

DIMINISHING NOISE MAINTAINS DIFFERENTIAL PRIVACY AND ENHANCES CONVERGENCE

Anonymous authors

Paper under double-blind review

ABSTRACT

Differential Privacy (DP) is a well-established framework for training models in distributed settings while safeguarding sensitive information. Although numerous DP algorithms exist, many current solutions inject noise with constant variance to the transmitted gradients, leading to convergence only to a neighborhood of the optimal solution. To address this limitation, we propose an error compensation technique that maintains linear convergence without compromising privacy guarantees. This is achieved through cautious adjusting the noise’s variance through the algorithm iterations. Experimental results validate the effectiveness of our approach.

1 INTRODUCTION

The increasing number of trainable parameters (Villalobos et al., 2022) and dataset sizes lead to more complex machine learning problems that cannot be solved on a single device. One possible solution to this challenge is to distribute the task across several smaller devices and find the solution via local computations and communications (Kairouz et al., 2021). This can be formalized by considering the following optimization problem:

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

In this context, the function $f(x)$ represents the loss function computed over the entire dataset, while $f_i(x)$ corresponds to its component computed on the i -th local device’s data.

Federated Learning (FL) has emerged as a prominent distributed learning paradigm where heterogeneous datasets are spread across numerous personal devices. Consequently, this approach faces several challenges. First, there is the requirement for efficient communication. This can be addressed by transmitting compressed gradients (Alistarh et al., 2017), constructing gradient-based sequences to compensate for approximation errors (Karimireddy et al., 2019) or performing several local steps before model averaging (Stich, 2018). These techniques reduce bandwidth consumption while maintaining overall performance. Another key challenge in FL is personalization, i.e., adapting models to local device states. At the same time, the primary objective remains ensuring **privacy** in personalization—enabling collaborative learning while protecting information stored on personal devices. Although FL inherently provides some privacy by exchanging model updates instead of raw data (Duchi et al., 2014), communications may still reveal sensitive information to adversaries (Zhu et al., 2019). To mitigate this, various approaches have been proposed, such as:

- Anonymization (Majeed & Lee, 2020) (which is cheap, but not always effective (Brasher, 2018)),
- Homomorphic encryption (Acar et al., 2018) (which provides fully privacy, but is drastically difficult to implement and expensive to utilize).

Because of this trade-off, new techniques are created and analyzed, that try to include both protectiveness and simplicity. For instance, flipping labels (Shen et al., 2023) or adding noise (Geng & Viswanath, 2015). These mechanisms can be formalized under **Differential Privacy (DP)** (Dwork, 2006), which provides a versatile framework for privacy-preserving methods that are relatively simple to analyze.

2 RELATED WORK AND CONTRIBUTION

Differential Privacy

The framework of DP (Dwork, 2006; 2008) was developed to address privacy concerns, initially as a formal criterion for dataset security. DP quantifies the risk of an adversary determining whether a specific data point was used in model training or included in a dataset. This framework has been extensively adopted to analyze the privacy guarantees of deep learning algorithms, both in centralized (Bassily et al., 2014; Abadi et al., 2016; Wang et al., 2020; Chen et al., 2021) and distributed settings (McMahan et al., 2018; Andrew et al., 2022; Li et al., 2022).

Its essence lies in privacy budget, which is defined beforehand. Though, we do not achieve fully protection, we determine a certain amount of privacy, which can vary depending on our requirements.

A key development in this field was the DP-SGD algorithm (Abadi et al., 2016), which enforces DP by adding Gaussian noise to gradient updates at each iteration. The noise magnitude required for achieving (ϵ, δ) -DP depends on the sensitivity of the gradients’ magnitude (Dwork, 2006). To bound this sensitivity, DP-SGD employs gradient clipping — a technique now widely used to simplify privacy analysis (Wang et al., 2018; Das et al., 2023).

Though, clipping tend to perform well in modern applications (Zhang et al., 2019), it cooperates poorly with noise injection during training under general assumptions SGD converges only to a neighborhood of the optimum (Koloskova et al., 2023). One may implement adaptive clipping (Andrew et al., 2021; Bu et al., 2023; Pichapati et al., 2019), where threshold diminishes during the training process. However, to our knowledge, existing convergence bounds do not coincide with clipping with the constant radius. Furthermore, in most cases this practice is not theoretically justified.

Recently, novel private algorithms have been introduced (Khirirat et al., 2023; Fatkhullin et al., 2023; Islamov et al., 2025), building on the error-compensating technique (Richtárik et al., 2021). Importantly, this technique was originally developed not for privacy, but for efficient message compression and acceleration in distributed optimization. In the highlighted works, however, the authors replace compression operators with clipping, a standard practice in differential privacy analysis, as noted above.

Compression

The demand for efficient communication during the training process lead to development of the schemes, that reduce the number of transmitted bits information (Khirirat et al., 2018). The main idea is to share not exact information (gradients, local states, etc.), but some compressed version of it.

One of the first well-analyzed algorithm for distributed optimization, QSGD (Alistarh et al., 2017), quantized the local gradients. The problem of this approach is in non-reducible variance of the compressed gradients. It can be illustrated as following: in the global optimum the sum of local model’s gradients equals to zero, but each of the gradients can be arbitrary large. Therefore, since variance of the quantized vector usually depends on the vector’s norm, even at the optimum point, sent gradients can deviate from the real ones by a large margin.

To mitigate this effect, error-compensating approaches were introduced, inspired by the variance reduction methods (Schmidt et al., 2017). The idea is simple. We need to compress not the gradient themselves, but the difference between the gradient and previous estimation. Then, these compressed errors are sent to the server. These approach result in diminishing gradient approximation errors, therefore, their compressed versions are also small.

Initially, these type of methods utilized unbiased compression operators (He et al., 2023), that do not change the compressed vector in expectation. This approach resulted as algorithm DIANA (Mishchenko et al., 2019). Later, biased compressors were also bridged with this approach in EF21 (Richtárik et al., 2021). These days, there are numerous methods that can be considered as error-compensating, for instance, MARINA (Gorbunov et al., 2021), DASHA (Tyurin & Richtárik, 2022) and EF-BV (Condat et al., 2023).

Our Contribution

The remarkable performance of error-compensating methods stems from their use of converging-to-zero compressed messages, whose sizes are bounded. Additionally, most DP techniques rely on additive Gaussian noise with constant variance, even though theoretical results suggest noise proportional to the message size should suffice. By connecting these approaches, we develop a

method that leverages the diminishing magnitude of communications to reduce the noise required for privacy preservation.

To summarize all the contribution in comparison with the SOTA results in this field, we:

1. **Propose a new DP version of EF21 with biased compressors.**

Previous works employed clipping operators in EF21, which facilitated differential privacy analysis but failed to preserve the communication efficiency properties offered by certain biased compressors. We bridge this gap by demonstrating that models can achieve both privacy and bandwidth efficiency during training.

2. **Introduce the concept of diminishing noise in DP setup, achieving the linear convergence.**

The foundation of our error compensation framework lies in the diminishing magnitude of transmitted messages throughout the optimization process. As these messages converge to zero, their sensitivity similarly decreases, allowing us to employ noise with progressively reducing variance without affecting convergence. We theoretically demonstrate that this diminishing noise preserves linear convergence under the PL condition, matching the convergence rate of the original EF21 algorithm. Notably, this convergence guarantee was not established in prior works on differentially private error compensation methods.

3. **Validate theoretical results experimentally.**

We compare our method with existing on CIFAR-10 dataset, a well-established benchmark for the optimization algorithms.

3 PRELIMINARIES AND DEFINITIONS

Notation. We use the standard Euclidean norm for vectors: $\|x\| \stackrel{\text{def}}{=} \langle x, x \rangle^{1/2}$, $x \in \mathbb{R}^d$. The objective functional $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is a differentiable function. We denote its global minimum by $f_* \stackrel{\text{def}}{=} \inf_{x \in \mathbb{R}^d} f(x) > -\infty$. We also introduce the gradient of f at point x as $\nabla f(x) \in \mathbb{R}^d$.

Below we introduce all the definitions that will be used throughout the manuscript.

Definition 1 (Smoothness). *Every f_i has L_i -Lipschitz gradient, i.e.*

$$\|\nabla f_i(x) - \nabla f_i(y)\| \leq L_i \|x - y\| \quad \forall x, y \in \mathbb{R}^d.$$

Furthermore, $L^2 \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n L_i^2$.

Definition 2 (PL condition). *There exists $\mu > 0$, such that*

$$f(x) - f_* \leq \frac{1}{2\mu} \|\nabla f(x)\|^2 \quad \forall x \in \mathbb{R}^d.$$

These assumptions are standard in stochastic optimization. Smoothness plays a central role in analyzing optimization algorithms in ML: it allows bounding the decrease in loss after each gradient update and enables proper step-size selection (Nesterov, 2013). Regarding PL condition, recently (Liu et al., 2022), it was shown, that over-parametrized neural networks are locally PL almost everywhere, which ensures faster convergence of gradient-based methods. It justifies our analysis, that heavily relies on this property.

All the necessary privacy prerequisites, that we refer to in the analysis, are introduced further:

Definition 3 ((ϵ, δ) -DP (Dwork, 2006)). *A randomized mechanism $\mathcal{M} : \mathcal{D} \rightarrow \mathcal{R}$ with domain \mathcal{D} and range \mathcal{R} is (ϵ, δ) -DP if for all adjacent (different within one element) datasets $D, D' \in \mathcal{D}$ and for all events $S \in \mathcal{R}$ in the output space of \mathcal{M} it holds*

$$\mathbb{P}\{\mathcal{M}(D) \in S\} \leq e^\epsilon \mathbb{P}\{\mathcal{M}(D') \in S\} + \delta.$$

Definition 4 (Sensitivity). *Function $f : \mathcal{D} \rightarrow \mathbb{R}^d$ is said to have sensitivity Δ , if for adjacent datasets $D' \sim D$ we have*

$$\Delta^2 = \max_{D \sim D'} \|\mathcal{M}(D) - \mathcal{M}(D')\|^2.$$

Lemma 1 (Gaussian Mechanism (GM) for DP). *Let $f : \mathcal{X}^n \rightarrow \mathbb{R}^d$ has sensitivity Δ . Define a randomized algorithm $\mathcal{M} : \mathcal{X}^n \rightarrow \mathbb{R}^d$ by $\mathcal{M}(x) = \mathcal{N}\left(f(x), \frac{2 \log 1.25/\delta}{\epsilon^2} \Delta^2 I_d\right)$. Then \mathcal{M} is (ϵ, δ) -DP.*

Definition 5 (Privacy Loss). Let P and Q be two probability distributions on \mathcal{X} . Define $f_{P||Q} : \mathcal{X} \rightarrow \mathbb{R}$ by $f_{P||Q}(y) = \log \frac{P(y)}{Q(y)}$. The privacy loss is a random variable $\text{PrivLoss}(P||Q) \stackrel{\text{def}}{=} f_{P||Q}(Y)$, where $Y \sim P$.

Definition 6 (Concentrated Differential Privacy). A randomized mechanism $\mathcal{M} : \mathcal{D} \rightarrow \mathcal{R}$ with domain \mathcal{D} and range \mathcal{R} is ρ -zCDP if for all adjacent (different within one sample) datasets $D, D' \in \mathcal{D}$ the privacy loss distribution is well-defined and

$$\mathbb{E}_{Z \sim \text{PrivLoss}(\mathcal{M}(D)||\mathcal{M}(D'))} \exp(tZ) \leq \exp(t(t+1) \cdot \rho), \quad \forall t \geq 0.$$

Lemma 2 (GM for Concentrated DP). Let $f : \mathcal{X}^n \rightarrow \mathbb{R}^d$ has sensitivity Δ . Define a randomized algorithm $\mathcal{M} : \mathcal{X}^n \rightarrow \mathbb{R}^d$ by $\mathcal{M}(x) = \mathcal{N}\left(f(x), \frac{\Delta^2}{2\rho} I_d\right)$. Then \mathcal{M} is ρ -zCDP.

