ ^2}{\alpha^2} \gamma_{C_d}\right)$	$\tilde{O}\left(\frac{\chi C_d G^2 \ \mathcal{C}\ ^2}{\alpha^2} \min\{\gamma_{C_d}, \gamma_{C_r}\}\right)$ (where $r = \left(\frac{\omega(C) \chi C_d G^2 \ \mathcal{C}\ ^2}{\alpha^2}\right)^2$ ) (From Corollary D.2, Part 3)
$\mu$ -smooth and diff. nonconvex $\alpha$ -FOSP	$O(\kappa_1 T d)$ Mokhtari et al. [2018]	$O(\kappa T \gamma_{C_d})$	$O(\kappa T \min\{\gamma_{C_d}, \gamma_{C_r}\})$ (where $r = \frac{\omega(C)^2}{\beta^2}$ ) (From Corollary D.2, Part 4)
$\mu$ -smooth and diff. nonconvex $\alpha$ -comp. grad. mapping	$O(\kappa_1 T d)$ Ghadimi et al. [2016]	$O(\kappa T \gamma_{C_d})$	$O(\kappa T \min\{\gamma_{C_d}, \gamma_{C_r}\})$ (where $r = \frac{\omega(C)^2}{\beta^2}$ ) (From Corollary D.2, Part 5)

Table 5: Comparison of the total communication costs of distributed synchronous SGD, distributed synchronous SGD with contraction operator  $C$ , and our proposed Algorithm COMPDISTSGD for solving (15). For the first 3 rows, we set  $\kappa = 1$ . The  $\tilde{O}(\cdot)$  notation hides some logarithmic terms. The settings of  $\kappa, \beta, T$  for the last two rows are from Corollary D.2, Parts 4 and 5 respectively. Here,  $\kappa_1$  is the value of  $\kappa$  obtained by setting  $\chi C_d = 1$ .

## D.1 STOCHASTIC ROUNDING TECHNIQUE OF [Alistarh et al., 2017]

As an illustration, we present the stochastic rounding technique which is a popular example of a contraction operator. Let  $s \geq 1$  be a tuning parameter. For any  $\mathbf{v} = (\mathbf{v}_1, \dots, \mathbf{v}_a) \in \mathbb{R}^a$  with  $\mathbf{v} \neq 0$ ,  $C_a(\mathbf{v}, s)$  is defined as:

$$C_a(\mathbf{v}, s) = \|\mathbf{v}\| \cdot \text{sign}(\mathbf{v}_i) \cdot \xi_i(\mathbf{v}, s),$$

where  $\xi_i(\mathbf{v}, s)$ 's are independent random variables defined as follows. Let  $0 \leq l < s$  be an integer such that  $|\mathbf{v}_i|/\|\mathbf{v}\| \in [l/s, (l+1)/s]$ . Then  $\xi_i(\mathbf{v}, s) = l/s$  with probability  $1 - p(|\mathbf{v}_i|/\|\mathbf{v}\|, s)$  and  $(l+1)/s$  otherwise. Here,  $p(a, s) = as - l$  for any  $a \in [0, 1]$ . If  $\mathbf{v} = 0$ , then  $C_a(\mathbf{v}, s) = 0$ .

From [Alistarh et al., 2017, Lemma 3.1], for any vector  $\mathbf{v} \in \mathbb{R}^a$ , we have  $\mathbb{E}[C_a(\mathbf{v}, s)] = \mathbf{v}$ ,  $\mathbb{E}[\|C_a(\mathbf{v}, s) - \mathbf{v}\|^2] \leq \min\{a/s^2, \sqrt{a}/s\} \|\mathbf{v}\|^2$ , and  $\mathbb{E}[\|C_a(\mathbf{v}, s)\|_0] \leq s(s + \sqrt{a})$  (sparsity). The authors in [Alistarh et al., 2017] combine this quantization with Elias coding to get better results.

