

MARKOVIAN COMPRESSION: LOOKING TO THE PAST HELPS ACCELERATE THE FUTURE

Anonymous authors

Paper under double-blind review

ABSTRACT

This paper deals with distributed optimization problems that use compressed communication to achieve efficient performance and mitigate the communication bottleneck. We propose a family of compression schemes in which operators transform vectors fed to their input according to a Markov chain, i.e., the stochasticity of the compressors depends on previous iterations. Intuitively, this should accelerate the convergence of optimization methods, as considering previous iterations seems more natural and robust. The compressors are implemented in the vanilla Quantized Stochastic Gradient Descent (QSGD) algorithm. To further improve efficiency and convergence rate, we apply the momentum acceleration method. We prove convergence results for our algorithms with Markovian compressors and show theoretically that the accelerated method converges faster than the basic version. The analysis covers non-convex, Polyak-Lojasiewicz (PL), and strongly convex cases. Experiments are conducted to demonstrate the applicability of the results to distributed data-parallel optimization problems. Practical results demonstrate the superiority of methods utilizing our compressors design over several existing optimization algorithms.

1 INTRODUCTION

The optimization problem is currently a key issue in many practical applications, such as optimization in neural network training, resource allocation in computational systems, and parameter tuning in algorithmic trading strategies.

In addition, a variety of algorithms for optimization on a single device, such as SGD Robbins & Monro (1951), Adam Kingma & Ba (2014), Lion Yazdani & Jolai (2016), have emerged and been subjected to theoretical analysis. However, in the contemporary landscape of deep learning, there is an increasing trend towards adopting intricate and expansive models that pose significant training challenges. Prominent among these challenges are advanced deep learning frameworks for image analysis, sophisticated natural language processing structures akin to transformers Vaswani et al. (2017), and complex reinforcement learning methodologies designed for autonomous system operations Kiran et al. (2021). As a result, the training of such models has become impractical for execution on a single device due to their requirement for extensive data sets for training, which are unfeasible to store on a single device. Consequently, optimization algorithms have been specifically developed for distributed training Verbraeken et al. (2020); Chen et al. (2021). These methods utilize a large number of devices, with each one processing distinct data subsets and participating in an effective data exchange mechanism, thereby aiding in the training of these computationally intensive models. Thus, the problem of classical optimization evolves into a distributed optimization form:

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

where f_i is a function, located on a device i . This formulation encompasses not only distributed learning, where data is dispersed across multiple devices to expedite training and facilitate the storage of large amounts of data, but also extends to federated learning Konečný et al. (2016); Li et al. (2020); Kairouz et al. (2021), where data distribution is motivated by the architecture of the system itself, allowing for decentralized model training while maintaining data privacy and integrity across diverse devices.

A downside of this approach manifests as the complexity associated with the transmission of large-scale data, a phenomenon often referred to as the "communication bottleneck" Gupta et al. (2021).

This bottleneck can significantly impede the efficiency of the system, particularly in scenarios involving extensive data exchange across distributed networks. The challenge intensifies in environments where the bandwidth is limited, requiring solutions to mitigate the impact of data transmission delays and ensure seamless data flow.

The primary solution at present is the compression of transmitted information Bekkerman et al. (2011); Chilimbi et al. (2014); Alistarh et al. (2017), wherein not a whole package is sent, but rather a selected subset. This method involves strategically selecting and compressing the most informative segments of data for transmission. By doing this way, it significantly reduces the volume of data that needs to be communicated across the network, thereby alleviating the communication bottleneck.

In recent times, a number of methods employing compression have been conceived and scrutinized Mishchenko et al. (2019); Gorbunov et al. (2021a); Richtárik et al. (2021). However, a lot of studies have utilized unbiased compression operators due to their simplicity and amenability to theoretical analysis. Such compression techniques, including methods as random sparsification and value rounding Nesterov (2012a); Alistarh et al. (2017); Horvath et al. (2022); Beznosikov et al. (2023a), fail to consider the integration of information conveyed in prior iterations. We hence highlight a potential research gap regarding the usage of previously transmitted data in compression operators and optimization algorithms.