Lemma 3 (Composition for Concentrated DP). Let $M_1, M_2, \dots, M_k : \mathcal{X}^n \rightarrow \mathcal{R}$ be randomized algorithms. Suppose, M_i is ρ_i -zCDP for each i . Define $\mathcal{M}(x) : \mathcal{X}^n \rightarrow \mathcal{R}^k$ by $\mathcal{M}(x) \stackrel{\text{def}}{=} (M_1(x), M_2(x), \dots, M_k(x))$, where each algorithm is run independently. Then, \mathcal{M} is ρ -zCDP for $\rho = \sum_{i=1}^k \rho_i$.

Lemma 4 (Conversion from Concentrated DP to DP). ρ -zCDP implies $(\varepsilon = \rho + 2\sqrt{\rho \log 1/\delta}, \delta)$ -DP for all $\delta > 0$. Also, to obtain a given target (ε, δ) -DP, it suffices to have ρ -zCDP with $\rho \in \left[\frac{\varepsilon^2}{4 \log 1/\delta + 4\varepsilon}, \frac{\varepsilon^2}{4 \log 1/\delta} \right]$

Let us discuss these claims. DP condition (Definition 3) bounds the probability, that random algorithm differs much on dataset and adjacent one. One of the simplest in implementation technique to provide the privacy is adding Gaussian noise (Lemma 1), which is anisotropic and has variance, proportional to the norm of sensitivity (Definition 4). This is the reason for utilizing clipping in existing DP approaches, since it naturally bounds the sensitivity.

When the adversary has access not only to one message, but instead to some finite numbers of them, he can combine them somehow, therefore we need to protect not the single messages only. It appears, that DP guarantees on single queries may be advanced to its union (Lemma 7).

The problem with Approximate DP is in sophisticated analysis of composition theorems. In previous works (Li et al., 2022; Islamov et al., 2025), they assumed the constant privacy budget per iteration. To consider a more general case, we investigate the framework of concentrated differential privacy (Definition 6), that has not been conducted before in optimization manuscripts. It is similar to approximate DP, as it also is provided via the Gaussian Mechanism (Lemma 6) and is convertible to (ε, δ) -DP (Lemma 8).

Having discussed the privacy preliminaries we continue to the compression property, as this is the cornerstone of efficient distributed methods, allowing to reduce the communication costs.

Our proposed algorithm relies on biased compressors satisfying the contraction property – a fundamental requirement for convergence in communication-efficient distributed optimization. Unlike unbiased compressors that preserve expectation but may increase the vector’s norm drastically, these operators provide controlled compression by consistently maintaining the original vector’s direction while bounding the relative magnitude reduction.

Definition 7 (Biased Compressor). Mapping $\mathcal{C} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is called a biased compressor, i.e.:

$$\|\mathcal{C}(x) - x\|^2 \leq (1 - \alpha) \|x\|^2, \quad \forall x \in \mathbb{R}^d$$

Some of the most practical biased compressors, are sparsifications, such as TopK , and quantizations, as Biased Roundings.

Example 1 (TopK (Stich et al., 2018)). Greedy sparsification is defined as

$$\mathcal{C}(x) = \sum_{i=d-k+1}^d x_{(i)} e_{(i)},$$

where coordinates are ordered by their absolute value, so that $|x_{(1)}| \leq |x_{(2)}| \leq \dots \leq |x_{(d)}|$ and e_1, \dots, e_d are a standard Euclidean basis. It satisfies Definition 7 with $\alpha = \frac{k}{d}$.

Example 2 (Biased Rounding (Beznosikov et al., 2024)). Let $\{a_k\}_{k \in \mathbb{Z}}$ be an arbitrary increasing sequence of positive numbers, such that $\inf_k a_k = 0$ and $\sup_k a_k = \infty$. Then, general biased rounding

is defined via

$$\mathcal{C}(x)_i = \text{sign}(x_i) \arg \min_{k \in \mathbb{Z}} |a_k - |x_i||, \quad i \in [d].$$

Then, this operator satisfies Definition 7 with $1/\alpha = \sup_{k \in \mathbb{Z}} \frac{(a_k + a_{k+1})^2}{4a_k a_{k+1}}$.

Examples above provide compressors, that allows to reduce the number of transmitted bits. The first one, transmits k coordinates instead of the problem dimension, d . Another one is used in the quantization, where we store the model more efficiently.

4 MAIN PART

4.1 NOISE AFTER COMPRESSION

Having established the necessary background, we can now proceed to the main theoretical contribution of our paper. We introduce the method DPd-EF21 (Differential Private diminishing EF21):

Algorithm 1 DPd-EF21

- 1: **Parameters:** starting point $x^0 \in \mathbb{R}^d$, $g_i^0 = 0 \in \mathbb{R}^d$, learning rate $\gamma > 0$.
- 2: **for** $t = 0, 1, 2, \dots, T - 1$ **do**
- 3: Send x^t to nodes
- 4: **for** all nodes $i = 1, 2, \dots, n$ in parallel **do**
- 5: Compress $\Delta_i^t = \mathcal{C}(\nabla f_i(x^t) - g_i^{t-1})$
- 6: Update local state $g_i^t = g_i^{t-1} + \Delta_i^t$
- 7: Send $\Delta_i^t + \mathcal{N}(0, \sigma_{i,t}^2 I_d)$ to the server
- 8: **end for**
- 9: Server computes

$$g^t = g^{t-1} + \frac{1}{n} \sum_{i=1}^n [\Delta_i^t + \mathcal{N}(0, \sigma_{i,t}^2 I_d)] \quad (1)$$

- 10: Update the model $x^{t+1} = x^t - \gamma_t g^t$
 - 11: **end for**
-

We utilize the EF21 approach (Richtárik et al., 2021), transmitting the compressed error between the local estimation and the exact local gradient (line 5). This allows to achieve better convergence, than by transferring gradients, since compressed estimation artifacts tend to converge to zero, that cannot be claimed for the local gradients.

The core idea of the convergence proof is analyzing the Lyapunov function

$$V^t = f(x^t) - f_* + \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^t) - g_i^{t-1}\|^2. \quad (2)$$

In the original EF21 paper, Lyapunov equation converge linearly to zero, due to the absence of the stochastic terms:

$$V^t \leq (1 - \gamma\mu)^t V^0. \quad (3)$$

Adding a noise with constant variance σ^2 (line 7) is not helpful in terms of convergence, since we will attract to some neighbourhood, instead of the exact solution (Fatkhullin et al., 2021):

$$\mathbb{E}V^t \leq (1 - \gamma\mu)^t \mathbb{E}V^0 + \frac{d\gamma\sigma^2}{\mu}. \quad (4)$$

The novelty we propose is adjusting $\sigma_{i,t}^2$ during the iteration process. To preserve a certain amount on privacy, variance of injected noise should linearly depends on the sensitivity of the mapping we are aiming to protect. In our scenario, this is

$$\Delta_i^2 = \max_{D \sim D'} \|\mathcal{C}(\nabla f_{\mathcal{D},i}(x^t) - g_{\mathcal{D},i}^{t-1}) - \mathcal{C}(\nabla f_{\mathcal{D}',i}(x^t) - g_{\mathcal{D}',i}^{t-1})\|^2. \quad (5)$$

Since we send $\mathcal{C}(\nabla f_i(x^t) - g_i^{t-1})$, we can bound this sensitivity from above by corresponding Lyapunov functions, that depend on their dataset:

$$\begin{aligned} \Delta_t^2 &\lesssim \left\| \mathcal{C}(\nabla f_{\mathcal{D},i}(x^t) - g_{\mathcal{D},i}^{t-1}) \right\|^2 + \max_{\mathcal{D} \sim \mathcal{D}'} \left\| \mathcal{C}(\nabla f_{\mathcal{D}',i}(x^t) - g_{\mathcal{D}',i}^{t-1}) \right\|^2 \\ &\leq (1 - \alpha) \left\| \nabla f_{\mathcal{D},i}(x^t) - g_{\mathcal{D},i}^{t-1} \right\|^2 + (1 - \alpha) \max_{\mathcal{D} \sim \mathcal{D}'} \left\| \nabla f_{\mathcal{D}',i}(x^t) - g_{\mathcal{D}',i}^{t-1} \right\|^2 \\ &\lesssim V_{\mathcal{D}}^t + V_{\mathcal{D}_1}^t. \end{aligned}$$

As we have $V^t \downarrow$ from the equation (4), we expect Δ_t^2 also to decrease. Since theory proposes, that sensitivity and noise's variance are correlated, we have $\sigma_i^2 \downarrow$. Thus, from the equation (4) we derive $V^t \rightarrow 0$ after carefully examining all the terms.

It turns out, that we can apply different strategies in terms of noise adding, depending on the privacy per iteration. This is obtained by advanced composition techniques, that allow various changing DP guarantees per iteration to get a needed privacy overall.

In our analysis, we aim to compound facts above into adding noise with reducing noise in order to maintain the linear convergence in PL scenario. We consider different strategies and end up with the following theorem:

Theorem 1. *Let function f be L -smooth and satisfy PL condition with constant μ . Define $c_{up} = \frac{4-\alpha}{4-2\alpha}$, γ_{EF} as following:*

$$\gamma_{EF} = \min \left\{ \frac{1}{2L \left(1 + \sqrt{\frac{2c_{up}}{8-2\alpha-4c_{up}\alpha}} \left(2 + \frac{4}{\alpha} \right) \right)}; \frac{\alpha}{2\mu} \right\}.$$

With $T = \Omega\left(\frac{C_1 \cdot C_2}{n}\right)$ there exists a sequence of stepsizes $\{\gamma_k\}_k$, satisfying

$$\frac{C_1 \cdot C_2 \cdot (1 + L\gamma_k + C_3\gamma_k^2)}{n \left(1 - \frac{\gamma_k\mu}{2} \right) (T + 1) \exp\left(\frac{(k+1)\gamma_k\mu}{4}\right)} \leq 1, \quad \gamma_k \leq \gamma_{EF}, \quad \gamma_{k-1} \leq \gamma_k \leq c_{up}\gamma_{k-1},$$

where $C_1 = C_1(\varepsilon, \delta)$ is the constant, depending on the privacy and $C_2 = C_2(L, \mu, \alpha, d, p)$, as well as $C_3 = C_3(L, \alpha, c_{up})$ – constants, depending on the problem, such that we have following convergence result:

$$V^{T+1} \leq 2 \prod_{k=0}^T \left(1 - \frac{\gamma_k\mu}{4} \right) \cdot \bar{V}^0$$

with high probability $1 - p$, where $V^k = f(x^k) - f_* + \frac{2\gamma_k}{\alpha n} \sum_{i=1}^n \|\nabla f_i(x^k) - g_i^k\|^2$ and $\bar{V}^0 =$

$$\max_{\mathcal{D}'} \left[f(x^0) - f_* + \frac{2\gamma_0}{\alpha n} \sum_{i=1}^n \|\nabla f_i(x^0) - g_i^0\|^2 \right], \text{ where maximum is taken across all considered}$$

datasets for this problem. Furthermore, Algorithm 1 will be (ε, δ) -DP.

We analyze the behaviour of the step sizes γ_k . After γ_0 is found, all others stepsizes inevitably exists – one can check, that γ_k is eligible for the $k + 1$ iteration. It can be noted, that our method have two phases of working: a warming-up DP regime and the main EF2.1 one, depending on the amount of iterations. Initial step sizes are in the warming-up regime, therefore, small. After several iterations, as we have $\exp\left(\frac{k\mu\gamma}{4}\right)$ in the denominator, the fraction will be close to zero, hence, we will operate in the main regime with $\gamma \sim \frac{1}{L}$. This is similar to the clipping procedure, where we intentionally reduce the steps. Moreover, as shown in (Khirirat et al., 2023), (Islamov et al., 2025), clipping is not active after a certain amount of iterations, therefore steps are not mitigated. Algorithm 1 behaves the same, where we end up with a constant step size after increasing for the warm-up phase. However, highlighted methods still proceeded to inject noise with constant variance, and not obtaining the exact solution, unlike our proposed framework.