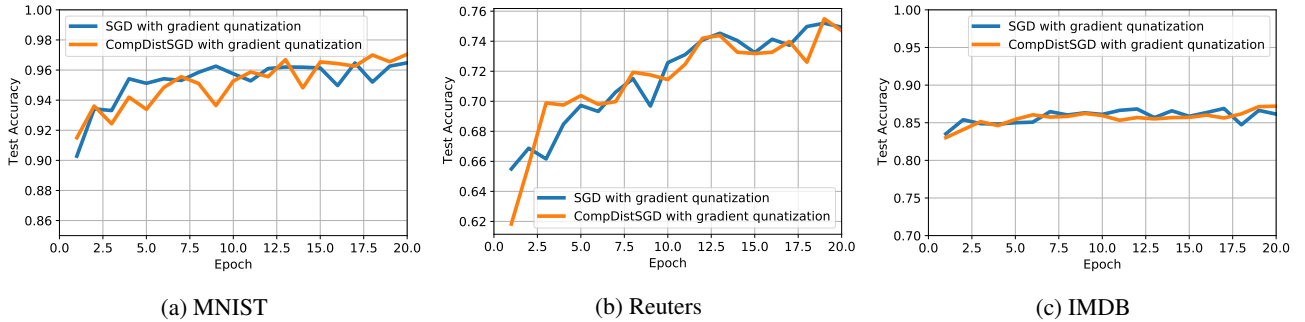


Figure 7: Performance comparison of training a MLP network with SGD vs. COMPDISTSGD on different datasets under the stochastic gradient quantization technique of Alistarh et al. [2017] (without the Elias encoding part). The setup is similar to Section 2.2. The dimensionality of gradients utilized in COMPDISTSGD experiments is a factor 10 smaller than that for SGD. See Figure 2 for results without the quantization.

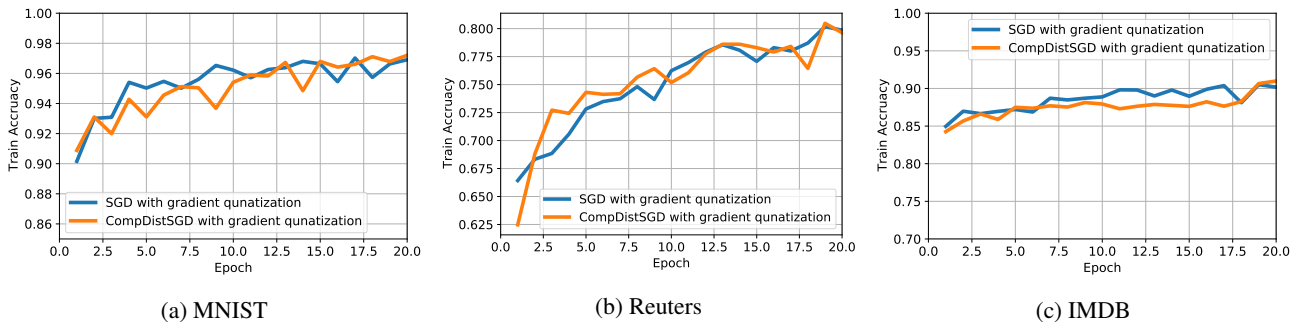


Figure 8: Training accuracy of a MLP network with SGD and COMPDISTSGD on different datasets under the stochastic gradient quantization technique of Alistarh et al. [2017] (without the Elias encoding part). The dimensionality of gradients utilized in COMPDISTSGD experiments is a factor 10 smaller than that for SGD.

## E FRANK-WOLFE WITH LOW-DIMENSIONAL GRADIENTS

The Frank-Wolfe (FW) optimization algorithm proposed by [Frank and Wolfe, 1956] (also known as *conditional gradient method*), is another popular first-order method to solve (1) while only requiring access to a linear optimization oracle (LOO) over  $\mathcal{C}$ , i.e., the ability to compute efficiently  $\text{LOO}(\mathbf{r}) := \text{argmin}_{\mathbf{s} \in \mathcal{C}} \langle \mathbf{s}, \mathbf{r} \rangle$ . Frank-Wolfe is a remarkably simple algorithm that given an initial guess  $\mathbf{w}_1$  constructs a sequence of estimates  $\mathbf{w}_2, \dots$  that converges towards a solution of the optimization problem. This algorithm has seen an impressive revival in recent years due to its low memory requirement and projection-free iterations.