This omission raises the following research questions that we address in the paper:

- *Is it possible to design compression operators that take into account information about what and how we forwarded in previous iterations?*
- *What methods can we integrate this kind of compression operators into? How does it affect the convergence rate of the methods, both in theory and in practice?*
- *Can the methods be made even more efficient, e.g., by using additional momentum acceleration techniques?*

In our paper, we focus on compression-based methods that take into account information collected across multiple preceding iterations, employing what are termed as Markovian compression operators. To the best of our knowledge, this approach emerges as novel and unexplored in the existing literature.

1.1 OUR CONTRIBUTIONS

New type of compression operators. We introduce a novel type of compressors that utilizes stochasticity transmitted over several previous iterations. We refer to this type of compressors as Markovian, because the states of these compressors can be viewed as a Markov chain. We examine two invented examples of such compressors: $\text{BanLast}(K, m)$ (Definition 5) and $\text{KAWASAKI}(K, b, \pi_\Delta, m)$ (Definition 6). The first new compressor operates on a more intuitive basis: it works as random sparsification, but prohibits the transmission of coordinates that were sent in the previous K iterations. The latter functions in terms of probabilities: it reduces the likelihood of transmitting coordinates that appeared in previous iterations. The $\text{KAWASAKI}(K, b, \pi_\Delta, m)$ compressor is more flexible and, in fact, modify the idea $\text{BanLast}(K, m)$, but it introduces two hyperparameters that will be discussed later in Section 2.1.

New algorithms. The compression operators described above give rise to new methods that utilize them. In this context, our paper outlines a general framework based on Alistarh et al. (2017) for distributed gradient descent algorithms that employ Markovian compression operators (MQSGD, see Algorithm 1). Subsequently, to make this basic algorithm faster we apply the multiple momentum technique Nesterov (2012a) and obtain the accelerated method AMQSGD. The formulation of such an algorithm is detailed in Algorithm 2. The basic and accelerated methods are explored both theoretically and experimentally throughout the paper. Furthermore, experiments utilizing Markovian operators in the DIANA Mishchenko et al. (2019) and SGD with momentum algorithms are conducted in Section 3.

Strongly convex and non-convex cases. Motivated by various applications primarily from machine learning, we provide the theoretical analysis in the strongly convex (Theorem 3) and non-convex / PL-condition (Theorem 2) cases of the target function f . Notably, we provide proper analysis for both setups with specific cases, which is rarely present in the field.

Numerical experiments. We conduct experiments with Markovian compressors in a data-parallel setup for several optimization problems and datasets. In particular, we analyze the proposed MQSGD and AMQSGD, as well as the DIANA and SGD optimizers for distributed optimization. In all setups, we observe an acceleration of convergence for methods employing the BanLast and KAWASAKI compressors compared to the baseline random sparsification.

1.2 RELATED WORK

Compressed communications. The use of compressed communications is a fairly well-known idea in distributed learning Seide et al. (2014). As soon as the main property of compressed messages is that they are much easier to transfer, it can be reached in different ways, such as by quantizing the entries of the input vector Alistarh et al. (2017); Mayekar & Tyagi (2019); Gandikota et al. (2020); Horvath et al. (2022), or by sparsifying it Richtárik & Takáč (2016); Alistarh et al. (2018), or even by combining these ideas Albasyoni et al. (2020); Beznosikov et al. (2023a). However, all of the compression operators could be roughly Condat et al. (2023) separated into two large groups: *unbiased* and *biased*.

The first group is much easier to analyze and is therefore more broadly represented in the literature. The basic method with unbiased compression was presented in Alistarh et al. (2017). Later this algorithms were modified using variance reduction technique with compression of gradient differences Mishchenko et al. (2019); Horváth et al. (2019); Gorbunov et al. (2021a) in order to improve the theoretical convergence guarantees. One can also note the works Gorbunov et al. (2019) and Khaled et al. (2020), where the authors developed a general theory for SGD-type methods with unbiased compression.

On the other hand, our understanding of distributed optimization with biased compressors is more complicated. In particular, biased compression implies the use of error compensation techniques Stich et al. (2018). Distributed SGD with biased compression and linear rate of convergence in a multi-node setting was first introduced in Beznosikov et al. (2023a). In the meantime, other error compensation techniques are being actively developed, Lin et al. (2022); Richtárik et al. (2021). The last approach called EF21 was later studied in Fatkhullin et al. (2021), Gruntkowska et al. (2023).