Therefore, the convergence process of our algorithm can be divided into two parts. In the warming up DP regime we will conduct less, than some T_0 steps, since the exponent will have the major impact after several iterations. The second regime – the EF one, inherits the linear convergence. Hence, we can derive the number of iterations, required to obtain ε -exact solution.

Corollary 1. *To achieve x^T with $f(x^T) - f_* \leq \varepsilon$ with high probability, Algorithm 1 needs*

$$T = T_0 \left(1 - \frac{\gamma_0}{\gamma_{EF}} \right) + \frac{4}{\gamma_{EF}\mu} \log \frac{\bar{V}^0}{\varepsilon}$$

iterations.

It is worth noticing, that our proposed method, unlike other methods, which incorporated differential privacy, converges to exact solution, rather than its neighbourhood.

One of the most important contributions is analyzing the added noise per algorithm step. We derive, that at every iteration the variance of added noise is proportional to $\sim \left(1 - \frac{\gamma\mu}{2}\right)^t$. It is in fact decreasing, which is crucial for achieving convergence to the exact optimum.

The next highlighted point is the choice of privacy levels per iteration. Prior approaches employed noise with constant variance throughout training, which resulted in same privacy budget per iteration. In this method, privacy at iteration t is proportional to $\sim \left(1 - \frac{\gamma\mu}{4}\right)^t \left(1 - \frac{\gamma\mu}{2}\right)^{-t}$. It can be noted, that the multiplier is greater, than one, therefore, this approach enables stronger information protection during later stages of convergence when model updates contain more valuable information about the optimum, as opposed to early iterations where updates primarily reflect the less informative starting point.

The last, but not least, is the number of iterations. Numerous existing methods (Wang et al., 2017; Li et al., 2022) both under nonconvex and PL assumptions bound the number of conducted iterations T , depending on the parameter’s properties and privacy budget. We have no upper bound on T , but the lower one. This makes sense - we can’t achieve adequate convergence with large privacy after one iteration - the noise term will be too impactful, and not distributed over many iterations. With given privacy budget per iteration, we need at least $\Omega\left(\frac{C_1 \cdot C_2}{n}\right)$ iterations. However, if we vary privacy levels, we can select $T = \Omega\left(\log \frac{C_1 \cdot C_2}{n}\right)$. More details are present in the Appendix.

4.2 NOISE BEFORE COMPRESSION

One may argue, that Algorithm 1 is not efficient, since we do not reduce the amount of transmitted information, due to the nature of Gaussian multivariate random vector. Frankly speaking, at line 7 of Algorithm 1 we transmit the uncompressed vector. To overcome this drawback we utilize the post-processing property of DP.

Lemma 5 (Post-processing). *Suppose $\mathcal{M} : \mathcal{D} \rightarrow \mathcal{R}$ is (ε, δ) -DP (ρ -zCDP) and $h : \mathcal{R} \rightarrow \mathcal{R}'$ be an arbitrary mapping. Then, $h \circ \mathcal{M}$ is (ε, δ) -DP (ρ -zCDP).*

One cannot compute a function of the output of a private algorithm \mathcal{M} and make it less private. This helps us a lot, since it allows us to utilize compression operators more efficient and to quantize the output of the Gaussian Mechanism. Therefore, we introduce Algorithm DPd-EF21-2:

Algorithm 2 DPd-EF21-2

- 1: **Parameters:** starting point $x^0 \in \mathbb{R}^d$, $g_i^0 = 0 \in \mathbb{R}^d$, learning rate $\gamma > 0$.
 - 2: **for** $t = 0, 1, 2, \dots, T - 1$ **do**
 - 3: Send x^t to nodes
 - 4: **for** all nodes $i = 1, 2, \dots, n$ in parallel **do**
 - 5: Compress $\Delta_i^t = \mathcal{C}(\nabla f_i(x^t) - g_i^{t-1} + \mathcal{N}(0, \sigma_{i,t}^2 I_d))$
 - 6: Update local state $g_i^t = g_i^{t-1} + \mathcal{C}(\nabla f_i(x^t) - g_i^{t-1})$
 - 7: Send Δ_i^t to the server
 - 8: **end for**
 - 9: Server computes $g^t = g^{t-1} + \frac{1}{n} \sum_{i=1}^n \Delta_i^t$
 - 10: Update the model $x^{t+1} = x^t - \gamma_t g^t$
 - 11: **end for**
-

Reasoning stays the same, as well, as proof’s techniques. Similar to the previous methods we end up with following corollary, concerning convergence process

Corollary 2. Selecting γ_k similarly as in Theorem 1, Algorithm 2 requires

$$T = \mathcal{O} \left(T_0 + \frac{1}{\gamma_{EF}\mu} \log \frac{\bar{V}_0}{\varepsilon} \right)$$

iterations to achieve x^T with $f(x^T) - f_* \leq \varepsilon$ with high probability.

One may find, that asymptotically rates in both cases coincide. This justifies and prefers applying compression after adding noise, due to the significantly less consumed bandwidth.

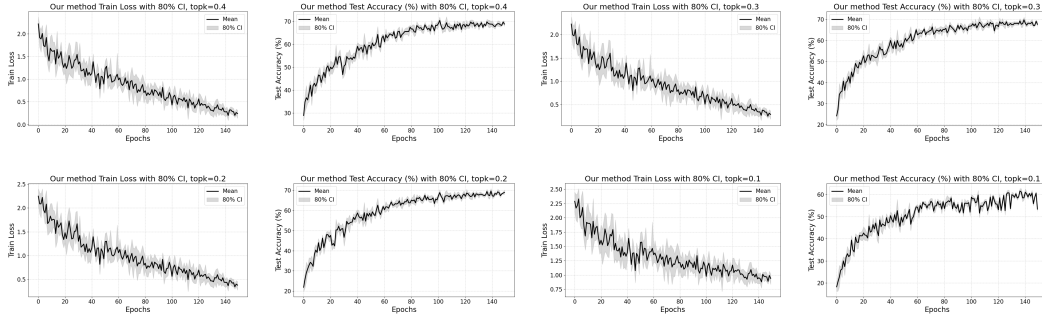


Figure 1: Top-k with varying k

5 NUMERICAL EXPERIMENTS

We evaluate our methods on the CIFAR-10 dataset to study the privacy–utility–communication trade-off. Unless otherwise noted, all methods use the same model and optimizer; full hyperparameters are provided in the appendix. For fairness, all curves are averaged over multiple seeds, and we report the mean with a shaded ± 1 std. band. We also report the realized privacy budget ε at a fixed $\delta = 10^{-5}$ using the zCDP accountant (Defs. 3–6; Lemmas 2–4), so comparisons are at *matched privacy*. We compare against DP-Clip21 and Clip21-SGD2M.

Noise schedules. We consider two diminishing-variance schedules: (i) a geometric schedule $\sigma_t^2 = (1 - \frac{\gamma\mu}{4})^t \sigma_0^2$; (ii) a theory-guided schedule with variance proportional to the Lyapunov term, $\sigma_t^2 = \kappa V_t$, where κ is chosen so that the per-round zCDP budget composes to the target (ε, δ) . Both schedules are compatible with our analysis, and we verify them empirically.

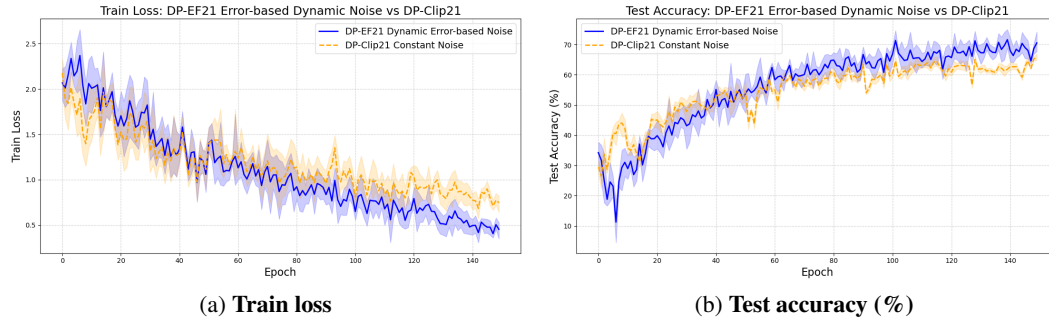


Figure 2: Our schedule with $\sigma_t^2 \propto V_t$ compared to DP-Clip21 at matched (ε, δ) . Diminishing noise attains comparable or better final accuracy and exhibits a lower loss plateau.

We observe that diminishing noise achieves at least the same peak accuracy as DP-Clip21 and yields a more favorable late-phase plateau in both loss and accuracy.

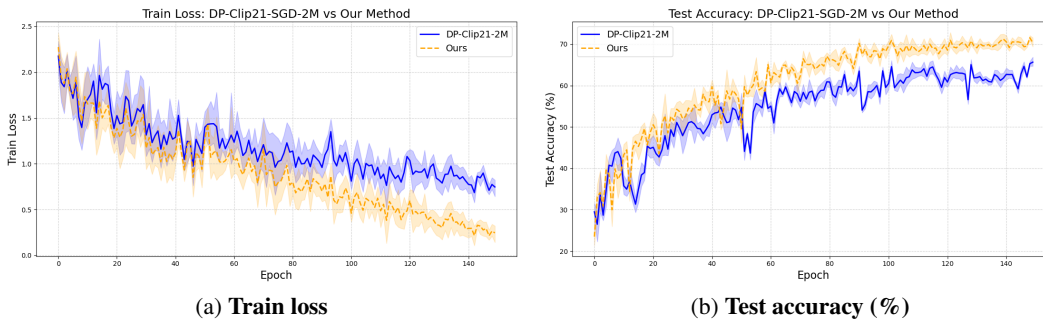


Figure 3: Comparison to Clip21-SGD2M (Islamov et al., 2025) at matched (ϵ, δ) . Our diminishing-noise variants are competitive while retaining communication efficiency.

Communication vs. accuracy. In practice, different compressors can be employed, and using fewer transmitted coordinates does not necessarily degrade accuracy. We vary the Top-K sparsification level with $k \in \{0.1d, 0.2d, 0.3d, 0.4d\}$ and also explore per-layer k . Figure 4 summarizes accuracy as a function of cumulative communicated coordinates (normalized).

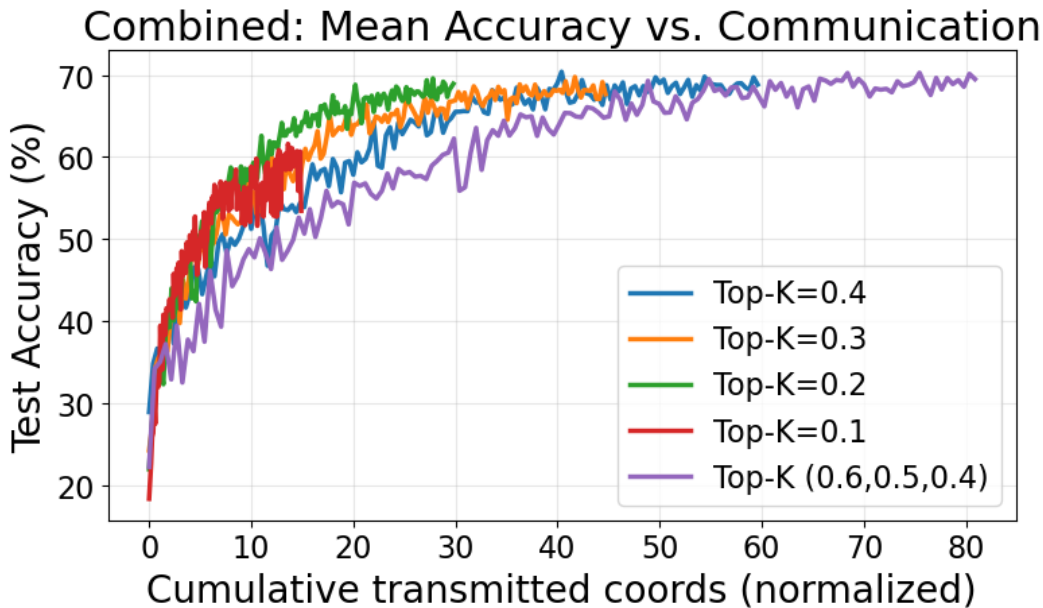


Figure 4: Accuracy vs. cumulative communication under different Top-K levels. Post-processing (noise then compression) preserves privacy while substantially reducing bandwidth.