In this section, we show how one could recover the standard convergence guarantees of Frank-Wolfe algorithm with only access to compressed gradients provided through the CSFO oracle. Algorithm COMPFW presents our scheme. The main difference, compared to a traditional Frank-Wolfe algorithm, is how we invoke the linear optimization oracle. A traditional (stochastic) FW update computes a stochastic approximation to the gradient at the current iterate  $\mathbf{w}_t$  and invokes the LOO oracle with it. This gives the element in  $\mathcal{C}$  that correlates the most with the steepest descent (the negative stochastic gradient). We use a similar idea but instead solve a linear minimization problem with the compressed gradients over the set  $\Phi_t \mathcal{C}$ . We then utilize the lifting idea from (2) to compute  $\mathbf{s}_t \in \mathcal{C}$  which is then used to take a step in the direction dictated by  $\mathbf{s}_t - \mathbf{w}_t$ . Algorithm COMPFW is also projection-free.



In Algorithm COMPFW, the linear minimization problem is over domain  $\Phi_t \mathcal{C}$ . Again, the fact that  $\Phi_t \mathcal{C}$  is a linear transform of  $\mathcal{C}$  could come in handy. For example, if  $\mathcal{C}$  is convex hull (polytope) over  $l$  vectors in  $\mathbb{R}^d$ , then  $\Phi_t \mathcal{C}$  is also a convex hull over  $l$  vectors in  $\mathbb{R}^{m_t}$ . Linear optimization over convex hulls is simple as the minimum is achieved at one of the extreme points.

We start here on the nonconvex case. The following theorem shows that for smooth nonconvex functions Algorithm COMPFW, after  $O(\alpha^{-2})$  iterations, converges to an  $\alpha$ -FOSP (with high probability). Now typical analyses of Frank-Wolfe algorithms rely on bounding the Frank-Wolfe (duality) gap that measures how close the iterate is to be a first-order stationary point. The Frank-Wolfe (FW) gap at iterate  $\mathbf{w}_t$  is defined as:  $\max_{\mathbf{s} \in \mathcal{C}} \langle \mathbf{s} - \mathbf{w}_t, -\nabla F(\mathbf{w}_t) \rangle$ . However, since we do not have access to the true (or even the stochastic gradients) computing this FW gap is not possible. We overcome this issue by constructing an approximation to the true Frank-Wolfe gap through the  $\langle \bar{\vartheta}_t, \Phi_t(\mathbf{s}_t - \mathbf{w}_t) \rangle$  term and using it as a stopping condition. In particular, we show that when  $\langle \bar{\vartheta}_t, \Phi_t(\mathbf{s}_t - \mathbf{w}_t) \rangle \geq -\alpha$ , then we do achieve an  $\alpha$ -FOSP (with high probability).

#### Algorithm COMPFW: Frank-Wolfe with Dimensionality Reduced Gradients

**Objective:**  $\min_{\mathbf{w} \in \mathcal{C}} F(\mathbf{w})$

**Input:** Convex set  $\mathcal{C}$ , learning rate parameters  $\{\eta_t\}$ , and projection dimension parameters  $\{\beta_t\}$ .

1. Pick  $\mathbf{w}_1$  as any point in  $\mathcal{C}$
2. for  $t = 1, \dots$  do
  - a. Set  $m_t \leftarrow \Omega(\omega(\mathcal{C})^2 / \beta_t^2)$
  - b. Let  $\Phi_t \in \mathbb{R}^{m_t \times d} \sim_{i.i.d} \mathcal{N}(0, 1/m_t)$
  - c. Let  $\bar{\vartheta}_t \leftarrow \frac{1}{b} \sum_{i=1}^b \text{CSFO}(\mathbf{w}_t, \Phi_t)$  ( $b$  independent invocations of CSFO oracle)
  - d. Let  $\theta_t \leftarrow \text{argmin}_{\theta \in \Phi_t \mathcal{C}} \langle \theta, \bar{\vartheta}_t \rangle^5$
  - e. Let  $\mathbf{s}_t \leftarrow$  pick any element from  $\mathcal{S}_t = \{\mathbf{w} : \Phi_t \mathbf{w} = \theta_t\}$  (e.g., by solving (2))
  - f.  $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \eta_t(\mathbf{s}_t - \mathbf{w}_t)$