Markovian stochasticity. Another recent trend in the literature is to design algorithms that use Markovian stochastic processes instead of *i.i.d.* random variables in various ways. For instance, Duchi et al. (2012) introduced a version of the Mirror Descent algorithm that yields optimal convergence rates for non-smooth and convex problems. Later, Doan et al. (2020a); Dorfman & Levy (2023); Beznosikov et al. (2023b) studied first-order methods in the Markovian noise setting. Alternatively, token algorithms Hendrikx (2022); Ayache et al. (2022) are also a popular area of research in Markovian stochasticity. In particular, Even (2023) obtained optimal rates of convergence, and Sun et al. (2022); Mao et al. (2019); Doan et al. (2020b) looked at the token algorithm from the angle of the Lagrangian duality and from variants of the ADMM method. At the same time, there exist particular results, e.g., Bresler et al. (2020), which provide a lower bound for the particular finite sum problems in the Markovian setting.

Despite all of the above, to the best of our knowledge, there are currently no works that combine compressed data communications and Markovian stochasticity of the compressors.

1.3 TECHNICAL PRELIMINARIES

Notations. We use $\langle x, y \rangle := \sum_{i=1}^d x_i y_i$ to denote standard inner product of vectors $x, y \in \mathbb{R}^d$ and $(x \odot y)_i = x_i y_i$ to denote Hadamard product of vectors $x, y \in \mathbb{R}^d$. We introduce l_2 -norm of vector $x \in \mathbb{R}^d$ as $\|x\| := \sqrt{\langle x, x \rangle}$. We define $x^* \in \mathbb{R}^d$ as a point, where we reach the minimum in the problem (1). We also denote $f^* > -\infty$ as a global (potentially not unique) minimum of f . We use a standard notation for $(d-1)$ -dimensional simplex $\Delta_d := \left\{ p \in \mathbb{R}^d \mid p_j \geq 0 \text{ and } \sum_{j=1}^d p_j = 1 \right\}$ and for a set of natural numbers $\overline{1, n} := \{1, 2, \dots, n\}$. We denote C_m^k as the binomial coefficient $\binom{m}{k}$.

Throughout the paper, we assume that the objective functions f_i and the function f from (1) satisfy the following assumptions.

Assumption 1 (L_i -smooth). *Every function f_i is L_i -smooth on \mathbb{R}^d with $L_i > 0$, i.e. it is differentiable and there exists a constant $L_i > 0$ such that for all $x, y \in \mathbb{R}^d$ it holds that $\|\nabla f_i(x) - \nabla f_i(y)\|^2 \leq L_i^2 \|x - y\|^2$. We define $L^2 := \frac{1}{n} \sum_{i=1}^n L_i^2$.*

Assumption 2 (μ -strongly convex). *The function f is μ -strongly convex on \mathbb{R}^d , i.e., it is differentiable and there is a constant $\mu > 0$ such that for all $x, y \in \mathbb{R}^d$ it holds that $(\mu/2) \|x - y\|^2 \leq f(x) - f(y) - \langle \nabla f(y), x - y \rangle$.*

Assumption 3 (PL-condition). *The function f satisfies the PL-condition, i.e., it is differentiable and there is a constant $\mu > 0$ such that for all $x \in \mathbb{R}^d$ it holds that $\|\nabla f(x)\|^2 \geq 2\mu (f(x) - f^*)$.*

Assumption 4 (Data similarity). *The functions f_i are similar on \mathbb{R}^d , i.e., there are constants $\delta, \sigma \geq 0$, such that the following inequality holds for all $x \in \mathbb{R}^d$: $\|\nabla f_i(x) - \nabla f(x)\|^2 \leq \delta^2 \|\nabla f(x)\|^2 + \sigma^2$.*

The equation above implies that the data stored at each worker does not differ significantly. This Assumption is quite standard in the literature Shamir et al. (2014); Arjevani & Shamir (2015); Khaled et al. (2020); Woodworth et al. (2020); Gorbunov et al. (2021b); Beznosikov et al. (2022; 2023b).

Now we introduce important definitions related to the theory of Markov processes.