Takeaways. Across matched privacy levels, diminishing-variance noise matches or exceeds the baselines’ accuracy while improving late-phase stability, and, when combined with compression after noising, delivers favorable accuracy–communication Pareto fronts.

REFERENCES

Martin Abadi, Andy Chu, Ian Goodfellow, H. Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security, CCS’16*. ACM, October 2016. doi: 10.1145/2976749.2978318. URL <http://dx.doi.org/10.1145/2976749.2978318>.

- 486 Abbas Acar, Hidayet Aksu, A Selcuk Uluogac, and Mauro Conti. A survey on homomorphic en-
487 crypton schemes: Theory and implementation. *ACM Computing Surveys (Csur)*, 51(4):1–35,
488 2018.
- 489 Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. Qsgd:
490 Communication-efficient sgd via gradient quantization and encoding, 2017.
- 491 Galen Andrew, Om Thakkar, Brendan McMahan, and Swaroop Ramaswamy. Differentially private
492 learning with adaptive clipping. *Advances in Neural Information Processing Systems*, 34:17455–
493 17466, 2021.
- 494 Galen Andrew, Om Thakkar, H. Brendan McMahan, and Swaroop Ramaswamy. Differentially
495 private learning with adaptive clipping, 2022.
- 496 Raef Bassily, Adam Smith, and Abhradeep Thakurta. Differentially private empirical risk minimiza-
497 tion: Efficient algorithms and tight error bounds, 2014.
- 498 Aleksandr Beznosikov, Samuel Horváth, Peter Richtárik, and Mher Safaryan. On biased compres-
499 sion for distributed learning, 2024.
- 500 Elizabeth A Brasher. Addressing the failure of anonymization: guidance from the european union’s
501 general data protection regulation. *Colum. Bus. L. Rev.*, pp. 209, 2018.
- 502 Zhiqi Bu, Yu-Xiang Wang, Sheng Zha, and George Karypis. Automatic clipping: Differentially pri-
503 vate deep learning made easier and stronger. *Advances in Neural Information Processing Systems*,
504 36:41727–41764, 2023.
- 505 Xiangyi Chen, Zhiwei Steven Wu, and Mingyi Hong. Understanding gradient clipping in private
506 sgd: A geometric perspective, 2021.
- 507 Laurent Condat, Kai Yi, and Peter Richtárik. Ef-bv: A unified theory of error feedback and variance
508 reduction mechanisms for biased and unbiased compression in distributed optimization, 2023.
- 509 Rudrajit Das, Satyen Kale, Zheng Xu, Tong Zhang, and Sujay Sanghavi. Beyond uniform lipschitz
510 condition in differentially private optimization, 2023.
- 511 John C Duchi, Michael I Jordan, and Martin J Wainwright. Privacy aware learning. *Journal of the*
512 *ACM (JACM)*, 61(6):1–57, 2014.
- 513 Cynthia Dwork. Differential privacy. In Michele Bugliesi, Bart Preneel, Vladimiro Sassone, and
514 Ingo Wegener (eds.), *Automata, Languages and Programming*, pp. 1–12, Berlin, Heidelberg,
515 2006. Springer Berlin Heidelberg. ISBN 978-3-540-35908-1.
- 516 Cynthia Dwork. Differential privacy: A survey of results. In Manindra Agrawal, Dingzhu Du,
517 Zhenhua Duan, and Angsheng Li (eds.), *Theory and Applications of Models of Computation*, pp.
518 1–19, Berlin, Heidelberg, 2008. Springer Berlin Heidelberg. ISBN 978-3-540-79228-4.
- 519 Ilyas Fatkhullin, Igor Sokolov, Eduard Gorbunov, Zhize Li, and Peter Richtárik. Ef21 with
520 bells & whistles: Practical algorithmic extensions of modern error feedback. *arXiv preprint*
521 *arXiv:2110.03294*, 2021.
- 522 Ilyas Fatkhullin, Alexander Tyurin, and Peter Richtárik. Momentum provably improves error feed-
523 back!, 2023.
- 524 Quan Geng and Pramod Viswanath. The optimal noise-adding mechanism in differential privacy.
525 *IEEE Transactions on Information Theory*, 62(2):925–951, 2015.
- 526 Eduard Gorbunov, Konstantin P Burlachenko, Zhize Li, and Peter Richtárik. Marina: Faster non-
527 convex distributed learning with compression. In *International Conference on Machine Learning*,
528 pp. 3788–3798. PMLR, 2021.
- 529 Yutong He, Xinmeng Huang, and Kun Yuan. Unbiased compression saves communication in dis-
530 tributed optimization: When and how much? *Advances in Neural Information Processing Sys-
531 tems*, 36:47991–48020, 2023.

- 540 Rustem Islamov, Samuel Horvath, Aurelien Lucchi, Peter Richtarik, and Eduard Gorbunov. Double
541 momentum and error feedback for clipping with fast rates and differential privacy. *arXiv preprint*
542 *arXiv:2502.11682*, 2025.
- 543
544 Peter Kairouz, H. Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin
545 Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, Rafael G. L.
546 D’Oliveira, Hubert Eichner, Salim El Rouayheb, David Evans, Josh Gardner, Zachary Garrett,
547 Adrià Gascón, Badih Ghazi, Phillip B. Gibbons, Marco Gruteser, Zaid Harchaoui, Chaoyang He,
548 Lie He, Zhouyuan Huo, Ben Hutchinson, Justin Hsu, Martin Jaggi, Tara Javidi, Gauri Joshi,
549 Mikhail Khodak, Jakub Konečný, Aleksandra Korolova, Farinaz Koushanfar, Sanmi Koyejo,
550 Tancrède Lepoint, Yang Liu, Prateek Mittal, Mehryar Mohri, Richard Nock, Ayfer Özgür, Rasmus
551 Pagh, Mariana Raykova, Hang Qi, Daniel Ramage, Ramesh Raskar, Dawn Song, Weikang Song,
552 Sebastian U. Stich, Ziteng Sun, Ananda Theertha Suresh, Florian Tramèr, Praneeth Vepakomma,
553 Jianyu Wang, Li Xiong, Zheng Xu, Qiang Yang, Felix X. Yu, Han Yu, and Sen Zhao. *Advances*
554 *and open problems in federated learning*, 2021.
- 555 Sai Praneeth Karimireddy, Quentin Rebjock, Sebastian U. Stich, and Martin Jaggi. Error feedback
556 fixes signsgd and other gradient compression schemes, 2019.
- 557 Sarit Khirirat, Hamid Reza Feyzmahdavian, and Mikael Johansson. Distributed learning with com-
558 pressed gradients, 2018.
- 559 Sarit Khirirat, Eduard Gorbunov, Samuel Horváth, Rustem Islamov, Fakhri Karray, and Peter
560 Richtárik. Clip21: Error feedback for gradient clipping, 2023.
- 561 Anastasia Koloskova, Hadrien Hendrikx, and Sebastian U. Stich. Revisiting gradient clipping:
562 Stochastic bias and tight convergence guarantees, 2023.
- 563
564 Zhize Li, Haoyu Zhao, Boyue Li, and Yuejie Chi. Soteriafl: A unified framework for private feder-
565 ated learning with communication compression, 2022.
- 566
567 Chaoyue Liu, Libin Zhu, and Mikhail Belkin. Loss landscapes and optimization in over-
568 parameterized non-linear systems and neural networks. *Applied and Computational Harmonic*
569 *Analysis*, 59:85–116, 2022.
- 570 Abdul Majeed and Sungchang Lee. Anonymization techniques for privacy preserving data publish-
571 ing: A comprehensive survey. *IEEE access*, 9:8512–8545, 2020.
- 572
573 H. Brendan McMahan, Daniel Ramage, Kunal Talwar, and Li Zhang. Learning differentially private
574 recurrent language models, 2018.
- 575 Konstantin Mishchenko, Eduard Gorbunov, Martin Takác, and Peter Richtárik. Distributed learning
576 with compressed gradient differences. 2019.
- 577
578 Yurii Nesterov. *Introductory lectures on convex optimization: A basic course*, volume 87. Springer
579 Science & Business Media, 2013.
- 580 Venkatadheeraj Pichapati, Ananda Theertha Suresh, Felix X Yu, Sashank J Reddi, and Sanjiv Kumar.
581 Adaclip: Adaptive clipping for private sgd. *arXiv preprint arXiv:1908.07643*, 2019.
- 582
583 Peter Richtárik, Igor Sokolov, and Ilyas Fatkhullin. Ef21: A new, simpler, theoretically better, and
584 practically faster error feedback, 2021.
- 585
586 Mark Schmidt, Nicolas Le Roux, and Francis Bach. Minimizing finite sums with the stochastic
587 average gradient. *Mathematical Programming*, 162(1):83–112, 2017.
- 588
589 Xicong Shen, Ying Liu, Fu Li, and Chunguang Li. Privacy-preserving federated learning against
590 label-flipping attacks on non-iid data. *IEEE Internet of Things Journal*, 11(1):1241–1255, 2023.
- 591
592 Sebastian U Stich. Local sgd converges fast and communicates little. *arXiv preprint*
593 *arXiv:1805.09767*, 2018.
- 594
595 Sebastian U Stich, Jean-Baptiste Cordonnier, and Martin Jaggi. Sparsified sgd with memory. *Ad-*
596 *vances in neural information processing systems*, 31, 2018.

594 Alexander Tyurin and Peter Richtárik. Dasha: Distributed nonconvex optimization with commu-
595 nication compression, optimal oracle complexity, and no client synchronization. *arXiv preprint*
596 *arXiv:2202.01268*, 2022.

597 Pablo Villalobos, Jaime Sevilla, Tamay Besiroglu, Lennart Heim, Anson Ho, and Marius Hobbhahn.
598 Machine learning model sizes and the parameter gap. *arXiv preprint arXiv:2207.02852*, 2022.

600 Di Wang, Minwei Ye, and Jinhui Xu. Differentially private empirical risk minimization revisited:
601 Faster and more general. *Advances in Neural Information Processing Systems*, 30, 2017.

602 Di Wang, Minwei Ye, and Jinhui Xu. Differentially private empirical risk minimization revisited:
603 Faster and more general, 2018.

605 Di Wang, Hanshen Xiao, Srini Devadas, and Jinhui Xu. On differentially private stochastic convex
606 optimization with heavy-tailed data, 2020.

607 Jingzhao Zhang, Tianxing He, Suvrit Sra, and Ali Jadbabaie. Why gradient clipping accelerates
608 training: A theoretical justification for adaptivity. *arXiv preprint arXiv:1905.11881*, 2019.