**Theorem E.1.** *Let  $F$  be  $\mu$ -smooth and continuously differentiable function over a convex set  $\mathcal{C}$ . Let the assumptions in (4) hold. Let  $\rho \in (0, 1)$  and  $\alpha > 0$ . Consider Algorithm COMPFW with output of first  $\mathbf{w}_\tau$  (if it exists) in  $(\mathbf{w}_1, \dots, \mathbf{w}_T)$  such that  $\langle \Phi_\tau \theta_\tau, \Phi_\tau(\mathbf{s}_\tau - \mathbf{w}_\tau) \rangle \geq -\alpha$ , and  $\perp$  (fail) otherwise. Set  $\eta_t = \eta = \frac{\alpha}{(16\|\mathcal{C}\|^2\mu)}$ ,  $\beta_t = \beta = \frac{\alpha}{G\|\mathcal{C}\|}$  for all  $t \in [T]$ , and batchsize  $b = \Omega((1/\rho) \cdot \max\{1, \frac{\zeta^2 \|\mathcal{C}\|^2}{\alpha^2}\})$ . Then, if  $T = \Omega(\frac{(F(\mathbf{w}_1) - F^*) \|\mathcal{C}\|^2 \mu}{(\alpha^2 \rho)})$ , with probability at least  $(1 - \rho)^2$ , this procedure outputs a  $\mathbf{w}_\tau$  that is an  $\alpha$ -FOSP for  $F$ .*

*Proof.* We use an idea similar to Theorem 2.4. Again, we establish this result over two steps.

1. If  $\mathbf{w}_\tau$  exists, then with probability at least  $1 - \rho$ ,  $\mathbf{w}_\tau$  is an  $\alpha$ -FOSP.
2. If  $T = \Omega(\frac{(F(\mathbf{w}_1) - F^*) \|\mathcal{C}\|^2 \mu}{(\alpha^2 \rho)})$ , then with probability at least  $1 - \rho$ ,  $\mathbf{w}_\tau$  exists in  $(\mathbf{w}_1, \dots, \mathbf{w}_T)$ .

Note that under our assumptions,  $\|\nabla F(\mathbf{w}_t)\| \leq G$  for all  $t \in [T]$ .

Let us start with the first part. For any  $t \in [T]$ , let  $\mathbf{g}_t = \nabla F(\mathbf{w}_t)$  and  $\bar{\mathbf{g}}_t = \frac{1}{b} \sum_{i=1}^b \hat{\mathbf{g}}_t^{(i)}$  where  $\hat{\mathbf{g}}_t^{(i)}$  are independent stochastic gradients of  $F(\mathbf{w}_t)$  used by CSFO oracle to compute  $\vartheta_t$  in Step 2c of Algorithm COMPFW. Under this notation  $\vartheta_t = \Phi_t \bar{\mathbf{g}}_t$ . Now for any  $\mathbf{w} \in \mathcal{C}$ ,

$$\langle \mathbf{g}_\tau, \mathbf{w} - \mathbf{w}_\tau \rangle = \langle \bar{\mathbf{g}}_\tau, \mathbf{w} - \mathbf{w}_\tau \rangle + \langle \mathbf{g}_\tau - \bar{\mathbf{g}}_\tau, \mathbf{w} - \mathbf{w}_\tau \rangle \geq \langle \bar{\mathbf{g}}_\tau, \mathbf{w} - \mathbf{w}_\tau \rangle - \|\mathcal{C}\| \|\mathbf{g}_\tau - \bar{\mathbf{g}}_\tau\|.$$

Now focusing on the  $\langle \bar{\mathbf{g}}_\tau, \mathbf{w} - \mathbf{w}_\tau \rangle$  term,

$$\begin{aligned} \langle \bar{\mathbf{g}}_\tau, \mathbf{w} - \mathbf{w}_\tau \rangle &\geq \mathbb{E}_{\Phi_\tau} [\langle \Phi_\tau \bar{\mathbf{g}}_\tau, \Phi_\tau(\mathbf{w} - \mathbf{w}_\tau) \rangle] - \beta_\tau \|\bar{\mathbf{g}}_\tau\| \|\mathbf{w} - \mathbf{w}_\tau\| \\ &\geq \mathbb{E}_{\Phi_\tau} [\langle \Phi_\tau \bar{\mathbf{g}}_\tau, \Phi_\tau(\mathbf{s}_\tau - \mathbf{w}_\tau) \rangle] - \beta_\tau \|\bar{\mathbf{g}}_\tau\| \|\mathcal{C}\| \\ &\geq -\alpha - \beta_\tau \|\bar{\mathbf{g}}_\tau\| \|\mathcal{C}\|, \end{aligned}$$

<sup>5</sup>As in common in the Frank-Wolfe literature, see e.g. [Jaggi, 2013], an approximate linear oracle suffices for our results too (omitted here for clarity).