Definition 1 (Markov chain). *Markov chain with a finite state space $\{\nu_n\}_{n=0}^N$ is a stochastic process $\{X_t\}_{t \geq 0}$, that satisfies Markov property, i.e. $\mathbb{P}\{X_t = \nu_t \mid X_{t-1} = \nu_{t-1}, X_{t-2} = \nu_{t-2}, \dots, X_0 = \nu_0\} = \mathbb{P}\{X_t = \nu_t \mid X_{t-1} = \nu_{t-1}\}$.*

Definition 2 (Ergodicity of Markov chain). *Markov chain $\{X_t\}_{t \geq 0}$ with a finite state space $\{\nu_n\}_{n=0}^N$ is referred to be ergodic if for any $n \in \overline{1, N}$ there exists $\lim_{t \rightarrow \infty} \mathbb{P}\{X_t = \nu_n \mid X_0 = \nu_0\} = p_n$, where*

$0 \leq p_n \leq 1$ does not depend on the ν_0 . If Markov chain is ergodic, then $\{p_n\}_{n=0}^N \in \Delta_N$ and there exist $0 < \rho < 1, C > 0$, such that $|\mathbb{P}\{X_t = \nu_n \mid X_0 = \nu_0\} - p_n| \leq C\rho^t$.

Definition 3 (Mixing time of the discrete Markov chain). *We say that $\tau_{\text{mix}}(\varepsilon)$ is the mixing time of the ergodic Markov chain $\{X_t\}_{t \geq 0}$ with stationary distribution $\{p_n\}_{n=0}^N$, if $\forall \varepsilon > 0, \forall t \geq \tau_{\text{mix}}(\varepsilon) \hookrightarrow \max_{n \in \overline{0, N}} \{|\mathbb{P}\{X_t = \nu_n \mid X_0 = \nu_0\} - p_n|\} \leq \varepsilon \cdot p_{\min}$, where $p_{\min} := \min_{n \in \overline{0, N}} \{p_n\}$. From the*

Definition 2, it follows that $\tau_{\text{mix}}(\varepsilon) \geq \frac{\log(C/p_{\min}\varepsilon)}{\log(1/\rho)}$.

These definitions are extremely important for further analysis of the Markovian compressors, which are presented in the next section.

2 MAIN RESULTS

2.1 MARKOVIAN COMPRESSORS

In this section, we introduce Markovian compressors that take into account the information transmitted in previous K operations. It is assumed that these compressors function within an iterative algorithm aimed at minimizing the problem (1), wherein a distinct discrete variable, denoted as the step t , is involved. Consequently, due to the dependence of the compressors on previous states, they exhibit a reliance on the step t . Let us narrow down the class of compressors to be discussed in this paper.

Definition 4 (Random sparsification). *$Q_t(x)$ is a random sparsification compressor, if it operates on the vector $x \in \mathbb{R}^d$ as $Q_t(x) = \frac{d}{m}x \odot \mathbb{1}(\nu_t)$, where ν_t is a set of m coordinates: $\nu_t \subseteq \overline{1, d}$.*

The classical `Randm` operator fits Definition 4, in particular, for this compressor subsets ν_t are generated uniformly at each step t , therefore it is unbiased, i.e., $\mathbb{E}_t[Q_t(x)] = x$ for all t . In this paper, we do not generate ν_t independently, but according to some Markov chain, i.e., compressors start to take into account past iterations. We formulate this idea as an assumption.

Assumption 5 (Asymptotic unbiasedness of Markovian compressors). *We assume that operator Q_t is a random sparsification compressor (Definition 4) and $\{\nu_t\}_{t \geq 0}$ are realizations of some ergodic Markov chain with uniform stationary distribution.*

Assumption 5 implies that in the limit as $t \rightarrow \infty$, the compressor Q_t is unbiased, i.e., $\mathbb{E}[Q_t(x)] \rightarrow x$ as $t \rightarrow \infty$, because the stationary distribution of the Markov chain is uniform. We are now ready to introduce two compressors that adhere to Assumption 5. The first compressor is called `BanLast`(K, m), it prohibits sending coordinates that have been sent at least once in the last K iterations.

Definition 5 (`BanLast`(K, m) compressor). *Let $Q_t(x)$ be a random sparsification compressor (Definition 4). The $j \in \nu_t$ are chosen according to the distribution $p^t \in \Delta_d$ and p^t is given by the formula:*

$$p_j^t = \begin{cases} 0, & \text{if } j \in \bigcup_{s=t-K}^{t-1} \nu_s, \\ \frac{1}{d-Km}, & \text{otherwise.} \end{cases}$$

The `BanLast`(K, m) compressor exhibits a limitation in its utility due to an application restriction: $d \geq (K+1)m$, since we need at least m coordinates to have a non-zero probability at each step t . In order to avoid these limitations, we introduce a more flexible Markovian compressor `KAWASAKI`(K, b, π_Δ, m).