610 Ligeng Zhu, Zhijian Liu, and Song Han. Deep leakage from gradients, 2019.

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

A DP STATEMENTS

Lemma 6 (GM for Concentrated DP). *Let $f : \mathcal{X}^n \rightarrow \mathbb{R}^d$ has sensitivity Δ . Define a randomized algorithm $\mathcal{M} : \mathcal{X}^n \rightarrow \mathbb{R}^d$ by $\mathcal{M}(x) = \mathcal{N}\left(f(x), \frac{\Delta^2}{2\rho} I_d\right)$. Then \mathcal{M} is ρ -zCDP.*

Lemma 7 (Composition for Concentrated DP). *Let $M_1, M_2, \dots, M_k : \mathcal{X}^n \rightarrow \mathcal{R}$ be randomized algorithms. Suppose, M_i is ρ_i -zCDP for each i . Define $\mathcal{M}(x) : \mathcal{X}^n \rightarrow \mathcal{R}^k$ by $\mathcal{M}(x) \stackrel{\text{def}}{=} (M_1(x), M_2(x), \dots, M_k(x))$, where each algorithm is run independently. Then, \mathcal{M} is ρ -zCDP for $\rho = \sum_{i=1}^k \rho_i$.*

Lemma 8 (Conversion from Concentrated DP to DP). *ρ -zCDP implies $(\varepsilon = \rho + 2\sqrt{\rho \log 1/\delta}, \delta)$ -DP for all $\delta > 0$. Also, to obtain a given target (ε, δ) -DP, it suffices to have ρ -zCDP with $\rho = \frac{\varepsilon^2}{4 \log 1/\delta}$.*

B DESCENT LEMMA

Lemma 9 (Compressor). *If \mathcal{C} is a biased compressor, then there is a following bound on g_i^t , generated by Algorithm 1*

$$\frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^{t+1}) - g_i^t\|^2 \leq (1 + 2/\alpha)L^2 \|x^{t+1} - x^t\|^2 + \left(1 - \frac{1}{2/\alpha}\right) \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^t) - g_i^{t-1}\|^2$$

Proof. Using the Young's inequality, we obtain

$$\begin{aligned} \|\nabla f_i(x^{t+1}) - g_i^t\|^2 &\leq (1+s) \|\nabla f_i(x^{t+1}) - \nabla f_i(x^t)\|^2 + (1+s^{-1}) \|\nabla f_i(x^t) - g_i^t\|^2 \\ &\leq (1+s)L_i^2 \|x^{t+1} - x^t\|^2 \\ &\quad + (1+s^{-1}) \|\nabla f_i(x^t) - \mathcal{C}(\nabla f_i(x^t) - g_i^{t-1}) + g_i^{t-1}\|^2 \\ &\leq (1+\alpha)L_i^2 \|x^{t+1} - x^t\|^2 + (1+s^{-1})(1-\alpha) \|\nabla f_i(x^t) - g_i^{t-1}\|^2, \end{aligned}$$

which holds $\forall s > 0$. Take $s = 2/\alpha$, hence

$$\begin{aligned} \|\nabla f_i(x^{t+1}) - g_i^t\|^2 &\leq (1+2/\alpha)L_i^2 \|x^{t+1} - x^t\|^2 + \left(1 + \frac{1}{2/\alpha}\right) (1-\alpha) \|\nabla f_i(x^t) - g_i^{t-1}\|^2 \\ &\leq (1+2/\alpha)L_i^2 \|x^{t+1} - x^t\|^2 + \left(1 - \frac{1}{2/\alpha}\right) \|\nabla f_i(x^t) - g_i^{t-1}\|^2. \end{aligned}$$

For the mean deviations we get

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^{t+1}) - g_i^t\|^2 &\leq (1+2/\alpha)L^2 \|x^{t+1} - x^t\|^2 \\ &\quad + \left(1 - \frac{1}{2/\alpha}\right) \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^t) - g_i^{t-1}\|^2 \end{aligned}$$

□

Lemma 10 (Descent). *Define $V^{t+1} = f(x^{t+1}) - f_* + \theta\gamma_t \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^{t+1}) - g_i^t\|^2$, and $w > 0$.*

Suppose, that $\gamma_{t+1} \leq c_{up}\gamma_t$, where $c_{up} \leq \frac{2M-2\alpha}{M(2-\alpha)}$ for some $M > \alpha$. Then, with

$$\gamma_t \leq \gamma_{EF} = \min \left\{ \frac{1}{L \left(2 + \sqrt{2\theta c_{up} \left(2 + \frac{4}{\alpha}\right)}\right)}; \frac{2\alpha}{M\mu} \right\}$$

and $\gamma_t \geq 2^{-1/w}$ Algorithm 1 iterations obtain

$$\begin{aligned} V^{t+1} &\leq \left(1 - \frac{\gamma_t \mu}{2}\right) V^t + n(\gamma_t) \xi_t \\ &= \left(1 - \frac{\gamma_t \mu}{2}\right) V^t + (a\gamma^{1+w} + b\gamma^2 + c\gamma^3) \xi_t, \end{aligned}$$

where $\xi_t \sim \sigma_t^2 \cdot \chi^2(d)$, and $\chi^2(d)$ is a chi-squared random variable with d degrees of freedom, and $a = 1, b = L, c = c_{up}\theta L^2 \left(2 + \frac{4}{\alpha}\right), \theta = \frac{M}{2M - M\alpha - 2Mc_{up} + Mc_{up}}$.

756 *Proof.* From L -smoothness and Young's inequality we obtain:

$$\begin{aligned}
757 & f(x^{t+1}) - f_* \leq f(x^t) - f_* + \langle \nabla f(x^t), x^{t+1} - x^t \rangle + \frac{L}{2} \|x^{t+1} - x^t\|^2 \\
758 & = f(x^t) - f_* - \gamma_t \langle \nabla f(x^t), g_{EF}^t + g_N^t \rangle + \frac{L\gamma_t^2}{2} \|g_{EF}^t + g_N^t\|^2 \\
759 & \leq f(x^t) - f_* - \gamma_t \langle \nabla f(x^t), g_{EF}^t + g_N^t \rangle + L\gamma_t^2 \|g_{EF}^t\|^2 + L\gamma_t^2 \|g_N^t\|^2 \\
760 & = f(x^t) - f_* - \gamma_t \langle \nabla f(x^t), g_{EF}^t \rangle - \gamma_t \langle \nabla f(x^t), g_N^t \rangle \\
761 & + L_t\gamma^2 \|g_{EF}^t\|^2 + L\gamma_t^2 \|g_N^t\|^2 \\
762 & = f(x^t) - f_* - \frac{\gamma_t}{2} \|\nabla f(x^t)\|^2 + \frac{\gamma_t}{2} \|\nabla f(x^t) - g_{EF}^t\|^2 \\
763 & - \gamma_t \langle \nabla f(x^t), g_N^t \rangle + L\gamma_t^2 \|g_{EF}^t\|^2 + \left(L\gamma_t^2 - \frac{\gamma_t}{2}\right) \|g_N^t\|^2 \\
764 & \leq f(x^t) - f_* - \frac{\gamma_t}{2} \|\nabla f(x^t)\|^2 + \frac{\gamma_t}{2} \|\nabla f(x^t) - g_{EF}^t\|^2 \\
765 & + \frac{\gamma_t}{4} \|\nabla f(x^t)\|^2 + \gamma_t \|g_N^t\|^2 + \left(L\gamma_t^2 - \frac{\gamma_t}{2}\right) \|g_{EF}^t\|^2 + L\gamma_t^2 \|g_{EF}^t\|^2 \\
766 & \leq \left(1 - \frac{\gamma_t\mu}{2}\right) (f(x^t) - f_*) + \frac{\gamma_t}{2} \|\nabla f(x^t) - g_{EF}^t\|^2 + \left(L\gamma_t^2 - \frac{\gamma_t}{2}\right) \|g_{EF}^t\|^2 \\
767 & + (\gamma_t + L\gamma_t^2) \|g_N^t\|^2.
\end{aligned}$$

777 From Lemma 9 we have

$$\begin{aligned}
778 & \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^{t+1}) - g_i^t\|^2 \leq (1 + 2/\alpha)L^2 \|x^{t+1} - x^t\|^2 \\
779 & + \left(1 - \frac{1}{2/\alpha}\right) \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^t) - g_i^{t-1}\|^2 \\
780 & = (1 + 2/\alpha)\gamma_t^2 L^2 \|g_{EF}^t + g_N^t\|^2 \\
781 & + \left(1 - \frac{1}{2/\alpha}\right) \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^t) - g_i^{t-1}\|^2 \\
782 & \leq (2 + 4/\alpha)\gamma_t^2 L^2 \|g_{EF}^t\|^2 + (2 + 4/\alpha)\gamma^2 L^2 \|g_N^t\|^2 \\
783 & + \left(1 - \frac{1}{2/\alpha}\right) \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^t) - g_i^{t-1}\|^2
\end{aligned}$$

792 Since $g_{EF}^t = \frac{1}{n} \sum_{i=1}^n g_i^t$, $\nabla f(x^t) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(x^t)$ from convexity we achieve

$$\begin{aligned}
793 & \|\nabla f(x^t) - g_{EF}^{t-1}\|^2 = \left\| \frac{1}{n} \sum_{i=1}^n \nabla f_i(x^t) - g_{EF}^{t-1} \right\|^2 \leq \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^t) - g_i^{t-1}\|^2. \\
794 & \\
795 & \\
796 & \\
797 & \\
798 & \\
799 & \\
800 & \\
801 & \\
802 & \\
803 & \\
804 & \\
805 & \\
806 & \\
807 & \\
808 & \\
809 &
\end{aligned}$$

810 Define $V^{t+1} = f(x^{t+1}) - f_* + \theta\gamma_t \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^{t+1}) - g_i^t\|^2$, where $\theta > 0$ and will be declared
811 further. Therefore, we have

$$812 \begin{aligned} 813 V^{t+1} &\leq \left(1 - \frac{\gamma_t \mu}{2}\right) (f(x^t) - f_*) + \frac{\gamma_t}{2} \|\nabla f(x^t) - g_{EF}^t\|^2 + \left(L\gamma_t^2 - \frac{\gamma_t}{2}\right) \|g_{EF}^t\|^2 \\ 814 &+ (\gamma_t + L\gamma_t^2) \|g_N^t\|^2 + (2 + 4/\alpha)\theta\gamma_t^2\gamma_{t+1}L^2 \|g_{EF}^t\|^2 + (2 + 4/\alpha)\theta\gamma_t^2\gamma_{t+1}L^2 \|g_N^t\|^2 \\ 815 &+ \left(1 - \frac{1}{2/\alpha}\right) \theta\gamma_{t+1} \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^t) - g_i^{t-1}\|^2 \\ 816 &\leq \left(1 - \frac{\gamma_t \mu}{2}\right) (f(x^t) - f_*) + \left(\frac{\gamma_t}{2} + \theta\gamma_{t+1} \left(1 - \frac{1}{2/\alpha}\right)\right) \sum_{i=1}^n \|\nabla f_i(x^t) - g_i^{t-1}\|^2 \\ 817 &+ \left(L\gamma_t^2 - \frac{\gamma_t}{2} + (2 + 4/\alpha)\theta\gamma_t^2\gamma_{t+1}L^2\right) \|g_{EF}^t\|^2 + (\gamma_t + \gamma_t^2(L + L^2(2 + 4/\alpha)\theta\gamma_{t+1})) \|g_N^t\|^2. \end{aligned}$$

818 With $\theta = \frac{M}{2M - M\alpha - 2Mc_{up} + Mc_{up}}$

$$819 \begin{aligned} 820 V^{t+1} &\leq \left(1 - \frac{\gamma_t \mu}{2}\right) (f(x^t) - f_*) + \theta\gamma_t \left(1 - \frac{\alpha}{M}\right) \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^t) - g_i^{t-1}\|^2 \\ 821 &+ \left(L\gamma_t^2 - \frac{\gamma_t}{2} + (2 + 4/\alpha)\theta c_{up}\gamma_t^3 L^2\right) \|g_{EF}^t\|^2 + (\gamma_t + \gamma_t^2(L + L^2(2 + 4/\alpha)\theta c_{up}\gamma_t)) \|g_N^t\|^2. \end{aligned}$$

822 For γ_{EF} , defined as following:

$$823 \gamma_{EF} = \min \left\{ \frac{1}{L \left(2 + \sqrt{2c_{up}\theta \left(2 + \frac{4}{\alpha}\right)}\right)}; \frac{2\alpha}{M\mu} \right\},$$