where the second inequality comes the minimization problem solved in Algorithm COMPFW and the last inequality comes by the assumption on  $\mathbf{w}_\tau$  that  $\langle \Phi_\tau \bar{\mathbf{g}}_\tau, \Phi_\tau(\mathbf{s}_\tau - \mathbf{w}_\tau) \rangle = \langle \vartheta_\tau, \Phi_\tau(\mathbf{s}_\tau - \mathbf{w}_\tau) \rangle \geq -\alpha$ . Hence, we have that

$$\langle \mathbf{g}_\tau, \mathbf{w} - \mathbf{w}_\tau \rangle \geq -\alpha - \beta_\tau \|\bar{\mathbf{g}}_\tau\| \|\mathcal{C}\| - \|\mathcal{C}\| \|\mathbf{g}_\tau - \bar{\mathbf{g}}_\tau\|.$$

Since the above discussion holds any  $\mathbf{w} \in \mathcal{C}$ , we get that

$$\min_{\mathbf{w} \in \mathcal{C}} \langle \mathbf{g}_\tau, \mathbf{w} - \mathbf{w}_\tau \rangle \geq -\alpha - \beta_\tau \|\bar{\mathbf{g}}_\tau\| \|\mathcal{C}\| - \|\mathcal{C}\| \|\mathbf{g}_\tau - \bar{\mathbf{g}}_\tau\|.$$

Since  $\mathbb{E}[\|\mathbf{g}_\tau - \bar{\mathbf{g}}_\tau\|^2] \leq \zeta^2/b$ , we obtain from Markov's inequality

$$\Pr[\|\mathbf{g}_\tau - \bar{\mathbf{g}}_\tau\| \leq \alpha/\|\mathcal{C}\|] \geq 1 - \frac{\zeta^2 \|\mathcal{C}\|^2}{b\alpha^2}.$$

Also,  $\|\bar{\mathbf{g}}_\tau\| \leq G$ . With our setting of  $\beta_\tau$ , we get that with probability at least  $1 - \frac{\zeta^2 \|\mathcal{C}\|^2}{b\alpha^2}$ , we have,

$$\min_{\mathbf{w} \in \mathcal{C}} \langle \mathbf{g}_\tau, \mathbf{w} - \mathbf{w}_\tau \rangle \geq -c''\alpha,$$

for some constant  $c''$ . By rescaling  $\alpha$  to  $\alpha/c''$ , and using the lower bound on  $b$ , we get that  $\mathbf{w}_\tau$  satisfies the  $\alpha$ -FOSP.

We now establish the second part which bounds the expected number of iterations  $\mathcal{T}$ , starting from any  $\mathbf{w}_1 \in \mathcal{C}$ , before reaching a  $\langle \vartheta_t, \Phi_t(\mathbf{s}_t - \mathbf{w}_\mathcal{T}) \rangle \geq -\alpha$ . Consider any  $t \leq \mathcal{T}$ . By our assumption,  $\langle \Phi_t \bar{\mathbf{g}}_t, \Phi_t(\mathbf{s}_t - \mathbf{w}_t) \rangle = \langle \vartheta_t, \Phi_t(\mathbf{s}_t - \mathbf{w}_t) \rangle < -\alpha$ . The smoothness condition implies,

$$\begin{aligned} F(\mathbf{w}_{t+1}) &\leq F(\mathbf{w}_t) + \langle \mathbf{g}_t, \mathbf{w}_{t+1} - \mathbf{w}_t \rangle + \frac{\mu}{2} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|^2 \\ &= F(\mathbf{w}_t) + \eta \langle \mathbf{g}_t, \mathbf{s}_t - \mathbf{w}_t \rangle + \frac{\mu}{2} \|\eta(\mathbf{s}_t - \mathbf{w}_t)\|^2 \\ &\leq F(\mathbf{w}_t) + \eta \langle \mathbf{g}_t, \mathbf{s}_t - \mathbf{w}_t \rangle + \frac{\mu\eta^2}{2} \|\mathcal{C}\|^2 \\ &= F(\mathbf{w}_t) + \eta \langle \bar{\mathbf{g}}_t, \mathbf{s}_t - \mathbf{w}_t \rangle + \eta \langle \mathbf{g}_t - \bar{\mathbf{g}}_t, \mathbf{s}_t - \mathbf{w}_t \rangle + \frac{\mu\eta^2}{2} \|\mathcal{C}\|^2 \\ &\leq F(\mathbf{w}_t) + \eta \langle \bar{\mathbf{g}}_t, \mathbf{s}_t - \mathbf{w}_t \rangle + \eta \|\mathbf{g}_t - \bar{\mathbf{g}}_t\| \|\mathcal{C}\| + \frac{\mu\eta^2}{2} \|\mathcal{C}\|^2. \end{aligned} \tag{18}$$