Definition 6 (KAWASAKI(K, b, π_Δ, m) compressor). Let $Q_t(x)$ be a random sparsification compressor (Definition 4). The $j \in \nu_t$ are chosen according to the distribution $p^t \in \Delta_d$, which is given by the formula:

$$\tilde{p}_j^t = \frac{1/d}{b^{\#\text{of choices } j \text{ for the last } K \text{ iterations}}}, \quad j \in \overline{1, d}; \quad p^t = \pi_\Delta(\tilde{p}^t),$$

where $b > 1$ is a forgetting rate and $\pi_\Delta : \mathbb{R}^d \rightarrow \Delta_d$ is an activation function.

The KAWASAKI(K, b, π_Δ, m) compressor is now applicable for arbitrary values of $d \geq m$, and K . However, it introduces two additional hyperparameters in comparison with $\text{BanLast}(K, m)$, namely b and π_Δ . The parameter b is responsible for the how strongly we penalize a coordinate if it was selected in previous iterations, the larger b is, the less likely we are to select a coordinate in step t if it was selected in steps $t - K$ to $t - 1$. The function π_Δ is required in order to obtain the probability vector p^t from the vector \tilde{p}^t , the necessary conditions for this function will be introduced later. The following examples illustrate potential selections for π_Δ :

$$(\pi_\Delta(\tilde{p}))_j = |\tilde{p}_j| / \|\tilde{p}\|_1, \quad \pi_\Delta(\tilde{p}) = \text{Softmax}(\tilde{p}), \quad \pi_\Delta(\tilde{p}) = \arg \min_{p \in \Delta_d} \{\|\tilde{p} - p\|^2\}.$$

We now provide an example where using the Markovian compressor $\text{BanLast}(K, m)$ (Definition 5) speeds up the optimization process by a factor of three compared to the unbiased compressor $\text{Rand}m$.

Example 1. Consider the QSGD algorithm (Algorithm 1), which solves the problem (1) in the case $n = 1$, of the form $x^{t+1} = x^t - \gamma Q(\nabla f(x^t))$. Assume that at some step t we observe gradient of the form $(1, 0, \dots, 0)^T \in \mathbb{R}^d$. In the QSGD algorithm, we compress the gradient at each step, therefore, we do not always send the first coordinate to the server, i.e. we do not move from the point x^t .

In the case of $m = 0.1 \cdot d$, i.e. we send 10% of all coordinates at each step, if we use the $\text{BanLast}(K, m)$ compressor, then the mathematical expectation of the number of steps to leave the point x^t is approximately 3.4 in the case of $K = 7$. For $\text{Rand}10\%$ this number is equal to 10, i.e. we speed up the optimization process by a factor of three. For arbitrary values of d and m , the formula for calculating the number of steps to leave the point x^t is provided in Appendix B.

Moreover, in Appendix B, we obtain more general results for an arbitrary value of $\alpha \in (0; 1]$ with $d = \alpha \cdot m$. In particular, we find the exact expression for the dependence of the number of steps to leave the point x^t . For each fixed α we can find the optimal value of $K^*(\alpha)$. It turns out that empirically this dependence is close to a linear one of the form $K^*(\alpha) \approx 0.73 \cdot \alpha$. Such a rule can be used as an automatic way of choosing K .

We now present a theorem demonstrating that our Markovian compressors from Definitions 5 and 6 satisfy the conditions outlined in Assumption 5.