824 for every $\gamma_t \leq \gamma_{EF}$ we have

$$825 V^{t+1} \leq \left(1 - \frac{\gamma_t \mu}{2}\right) V^t + (\gamma_t + \gamma_t^2(L + L^2(2 + 4/\alpha)\theta c_{up}\gamma_t)) \|g_N^t\|^2.$$

826 as $g_N \sim \mathcal{N}(0, \sigma_t^2 \cdot I)$, we obtain the needed. \square

827 **Corollary 3.** With $M = 4$ we get $\theta = \frac{4}{8 - 2\alpha - 4c_{up}\alpha}$, $c_{up} \leq \frac{4 - \alpha}{4 - 2\alpha}$,

$$828 \gamma_{EF} = \min \left\{ \frac{1}{2L \left(1 + \sqrt{\frac{2c_{up}}{8 - 2\alpha - 4c_{up}\alpha} \left(2 + \frac{4}{\alpha}\right)}\right)}; \frac{\alpha}{2\mu} \right\}$$

C MAIN THEOREM

Theorem 2. Suppose, that V^t is a Lyapunov function, defined by

$$V^t = f(x^t) - f_* + \frac{\theta\gamma_t}{n} \sum_{i=1}^n \|\nabla f_i(x^t) - g^{t-1}\|^2,$$

and let $\sigma_{i,t}^2$ be a variance of DP noise, injected on device i . Then, with proper choice of $\gamma_t \geq \gamma_{t-1}$ we might achieve linear convergence for V^t on any dataset:

$$V^t(\mathcal{D}) \leq 2 \prod_{k=0}^{t-1} \left(1 - \frac{\gamma_k \mu}{4}\right) \cdot \bar{V}^0 \quad \forall \mathcal{D},$$

where $\bar{V}^0 = \max_{\mathcal{D}'} V^0(\mathcal{D}')$.

Proof. We will proof this by induction. Base is true:

$$V^0(\mathcal{D}) \leq 2 \max_{\mathcal{D}'} V^0(\mathcal{D}')$$

Then we prove bound on sensitivity, We have

$$\begin{aligned} \Delta_{i,t}^2 &= \max_{\mathcal{D}_i \sim \mathcal{D}'_i} \left\| \mathcal{C}(\nabla f_i(x^t) - g_i^{t-1}) - \mathcal{C}(\nabla \hat{f}_i(x^t) - \hat{g}_i^{t-1}) \right\|^2 \\ &\leq 2 \left\| \mathcal{C}(\nabla f_i(x^t) - g_i^{t-1}) \right\|^2 + 2 \max_{\mathcal{D}_i \sim \mathcal{D}'_i} \left\| \mathcal{C}(\nabla \hat{f}_i(x^t) - \hat{g}_i^{t-1}) \right\|^2 \\ &\leq 2(1 + \alpha) \left\| \mathcal{C}(\nabla f_i(x^t) - g_i^{t-1}) - (\nabla f_i(x^t) - g_i^{t-1}) \right\|^2 \\ &\quad + 2(1 + \alpha^{-1}) \left\| \nabla f_i(x^t) - g_i^{t-1} \right\|^2 \\ &\quad + \max_{\mathcal{D}_i \sim \mathcal{D}'_i} \left[2(1 + \beta) \left\| \mathcal{C}(\nabla \hat{f}_i(x^t) - \hat{g}_i^{t-1}) - (\nabla \hat{f}_i(x^t) - \hat{g}_i^{t-1}) \right\|^2 \right. \\ &\quad \left. + 2(1 + \beta^{-1}) \left\| \nabla \hat{f}_i(x^t) - \hat{g}_i^{t-1} \right\|^2 \right] \\ &\leq 2((1 + \alpha)(1 - \alpha) + (1 + \alpha^{-1})) \left\| \nabla f_i(x^t) - g_i^{t-1} \right\|^2 \\ &\quad + 2((1 + \beta)(1 - \alpha) + (1 + \beta^{-1})) \max_{\mathcal{D}_i \sim \mathcal{D}'_i} \left\| \nabla \hat{f}_i(x^t) - \hat{g}_i^{t-1} \right\|^2 \end{aligned}$$

Take $\alpha = \beta = (1 - \alpha)^{-1/2}$, then

$$\Delta_{i,t}^2 \leq 2(1 + \sqrt{1 - \alpha})^2 \left(\left\| \nabla f_i(x^t) - g_i^{t-1} \right\|^2 + \max_{\mathcal{D}_i \sim \mathcal{D}'_i} \left\| \nabla \hat{f}_i(x^t) - \hat{g}_i^{t-1} \right\|^2 \right).$$

One can notice, that $\frac{1}{n} \sum_{i=1}^n \left\| \nabla f_i(x^t) - g_i^{t-1} \right\|^2 \leq \frac{1}{\theta\gamma_t} V^t$. Then,

$$\frac{1}{n} \sum_{i=1}^n \Delta_{i,t}^2 \leq \frac{2}{\theta\gamma_t} (1 + \sqrt{1 - \alpha})^2 \left(V^t(\mathcal{D}) + \frac{1}{n} \sum_{i=1}^n \max_{\mathcal{D}_i \sim \mathcal{D}'_i} \left\| \nabla \hat{f}_i(x^t) - \hat{g}_i^{t-1} \right\|^2 \right).$$

As choice of datasets \mathcal{D}'_i depends the device i , bound is less strict: $\left\| \nabla \hat{f}_i(x^t) - \hat{g}_i^{t-1} \right\|^2 \leq \frac{n}{\theta} V^t(\mathcal{D}')$.

Therefore, applying the induction presumption we have

$$\frac{1}{n} \sum_{i=1}^n \Delta_{i,t}^2 \leq \frac{4}{\theta\gamma_t} (1 + \sqrt{1 - \alpha})^2 \left(V^t(\mathcal{D}) + \sum_{i=1}^n V^t(\mathcal{D}') \right) \leq \frac{A}{\gamma_t} \prod_{k=0}^{t-1} \left(1 - \frac{\gamma_k \mu}{4}\right) \cdot \bar{V}^0,$$

where $A = \frac{4(n+1)}{\theta} (1 + \sqrt{1 - \alpha})^2$. Then we unroll the inequalities from the descent lemma:

$$\begin{aligned} V^{t+1} &\leq \left(1 - \frac{\gamma_t \mu}{2}\right) V^t + n(\gamma_t) \xi_t \leq \prod_{k=0}^t \left(1 - \frac{\gamma_k \mu}{2}\right) \cdot V^0 + \sum_{k=0}^t n(\gamma_k) \xi_k \prod_{i=k+1}^t \left(1 - \frac{\gamma_i \mu}{2}\right) \\ &\leq \prod_{k=0}^t \left(1 - \frac{\gamma_k \mu}{2}\right) \cdot \bar{V}^0 + \sum_{k=0}^t n(\gamma_k) \xi_k \prod_{i=k+1}^t \left(1 - \frac{\gamma_i \mu}{2}\right). \end{aligned}$$

For ξ_t we have following concentration inequality:

$$\mathbb{P}[\xi_t \geq \sigma_t^2 d(1 + \varepsilon_t)] \leq \exp\left(-\frac{d}{4} \min(\varepsilon_t, \varepsilon_t^2)\right).$$

Define $\varepsilon_k = \max\{1, \beta(k+1)\}$, where β will be determined further. Therefore, always $\min\{\varepsilon_k, \varepsilon_k^2\} = \varepsilon_k$. Then,

$$\begin{aligned} \varepsilon_k &\geq \beta(k+1), \\ -\varepsilon_k &\leq -\beta(k+1), \\ \exp\left(-\frac{d\varepsilon_k}{4}\right) &\leq \exp\left(-\frac{d\beta(k+1)}{4}\right) \end{aligned}$$

and

$$\sum_{k=0}^t \exp\left(-\frac{d\varepsilon_k}{4}\right) \leq \sum_{k=0}^t \exp\left(-\frac{d\beta(k+1)}{4}\right) \leq \frac{\exp(-d\beta/4)}{1 - \exp(-d\beta/4)}.$$

To bound this with p we need to take

$$\beta = \frac{4}{d} \ln\left(1 + \frac{1}{p}\right)$$

Therefore, with probability at least $1 - p$ we have

$$V^{t+1} \leq \prod_{k=0}^t \left(1 - \frac{\gamma_k \mu}{2}\right) \cdot \bar{V}^0 + d \sum_{k=0}^t n(\gamma_k) \sigma_k^2 (1 + \varepsilon_k) \prod_{i=k+1}^t \left(1 - \frac{\gamma_i \mu}{2}\right)$$

Next, using this derived inequality we will derive the overall convergence. As shown earlier, if iteration k satisfies ρ_k -zCDP for device i , if added Gaussian noise have variance $\sigma_{i,k}^2 = \frac{\Delta_{i,k}^2}{2\rho_k}$. If every iteration is ρ_k -zCDP, then, adaptive composition will be $\sum_{k=0}^t \rho_k$ -zCDP. Then, $\rho_k \stackrel{\text{def}}{=} \frac{\nu_k}{Z} \frac{\varepsilon^2}{4 \log 1/\delta}$,

where $Z = \sum_{k=0}^t \nu_k$. Then, according to previous lemma, adaptive composition will be (ε, δ) -DP.

Note, that we analyzed $\sigma_{i,t}$ instead of σ_t . However, we have $\mathcal{N}(0, \sigma_t^2 I) = \frac{1}{n} \sum_{i=1}^n \mathcal{N}(0, \sigma_{i,t}^2 I)$. Hence,

$\sigma_t^2 = \frac{1}{n^2} \sum_{i=1}^n \sigma_{i,t}^2$. Therefore, $\sigma_t^2 = \frac{1}{2\rho_t n^2} \sum_{i=1}^n \Delta_{i,t}^2 \leq \frac{1}{2\rho_t} \frac{A}{n\gamma_t} \prod_{k=0}^t \left(1 - \frac{\gamma_k \mu}{4}\right) \bar{V}^0$. Therefore,

$$\begin{aligned} V^{t+1} &\leq \prod_{k=0}^t \left(1 - \frac{\gamma_k \mu}{2}\right) \cdot \bar{V}^0 + \frac{Ad\bar{V}^0}{2n} \sum_{k=0}^t \frac{n(\gamma_k)}{\gamma_k \rho_k} (1 + \varepsilon_k) \prod_{i=0}^{k-1} \left(1 - \frac{\gamma_i \mu}{4}\right) \cdot \prod_{i=k+1}^t \left(1 - \frac{\gamma_i \mu}{2}\right) \\ &\leq \prod_{k=0}^t \left(1 - \frac{\gamma_k \mu}{2}\right) \cdot \bar{V}^0 + \frac{Ad\bar{V}^0}{2n} \sum_{k=0}^t \frac{n(\gamma_k)}{\gamma_k \rho_k} (1 + \varepsilon_k) \prod_{i=0}^{k-1} \left(1 - \frac{\gamma_i \mu}{4}\right) \cdot \prod_{i=k}^{t-1} \left(1 - \frac{\gamma_i \mu}{2}\right) \\ &= \prod_{k=0}^t \left(1 - \frac{\gamma_k \mu}{2}\right) \cdot \bar{V}^0 + \frac{Ad\bar{V}^0}{2n} \sum_{k=0}^t \frac{n(\gamma_k)}{\gamma_k \rho_k} (1 + \varepsilon_k) \prod_{i=0}^{k-1} \left(1 - \frac{\gamma_i \mu}{4}\right) \left(1 - \frac{\gamma_i \mu}{2}\right)^{-1} \cdot \prod_{i=0}^{t-1} \left(1 - \frac{\gamma_i \mu}{2}\right) \\ &= \prod_{k=0}^t \left(1 - \frac{\gamma_k \mu}{2}\right) \cdot \bar{V}^0 \cdot \left[1 + \frac{Ad/(2n)}{1 - \gamma_t \mu/2} \sum_{k=0}^t \frac{n(\gamma_k)}{\gamma_k \rho_k} (1 + \varepsilon_k) \prod_{i=0}^{k-1} \left(\frac{1 - \gamma_i \mu/4}{1 - \gamma_i \mu/2}\right)\right]. \end{aligned}$$