Let us focus on the  $\langle \bar{\mathbf{g}}_t, \mathbf{s}_t - \mathbf{w}_t \rangle$  term.

$$\langle \bar{\mathbf{g}}_t, \mathbf{s}_t - \mathbf{w}_t \rangle \leq \mathbb{E}_{\Phi_t}[\langle \Phi_t \bar{\mathbf{g}}_t, \Phi_t(\mathbf{s}_t - \mathbf{w}_t) \rangle] + \beta \|\mathbf{s}_t - \mathbf{w}_t\| \|\bar{\mathbf{g}}_t\| \leq -\alpha + \beta \|\mathbf{s}_t - \mathbf{w}_t\| \|\bar{\mathbf{g}}_t\| \leq -\alpha + \beta \|\mathcal{C}\| \|\bar{\mathbf{g}}_t\|.$$

Plugging the above inequality in (18),

$$F(\mathbf{w}_{t+1}) \leq F(\mathbf{w}_t) - \eta\alpha + \eta\beta \|\mathcal{C}\| \mathbb{E}[\|\bar{\mathbf{g}}_t\|] + \eta \mathbb{E}[\mathbf{g}_t - \bar{\mathbf{g}}_t] \|\mathcal{C}\| + \frac{\mu\eta^2}{2} \|\mathcal{C}\|^2.$$

Consider  $\mathcal{F}_t$  as the sigma algebra that measures all sources of randomness up to iteration  $t$ . Then taking expectation on both sides,

$$\mathbb{E}[F(\mathbf{w}_{t+1}) | \mathcal{F}_t] \leq F(\mathbf{w}_t) - \eta\alpha + \eta\beta \|\mathcal{C}\| G + \eta \frac{\zeta}{\sqrt{b}} \|\mathcal{C}\| + \frac{\mu\eta^2}{2} \|\mathcal{C}\|^2.$$

Under our setting of  $b$ ,  $\eta$ , and  $\beta$ ,

$$\mathbb{E}[F(\mathbf{w}_{t+1}) | \mathcal{F}_t] \leq F(\mathbf{w}_t) - \frac{\alpha^2}{8\|\mathcal{C}\|^2\mu}.$$

Starting from the above expression, we analyze the expected value of  $\mathbb{E}[F(\mathbf{w}_1) - F(\mathbf{w}_{\mathcal{T}+1})]$  as in (10). We get that

$$\mathbb{E}[F(\mathbf{w}_1) - F(\mathbf{w}_{\mathcal{T}+1})] \geq \frac{\alpha^2}{8\|\mathcal{C}\|^2\mu} \mathbb{E}[\mathcal{T}].$$

This implies that  $\mathbb{E}[\mathcal{T}] \leq \mathbb{E}[F(\mathbf{w}_1) - F(\mathbf{w}_{\mathcal{T}+1})](8\|\mathcal{C}\|^2\mu/\alpha^2)$ . Since  $F^* \leq F(\mathbf{w}_{\mathcal{T}+1})$ , we get that

$$\mathbb{E}[\mathcal{T}] \leq (F(\mathbf{w}_1) - F^*) \frac{8\|\mathcal{C}\|^2\mu}{\alpha^2}.$$

Using Markov's inequality

$$\Pr \left[ \mathcal{T} \leq \frac{8(F(\mathbf{w}_1) - F^*)\|\mathcal{C}\|^2\mu}{\alpha^2\rho} \right] \geq 1 - \rho.$$

This completes both parts of the theorem. Putting them together gives the claimed result.  $\square$

For convex functions too the convergence guarantees of Algorithm COMPFW matches the known results with stochastic Frank-Wolfe algorithm [Hazan and Luo, 2016, Theorem 3].