Theorem 1 (Asymptotic unbiasedness of $\text{BanLast}(K, m)$ and $\text{KAWASAKI}(K, b, \pi_\Delta, m)$). Compressors from Definitions 5 and 6 can be described using Markov chains with states $\{\nu_1, \nu_2, \dots, \nu_K\}_{\nu_1, \dots, \nu_K \in M}$, where M is the set of all subsets of $\overline{1, d}$ of size m . Moreover,

- $\text{BanLast}(K, m)$ (Definition 5) is ergodic with a uniform stationary distribution, if $d > (K+1)m$.
- If $d > (2K+1)m$, then for $\text{BanLast}(K, m)$ we get

$$\rho = \sqrt{1 - \left(\frac{C_{d-2Km}^m}{(C_{d-Km}^m)^2} \right)^K} \quad \text{and} \quad C = \left(1 - \left(\frac{C_{d-2Km}^m}{(C_{d-Km}^m)^2} \right)^K \right)^{-1}.$$

- If for all permutations ϕ of the set $\overline{1, d}$ it holds that $\pi_\Delta(\phi(\tilde{p})) = \phi(\pi_\Delta(\tilde{p}))$, then KAWASAKI(K, b, π_Δ, m) (Definition 6) is ergodic with a uniform stationary distribution.
- If $(\pi_\Delta(\tilde{p}))_j = |\tilde{p}_j| / \|\tilde{p}\|_1$, then

$$\rho = 1 - [db^K - m(b^K - 1)]^{-mK} \quad \text{and} \quad C = \left(1 - [db^K - m(b^K - 1)]^{-mK} \right)^{-1}. \quad (2)$$

The proof of Theorem 1 is provided in Appendix C. The outcomes of Theorem 1 hold significant importance for the subsequent investigation of algorithms aimed at solving problem (1) employing Markovian compressors. Note that the examples of activation functions π_Δ provided above satisfy the conditions of Theorem 1.

2.2 DISTRIBUTED GRADIENT DESCENT WITH MARKOVIAN COMPRESSORS

In this section, we propose a new algorithm `Markovian QSGD` (Algorithm 1). This algorithm is similar to the vanilla QSGD Alistarh et al. (2017), but in line 7 of Algorithm 1 we use Markovian compressor Q_t^i , that we introduced in Section 2.1, i.e., Q_t^i can be either `BanLast`(K, m) (Definition 5) or `KAWASAKI`(K, b, π_Δ, m) (Definition 6).

Theorem 2 (Convergence of MQSGD (Algorithm 1)). *Consider Assumptions 1, 4 and 5. Let the problem (1) be solved by Algorithm 1.*

- For any $\varepsilon, \gamma > 0, T > \tau > \tau_{\text{mix}}(\varepsilon)$ satisfying conditions, described in Appendix E.1, it holds that

$$\mathbb{E} \left[\|\nabla f(\hat{x}^T)\|^2 \right] = \mathcal{O} \left(\frac{F_\tau}{\gamma T} + \frac{\gamma L \tau d^2}{m^2} \sigma^2 \right),$$

where \hat{x}^T is chosen uniformly from $\{x^t\}_{t=0}^T$.

- If f additionally verifies the PL-condition (Assumption 3), then for any $\varepsilon > 0, \gamma > 0, \tau > \tau_{\text{mix}}(\varepsilon)$ and $T > \tau$ satisfying conditions, described in Appendix E.1, it holds that

$$F_T = \mathcal{O} \left(\left(1 - \frac{\mu\gamma}{12}\right)^{T-\tau} F_\tau + \frac{\gamma d^2 L \tau}{\mu m^2} \sigma^2 \right).$$

Here we use the notations $F_t := \mathbb{E}[f(x^t) - f(x^*)]$ and $F_\tau := \mathbb{E}[f(x^\tau) - f(x^*)]$.

The proof of Theorem 2 is provided in Appendix E.3, E.4. If Assumption 4 does not hold we observe different results, which are provided in the Appendix F.

Usually in convergence evaluations of various methods, expressions with the term of F_0 , i.e., something that depends on the initial choice, arise as constants, but in Theorem 2, a term of the form F_τ appears. This can be explained by the fact that at iterations from $t = 0 \rightarrow \tau$ the Markov chain has not yet been stabilized, and the initial state can be taken as $t = \tau$.