Let $h(y) = \frac{1-y/4}{1-y/2}$. We have $h'(y) = \frac{1}{(y-2)^2}$, therefore, h is increasing. Define $\nu_k =$

$$\prod_{i=0}^{k-1} \left(\frac{1-\frac{\gamma_i\mu}{4}}{1-\frac{\gamma_i\mu}{2}} \right)^{-1} \cdot \left(\frac{1-\frac{\gamma_{EF}\mu}{2}}{1-\frac{\gamma_{EF}\mu}{4}} \right)^k. \text{ Then,}$$

$$Z = \sum_{k=0}^t \prod_{i=0}^{k-1} \left(\frac{1-\frac{\gamma_i\mu}{4}}{1-\frac{\gamma_i\mu}{2}} \right)^{-1} \cdot \left(\frac{1-\frac{\gamma_{EF}\mu}{2}}{1-\frac{\gamma_{EF}\mu}{4}} \right)^k \geq \sum_{k=0}^t \left(\frac{1-\frac{\gamma_{EF}\mu}{4}}{1-\frac{\gamma_{EF}\mu}{2}} \right)^{-k} \cdot \left(\frac{1-\frac{\gamma_{EF}\mu}{2}}{1-\frac{\gamma_{EF}\mu}{4}} \right)^k = t+1$$

and

$$\frac{1}{Z} \leq \frac{1}{t+1}$$

So, we finally obtain

$$V^{t+1} \leq \prod_{k=0}^t \left(1 - \frac{\gamma_k\mu}{2} \right) \cdot \bar{V}^0 \cdot \left[1 + \frac{Ad/(2n)}{1-\gamma_t\mu/2} \frac{1}{t+1} \frac{\epsilon^2}{4 \log 1/\delta} \sum_{k=0}^t \frac{n(\gamma_k)}{\gamma_k} \left(\frac{1-\frac{\gamma_{EF}\mu}{2}}{1-\frac{\gamma_{EF}\mu}{4}} \right)^k (1 + \epsilon_k) \right]$$

Use the fact, that $\frac{n(\gamma_k)}{\gamma_k} = \frac{a\gamma_k + b\gamma_k^2 + c\gamma_k^3}{\gamma_k} = a + b\gamma_k + c\gamma_k^2 = m(\gamma_k) \leq m(\gamma_t)$ for $k \leq t$. Then,

$$\begin{aligned} V^{t+1} &\leq \prod_{k=0}^t \left(1 - \frac{\gamma_k\mu}{2} \right) \cdot \bar{V}^0 \cdot \left[1 + \frac{Ad/(2n)}{1-\gamma_t\mu/2} \frac{1}{t+1} \frac{\epsilon^2}{4 \log 1/\delta} m(\gamma_t) \sum_{k=0}^t \left(\frac{1-\frac{\gamma_{EF}\mu}{2}}{1-\frac{\gamma_{EF}\mu}{4}} \right)^k (1 + \epsilon_k) \right] \\ &\leq \prod_{k=0}^t \left(1 - \frac{\gamma_k\mu}{2} \right) \cdot \bar{V}^0 \cdot \left[1 + \frac{Ad/(2n)}{1-\gamma_t\mu/2} \frac{1}{t+1} \frac{\epsilon^2}{4 \log 1/\delta} m(\gamma_t) \left(2 \sum_{k=0}^t \left(\frac{1-\frac{\gamma_{EF}\mu}{2}}{1-\frac{\gamma_{EF}\mu}{4}} \right)^k + \beta \sum_{k=0}^t k \left(\frac{1-\frac{\gamma_{EF}\mu}{2}}{1-\frac{\gamma_{EF}\mu}{4}} \right)^k \right) \right] \\ &\leq \prod_{k=0}^t \left(1 - \frac{\gamma_k\mu}{2} \right) \cdot \bar{V}^0 \cdot \left[1 + \frac{Ad/(2n)}{1-\gamma_t\mu/2} \frac{1}{t+1} \frac{\epsilon^2}{4 \log 1/\delta} m(\gamma_t) \left(2 \sum_{k=0}^{\infty} \left(\frac{1-\frac{\gamma_{EF}\mu}{2}}{1-\frac{\gamma_{EF}\mu}{4}} \right)^k + \beta \sum_{k=0}^{\infty} k \left(\frac{1-\frac{\gamma_{EF}\mu}{2}}{1-\frac{\gamma_{EF}\mu}{4}} \right)^k \right) \right] \\ &\leq \prod_{k=0}^t \left(1 - \frac{\gamma_k\mu}{2} \right) \cdot \bar{V}^0 \cdot \left[1 + \frac{Ad/(2n)}{1-\gamma_t\mu/2} \frac{1}{t+1} \frac{\epsilon^2}{4 \log 1/\delta} m(\gamma_t) \left(\frac{2}{1-\frac{1-\frac{\gamma_{EF}\mu}{2}}{1-\frac{\gamma_{EF}\mu}{4}}} + \frac{\beta \cdot \frac{1-\frac{\gamma_{EF}\mu}{2}}{1-\frac{\gamma_{EF}\mu}{4}}}{\left(1-\frac{1-\frac{\gamma_{EF}\mu}{2}}{1-\frac{\gamma_{EF}\mu}{4}} \right)^2} \right) \right] \\ &= \prod_{k=0}^t \left(1 - \frac{\gamma_k\mu}{2} \right) \cdot \bar{V}^0 \cdot \left[1 + \frac{Ad/(2n)}{1-\gamma_t\mu/2} \frac{1}{t+1} \frac{\epsilon^2}{4 \log 1/\delta} m(\gamma_t) \left(\frac{8(1-\frac{\gamma_{EF}\mu}{4})}{\gamma_{EF}\mu} + \frac{16\beta(1-\frac{\gamma_{EF}\mu}{2})(1-\frac{\gamma_{EF}\mu}{4})}{(\gamma_{EF}\mu)^2} \right) \right] \end{aligned}$$

Overall, we have following situation:

$$V^{t+1} \leq \prod_{k=0}^t \left(1 - \frac{\gamma_k\mu}{2} \right) \cdot \bar{V}^0 \cdot \text{func}(\gamma_t).$$

We want

$$V^{t+1} \leq 2 \prod_{k=0}^t \left(1 - \frac{\gamma_k\mu}{4} \right) \cdot \bar{V}^0,$$

therefore, we need

$$\text{func}(\gamma_t) \leq 2 \prod_{k=0}^t \left(\frac{1-\frac{\gamma_k\mu}{4}}{1-\frac{\gamma_k\mu}{2}} \right).$$

We have

$$\frac{1-\frac{\gamma_k\mu}{4}}{1-\frac{\gamma_k\mu}{2}} \geq \frac{1-\frac{\gamma_t\mu}{4}}{1-\frac{\gamma_t\mu}{2}},$$

therefore, it is sufficient to proof

$$\text{func}(\gamma_t) \leq 2 \left(\frac{1-\frac{\gamma_t\mu}{4}}{1-\frac{\gamma_t\mu}{2}} \right)^{t+1} = 2 \left(\frac{1-\frac{\gamma_t\mu}{2}}{1-\frac{\gamma_t\mu}{4}} \right)^{-t-1} = 2 \left(1 - \frac{\frac{\gamma_t\mu}{4}}{1-\frac{\gamma_t\mu}{4}} \right)^{-t-1},$$

or

$$\text{func}(\gamma_t) \left(1 - \frac{\frac{\gamma_t\mu}{4}}{1-\frac{\gamma_t\mu}{4}} \right)^{t+1} \leq 2.$$

1026 As we have $(1 - x)^k \leq e^{-kx}$, for $x < 1$ we obtain

1027
 1028
$$\text{func}(\gamma_t) \left(1 - \frac{\frac{\gamma_t \mu}{4}}{1 - \frac{\gamma_t \mu}{4}}\right)^{t+1} \leq \text{func}(\gamma_t) \exp\left(- (t+1) \frac{\gamma_t \mu}{4 - \gamma_t \mu}\right) \leq \text{func}(\gamma_t) \exp\left(- \frac{(t+1)\gamma_t \mu}{4}\right).$$

1029
 1030 Finally, it is sufficient to guarantee, that

1031
 1032
$$\frac{\text{func}(\gamma_t)}{\exp\left(\frac{(t+1)\gamma_t \mu}{4}\right)} \leq 2$$

1033
 1034 Putting all the constants together we obtain the needed. \square

1035
 1036 **Remark 1.** We can take non-decreasing γ_y , since if γ_t satisfies the constraints and iteration t , γ_t
 1037 will satisfy constraints at iteration $t + 1$ too. Due to the inevitable increase of the exponent, γ_{EF}
 1038 will be eligible at some point, therefore, we will end up with linear convergence.

1039
 1040
 1041
 1042
 1043
 1044
 1045
 1046
 1047
 1048
 1049
 1050
 1051
 1052
 1053
 1054
 1055
 1056
 1057
 1058
 1059
 1060
 1061
 1062
 1063
 1064
 1065
 1066
 1067
 1068
 1069
 1070
 1071
 1072
 1073
 1074
 1075
 1076
 1077
 1078
 1079

D VARYING PRIVACY LEVELS

One can notice, that careful choice of constant q will directly influence the privacy levels per iteration.

Define $\nu_k = \prod_{i=0}^{k-1} \left(\frac{1 - \frac{\gamma_i \mu}{4}}{1 - \frac{\gamma_i \mu}{2}} \right)^{-1} \cdot q^{-k}$, where $q \in \left(\frac{1 - \frac{\gamma_{EF} \mu}{4}}{1 - \frac{\gamma_{EF} \mu}{2}}; 1 \right)$. Then,

$$Z = \sum_{k=0}^t \prod_{i=0}^{k-1} \left(\frac{1 - \frac{\gamma_i \mu}{4}}{1 - \frac{\gamma_i \mu}{2}} \right)^{-1} \cdot q^{-k} \geq \sum_{k=0}^t \left(\frac{1 - \frac{\gamma_{EF} \mu}{4}}{1 - \frac{\gamma_{EF} \mu}{2}} \right)^{-k} \cdot q^{-k} = \frac{1 - \left(\frac{1 - \frac{\gamma_{EF} \mu}{4}}{1 - \frac{\gamma_{EF} \mu}{2}} \right)^{-t-1} q^{-t-1}}{1 - \left(\frac{1 - \frac{\gamma_{EF} \mu}{4}}{1 - \frac{\gamma_{EF} \mu}{2}} \right)^{-1} q^{-1}}$$

and

$$\frac{1}{Z} \leq \frac{1 - \left(\frac{1 - \frac{\gamma_{EF} \mu}{4}}{1 - \frac{\gamma_{EF} \mu}{2}} \right)^{-1} q^{-1}}{1 - \left(\frac{1 - \frac{\gamma_{EF} \mu}{4}}{1 - \frac{\gamma_{EF} \mu}{2}} \right)^{-t-1} q^{-t-1}}$$

So, we finally obtain

$$V^{t+1} \leq \prod_{k=0}^t \left(1 - \frac{\gamma_k \mu}{2} \right) \cdot \bar{V}^0 \cdot \left[1 + \frac{Ad/(2n)}{1 - \gamma_t \mu / 2} \frac{1 - \left(\frac{1 - \frac{\gamma_{EF} \mu}{4}}{1 - \frac{\gamma_{EF} \mu}{2}} \right)^{-1} q^{-1}}{1 - \left(\frac{1 - \frac{\gamma_{EF} \mu}{4}}{1 - \frac{\gamma_{EF} \mu}{2}} \right)^{-t-1} q^{-t-1}} \frac{\varepsilon^2}{4 \log 1/\delta} \sum_{k=0}^t \frac{n(\gamma_k)}{\gamma_k} q^k (1 + \varepsilon_k) \right]$$