**Theorem E.2.** *Let  $F$  be a  $\mu$ -smooth convex differentiable function over a convex set  $\mathcal{C}$ . Let the assumptions in (4) hold. Then, with  $\eta_t = \frac{2}{(t+1)}$ ,  $\beta_t = \min\{1/4, \frac{\mu\eta_t\|\mathcal{C}\|}{G}\}$ , and batchsize  $b = \Omega(\max\{1, (\frac{\zeta}{\|\mathcal{C}\|\mu\eta_t})^2\})$ , Algorithm COMPFW satisfies:*

$$\mathbb{E}[F(\mathbf{w}_T) - F(\mathbf{w}^*)] \leq \frac{4\mu\|\mathcal{C}\|^2}{T}.$$

*Proof.* For any  $t \in [T]$ , let  $\mathbf{g}_t = \nabla F(\mathbf{w}_t)$  and  $\bar{\mathbf{g}}_t = \frac{1}{b} \sum_{i=1}^b \hat{\mathbf{g}}_t^{(i)}$  where  $\hat{\mathbf{g}}_t^{(i)}$  are independent stochastic gradients of  $F(\mathbf{w}_t)$  used by CSFO oracle to compute  $\vartheta_t$  in Step 2c of Algorithm COMPFW. Under this notation  $\vartheta_t = \Phi_t \bar{\mathbf{g}}_t$ . We start by utilizing the smoothness assumption on  $F$ .

$$\begin{aligned} F(\mathbf{w}_{t+1}) &\leq F(\mathbf{w}_t) + \langle \mathbf{g}_t, \mathbf{w}_{t+1} - \mathbf{w}_t \rangle + \frac{\mu}{2} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|^2 \\ &= F(\mathbf{w}_t) + \eta_t \langle \mathbf{g}_t, \mathbf{s}_t - \mathbf{w}_t \rangle + \frac{\mu}{2} \|\eta_t (\mathbf{s}_t - \mathbf{w}_t)\|^2 \\ &= F(\mathbf{w}_t) + \eta_t \langle \bar{\mathbf{g}}_t, \mathbf{s}_t - \mathbf{w}_t \rangle + \eta_t \langle \mathbf{g}_t - \bar{\mathbf{g}}_t, \mathbf{s}_t - \mathbf{w}_t \rangle + \frac{\mu}{2} \|\eta_t (\mathbf{s}_t - \mathbf{w}_t)\|^2 \\ &\leq F(\mathbf{w}_t) + \eta_t \langle \bar{\mathbf{g}}_t, \mathbf{s}_t - \mathbf{w}_t \rangle + \eta_t \|\mathbf{g}_t - \bar{\mathbf{g}}_t\| \|\mathcal{C}\| + \frac{\mu\eta_t^2}{2} \|\mathcal{C}\|^2. \end{aligned}$$

Let us focus on the  $\langle \bar{\mathbf{g}}_t, \mathbf{s}_t - \mathbf{w}_t \rangle$  term.

$$\begin{aligned} \langle \bar{\mathbf{g}}_t, \mathbf{s}_t - \mathbf{w}_t \rangle &\leq \mathbb{E}_{\Phi_t} [\langle \Phi_t \bar{\mathbf{g}}_t, \Phi_t (\mathbf{s}_t - \mathbf{w}_t) \rangle] + \beta_t \|\mathbf{s}_t - \mathbf{w}_t\| \|\bar{\mathbf{g}}_t\| \leq \mathbb{E}_{\Phi_t} [\langle \Phi_t \bar{\mathbf{g}}_t, \Phi_t (\mathbf{w}^* - \mathbf{w}_t) \rangle] + \beta_t \|\mathbf{s}_t - \mathbf{w}_t\| \|\bar{\mathbf{g}}_t\| \\ &= \langle \bar{\mathbf{g}}_t, \mathbf{w}^* - \mathbf{w}_t \rangle + \beta_t \|\mathbf{s}_t - \mathbf{w}_t\| \|\bar{\mathbf{g}}_t\| \leq \langle \bar{\mathbf{g}}_t, \mathbf{w}^* - \mathbf{w}_t \rangle + \beta_t \|\mathcal{C}\| \|\bar{\mathbf{g}}_t\|, \end{aligned}$$

where we used the fact  $\mathbb{E}_{\Phi_t} [\langle \Phi_t \bar{\mathbf{g}}_t, \Phi_t (\mathbf{w}^* - \mathbf{w}_t) \rangle] = \langle \bar{\mathbf{g}}_t, \mathbf{w}^* - \mathbf{w}_t \rangle$  as  $\bar{\mathbf{g}}_t, \mathbf{w}^*$  and  $\mathbf{w}_t$  are all independent of  $\Phi_t$ . Also, we used that  $\Phi_t \bar{\mathbf{g}}_t = \vartheta_t$  and  $\Phi_t \bar{\mathbf{g}}_t = \theta_t$  and therefore  $\langle \vartheta_t, \theta_t \rangle = \langle \Phi_t \bar{\mathbf{g}}_t, \Phi_t \mathbf{s}_t \rangle \leq \langle \Phi_t \bar{\mathbf{g}}_t, \Phi_t \mathbf{w}^* \rangle$  due to the minimization problem solved in Algorithm COMPFW. Using this inequality, in the above bound on  $F(\mathbf{w}_{t+1})$  yields

$$F(\mathbf{w}_{t+1}) \leq F(\mathbf{w}_t) + \eta_t (\langle \bar{\mathbf{g}}_t, \mathbf{w}^* - \mathbf{w}_t \rangle + \beta_t \|\mathcal{C}\| \|\bar{\mathbf{g}}_t\|) + \eta_t \|\mathbf{g}_t - \bar{\mathbf{g}}_t\| \|\mathcal{C}\| + \frac{\mu\eta_t^2}{2} \|\mathcal{C}\|^2.$$

Taking expectation on both sides, and using convexity on  $F$ ,

$$\begin{aligned} \mathbb{E}[F(\mathbf{w}_{t+1})] &\leq \mathbb{E}[F(\mathbf{w}_t)] + \eta_t \mathbb{E}[\langle \mathbf{g}_t, \mathbf{w}^* - \mathbf{w}_t \rangle] + \eta_t \beta_t \|\mathcal{C}\| \mathbb{E}[\|\bar{\mathbf{g}}_t\|] + \eta_t \mathbb{E}[\|\mathbf{g}_t - \bar{\mathbf{g}}_t\|] \|\mathcal{C}\| + \frac{\mu\eta_t^2}{2} \|\mathcal{C}\|^2 \\ &= \mathbb{E}[F(\mathbf{w}_t)] + \eta_t \mathbb{E}[F(\mathbf{w}^*) - F(\mathbf{w}_t)] + \eta_t \beta_t \|\mathcal{C}\| \mathbb{E}[\|\bar{\mathbf{g}}_t\|] + \eta_t \mathbb{E}[\|\mathbf{g}_t - \bar{\mathbf{g}}_t\|] \|\mathcal{C}\| + \frac{\mu\eta_t^2}{2} \|\mathcal{C}\|^2 \\ &\leq \mathbb{E}[F(\mathbf{w}_t)] + \eta_t \mathbb{E}[F(\mathbf{w}^*) - F(\mathbf{w}_t)] + \eta_t \beta_t \|\mathcal{C}\| G + \eta_t \frac{\zeta}{\sqrt{b}} \|\mathcal{C}\| + \frac{\mu\eta_t^2}{2} \|\mathcal{C}\|^2. \end{aligned}$$

This can be re-expressed as

$$\mathbb{E}[F(\mathbf{w}_{t+1}) - F(\mathbf{w}^*)] \leq (1 - \eta_t) \mathbb{E}[F(\mathbf{w}_t) - F(\mathbf{w}^*)] + \eta_t \beta_t \|\mathcal{C}\| G + \eta_t \frac{\zeta}{\sqrt{b}} \|\mathcal{C}\| + \frac{\mu\eta_t^2}{2} \|\mathcal{C}\|^2.$$

If  $b = \Omega(\max\{1, (\zeta/(\|\mathcal{C}\|\mu\eta_t))^2\})$  and  $\beta_t = \min\{1/4, \mu\eta_t\|\mathcal{C}\|/G\}$ , then

$$\mathbb{E}[F(\mathbf{w}_{t+1}) - F(\mathbf{w}^*)] \leq (1 - \eta_t) \mathbb{E}[F(\mathbf{w}_t) - F(\mathbf{w}^*)] + \mu\eta_t^2\|\mathcal{C}\|^2.$$

A simple inductive argument now shows that  $\eta_t = 2/(t + 1)$ , we get

$$\mathbb{E}[F(\mathbf{w}_{T+1}) - F(\mathbf{w}^*)] \leq 4\mu\|\mathcal{C}\|^2/(T + 1).$$

□