Use the fact, that $\frac{n(\gamma_k)}{\gamma_k} = \frac{a\gamma_k + b\gamma_k^2 + c\gamma_k^3}{\gamma_k} = a + b\gamma_k + c\gamma_k^2 = m(\gamma_k) \leq m(\gamma_t)$ for $k \leq t$. Then,

$$\begin{aligned} V^{t+1} &\leq \prod_{k=0}^t \left(1 - \frac{\gamma_k \mu}{2} \right) \cdot \bar{V}^0 \cdot \left[1 + \frac{Ad/(2n)}{1 - \gamma_t \mu / 2} \frac{1 - \left(\frac{1 - \frac{\gamma_{EF} \mu}{4}}{1 - \frac{\gamma_{EF} \mu}{2}} \right)^{-1} q^{-1}}{1 - \left(\frac{1 - \frac{\gamma_{EF} \mu}{4}}{1 - \frac{\gamma_{EF} \mu}{2}} \right)^{-t-1} q^{-t-1}} \frac{\varepsilon^2}{4 \log 1/\delta} m(\gamma_t) \sum_{k=0}^t q^k (1 + \varepsilon_k) \right] \\ &\leq \prod_{k=0}^t \left(1 - \frac{\gamma_k \mu}{2} \right) \cdot \bar{V}^0 \cdot \left[1 + \frac{Ad/(2n)}{1 - \gamma_t \mu / 2} \frac{1 - \left(\frac{1 - \frac{\gamma_{EF} \mu}{4}}{1 - \frac{\gamma_{EF} \mu}{2}} \right)^{-1} q^{-1}}{1 - \left(\frac{1 - \frac{\gamma_{EF} \mu}{4}}{1 - \frac{\gamma_{EF} \mu}{2}} \right)^{-t-1} q^{-t-1}} \frac{\varepsilon^2}{4 \log 1/\delta} m(\gamma_t) \left(2 \sum_{k=0}^t q^k + \beta \sum_{k=0}^t k q^k \right) \right] \\ &\leq \prod_{k=0}^t \left(1 - \frac{\gamma_k \mu}{2} \right) \cdot \bar{V}^0 \cdot \left[1 + \frac{Ad/(2n)}{1 - \gamma_t \mu / 2} \frac{1 - \left(\frac{1 - \frac{\gamma_{EF} \mu}{4}}{1 - \frac{\gamma_{EF} \mu}{2}} \right)^{-1} q^{-1}}{1 - \left(\frac{1 - \frac{\gamma_{EF} \mu}{4}}{1 - \frac{\gamma_{EF} \mu}{2}} \right)^{-t-1} q^{-t-1}} \frac{\varepsilon^2}{4 \log 1/\delta} m(\gamma_t) \left(2 \sum_{k=0}^{\infty} q^k + \beta \sum_{k=0}^{\infty} k q^k \right) \right] \\ &\leq \prod_{k=0}^t \left(1 - \frac{\gamma_k \mu}{2} \right) \cdot \bar{V}^0 \cdot \left[1 + \frac{Ad/(2n)}{1 - \gamma_t \mu / 2} \frac{1 - \left(\frac{1 - \frac{\gamma_{EF} \mu}{4}}{1 - \frac{\gamma_{EF} \mu}{2}} \right)^{-1} q^{-1}}{1 - \left(\frac{1 - \frac{\gamma_{EF} \mu}{4}}{1 - \frac{\gamma_{EF} \mu}{2}} \right)^{-t-1} q^{-t-1}} \frac{\varepsilon^2}{4 \log 1/\delta} m(\gamma_t) \left(\frac{2}{1-q} + \frac{\beta q}{(1-q)^2} \right) \right]. \end{aligned}$$

The rest of the proof is similar With various q we can influence the privacy budget per iteration.

E DESCENT LEMMA 2

Now we will derive the decent lemma for the Algorithm 2.

From L -smoothness we obtain

$$\begin{aligned} f(x^{t+1}) - f_* &\leq f(x^t) - f_* - \frac{\gamma_t}{2} \|\nabla f(x^t)\|^2 - \frac{\gamma_t}{2} \|\tilde{g}^t\|^2 + \frac{\gamma_t}{2} \|\nabla f(x^t) - \tilde{g}^t\|^2 \\ &\leq f(x^t) - f_* - \frac{\gamma_t}{2} \|\nabla f(x^t)\|^2 - \frac{\gamma_t}{2} \|\tilde{g}^t\|^2 + \gamma_t \|\nabla f(x^t) - g^t\|^2 + \gamma_t \|g^t - \tilde{g}^t\|^2. \end{aligned}$$

For the EF part the derivation stays the same:

$$\|\nabla f(x^{t+1}) - g^{t+1}\|^2 \leq \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^{t+1}) - g_i^{t+1}\|^2 \leq \left(1 - \frac{\alpha}{2}\right) \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^t) - g_i^t\|^2 + \left(1 + \frac{2}{\alpha}\right) \gamma_t^2 L^2 \|\tilde{g}^t\|^2.$$

Our interest is in the second term, where we can obtain following:

$$\begin{aligned} \gamma_t \|g^t - \tilde{g}^t\|^2 &\leq \frac{\gamma_t}{n} \sum_{i=1}^n \|\mathcal{C}(\nabla f_i(x^t) - g_i^{t-1}) - \mathcal{C}(\nabla f_i(x^t) - g_i^{t-1} + \eta_{i,t})\|^2 \\ &\leq \frac{2\gamma_t}{n} \sum_{i=1}^n \|\mathcal{C}(\nabla f_i(x^t) - g_i^{t-1})\|^2 + \|\mathcal{C}(\nabla f_i(x^t) - g_i^{t-1} + \eta_{i,t})\|^2 \\ &\leq \frac{2\gamma_t (1 + \sqrt{1 - \alpha})^2}{n} \sum_{i=1}^n \|\nabla f_i(x^t) - g_i^{t-1}\|^2 + \|\nabla f_i(x^t) - g_i^{t-1} + \eta_{i,t}\|^2 \\ &\leq \frac{\gamma_t (1 + \sqrt{1 - \alpha})^2}{n} \sum_{i=1}^n 4 \|\nabla f_i(x^t) - g_i^{t-1}\|^2 + 2 \|\eta_{i,t}\|^2 \end{aligned}$$

Similarly, we obtain the same asymptotical rates, as above

F ADDITIONAL NUMERICAL EXPERIMENTS

We compare Algorithm 2 to 1 and obtain, that results are similar, therefore we may significantly reduce the number of bits sent. All experiments are done on A100.

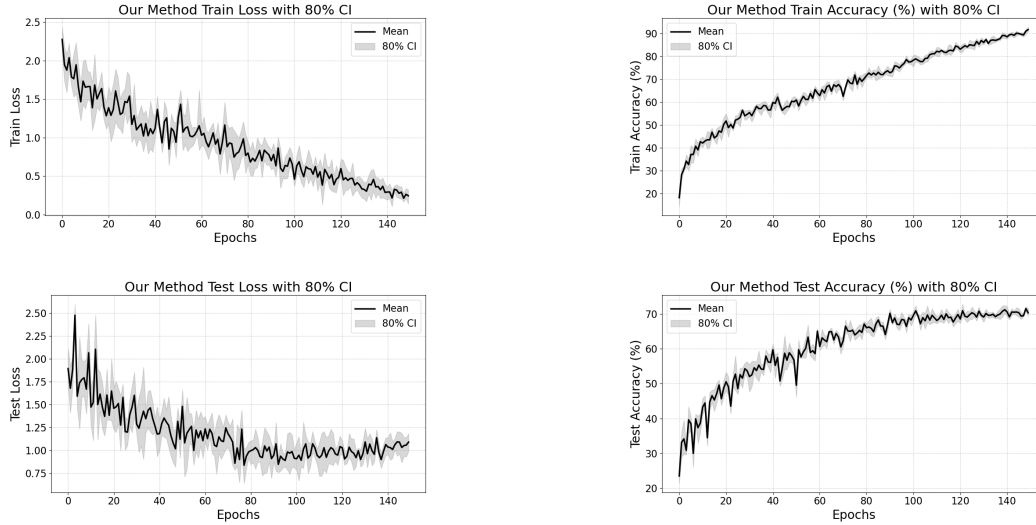


Figure 5: Compression before adding noise

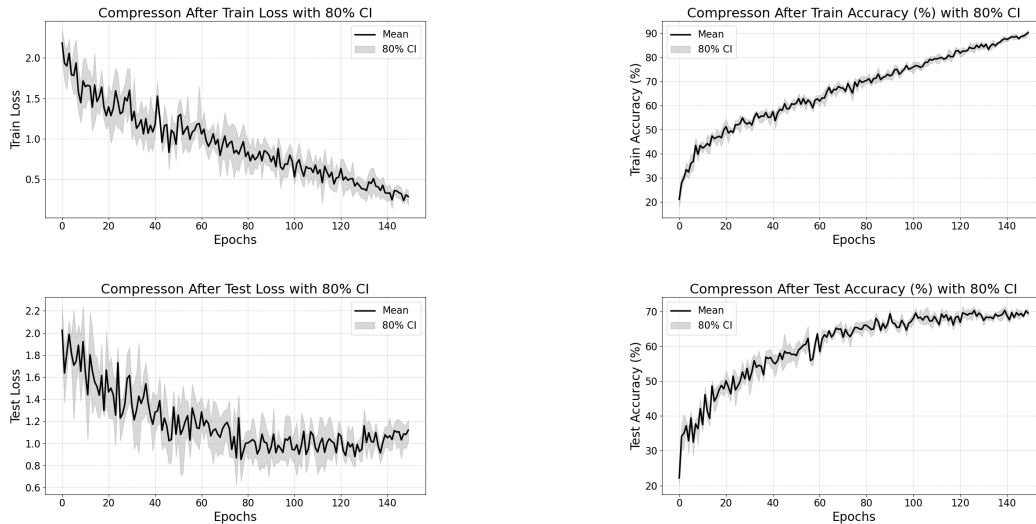


Figure 6: Compression after adding noise

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

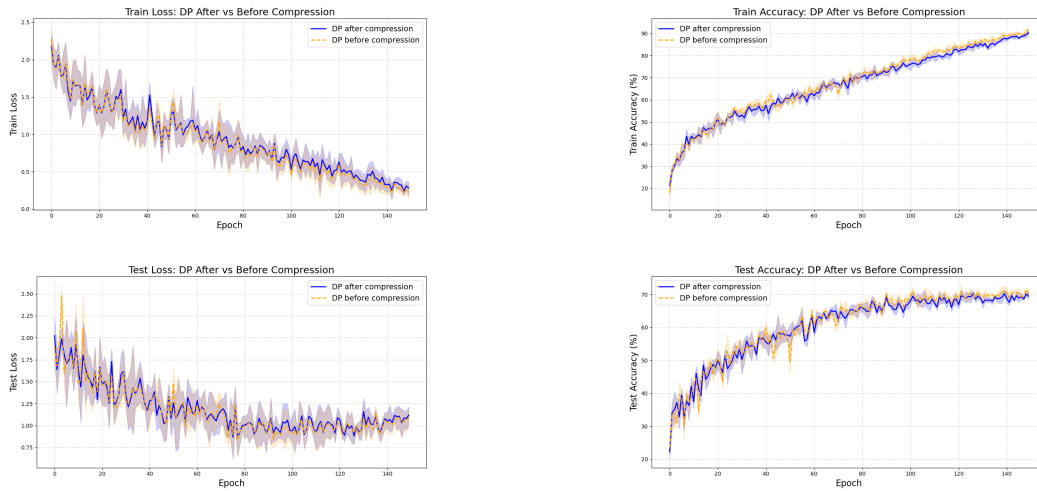


Figure 7: Comparison of compressing before and after adding noise

G DECLARATION OF LLM USAGE

We employed Large Language Models to improve the clarity and style of the text.