Sketch proof of Theorem 2. Let us write out a descent lemma of the form

$$\begin{aligned} \mathbb{E} \left[\|x^{t+1} - x^*\|^2 \right] &= \mathbb{E} \left[\|x^t - x^*\|^2 \right] - 2\mathbb{E} \left[\gamma \langle \nabla f(x^t), x^t - x^* \rangle \right] \\ &\quad - \underbrace{\frac{2\gamma}{n} \sum_{i=1}^n \mathbb{E} \left[\langle Q_t^i(\nabla f(x^t)) - \nabla f_i(x^t), x^t - x^* \rangle \right]}_{\textcircled{1}} + \gamma^2 \mathbb{E} \left[\left\| \frac{1}{n} \sum_{i=1}^n Q_t^i(\nabla f_i(x^t)) \right\|^2 \right]. \end{aligned} \quad (3)$$

The expression $\textcircled{1}$ in (3) is zero if Q_t^i are unbiased and independent from iteration t , because $\mathbb{E} \left[\langle Q_t^i(\nabla f(x^t)) - \nabla f_i(x^t), x^t - x^* \rangle \right] = \mathbb{E} \left[\langle \mathbb{E}_t [Q_t^i(\nabla f(x^t)) - \nabla f_i(x^t)], x^t - x^* \rangle \right] = 0$, where $\mathbb{E}_t[\cdot]$ is the conditional expectation at a step t . Therefore, the theory for such compressors is highly developed. In our case, $Q_t^i(x^s)$ are unbiased only if $t - s \rightarrow \infty$, which follows from asymptotic unbiasedness of our Markovian compressors obtained from Assumption 5. However, we can use some coarsening rather than unbiasedness when $t - s = \tau$, where $\tau > \tau_{\text{mix}}(\varepsilon)$, using the technique of "stepping back" as follows:

$$\mathbb{E} \left[\langle Q_t^i(a^{t-\tau}) - a^{t-\tau}, b^{t-\tau} \rangle \right] \leq \frac{\varepsilon d}{m} \mathbb{E} \left[\|a^{t-\tau}\| \|b^{t-\tau}\| \right]. \quad (4)$$

Importantly, we must apply the compressor Q_t at step t to the vector $a^{t-\tau}$ at step $t - \tau$, since if we apply it to the vector a^t at step t , we will not be able to uncover the conditional expectation, since we will have randomness in a^t (see details in Appendix D). As can be seen from (3) we need to apply the last inequality with $a^{t-\tau} = \nabla f_i(x^{t-\tau})$ and $b^{t-\tau} = x^{t-\tau} - x^*$, but in (3) we only obtain expression with variables at step t , therefore, it has to be handled in some way. In order to resolve this issue we use a straightforward algebra:

$$\begin{aligned}
& \mathbb{E} [\langle Q_t^i (\nabla f_i(x^t)) - \nabla f_i(x^t), x^t - x^* \rangle] = \mathbb{E} [\langle Q_t^i (\nabla f_i(x^{t-\tau})) - \nabla f_i(x^{t-\tau}), x^{t-\tau} - x^* \rangle] \\
& - \mathbb{E} \left[\left\langle Q_t^i (\nabla f_i(x^t) - \nabla f_i(x^{t-\tau})) - \nabla f_i(x^t) + \nabla f_i(x^{t-\tau}), x^t - x^{t-\tau} \right\rangle \right] \\
& + \mathbb{E} \left[\left\langle Q_t^i (\nabla f_i(x^t) - \nabla f_i(x^{t-\tau})) - \nabla f_i(x^t) + \nabla f_i(x^{t-\tau}), x^t - x^* \right\rangle \right] \\
& + \mathbb{E} [\langle Q_t^i (\nabla f_i(x^t)) - \nabla f_i(x^t), x^t - x^{t-\tau} \rangle].
\end{aligned} \tag{5}$$

The first term in the last inequality (5) is solved with the ε -inequality (4), other scalar products are solved using the Fenchel-Young inequality. Terms with $\mathbb{E} \|x^t - x^{t-\tau}\|^2$ are evaluated using line 9 of Algorithm 1: $x^t - x^{t-\tau} = -\gamma \sum_{s=t-\tau}^{t-1} g^s$. Terms with $\mathbb{E} \|Q_t^i (\nabla f_i(x^t) - \nabla f_i(x^{t-\tau}))\|^2$ are obtained from the following inequalities (see details in Appendix E):

$$\|Q_t^i (\nabla f(x) - \nabla f(y))\|^2 \leq \frac{d^2}{m^2} \|\nabla f(x) - \nabla f(y)\|^2 \leq \frac{d^2 L^2}{m^2} \|x - y\|^2,$$

Since the evaluation of $\mathbb{E} \|x^{t+1} - x^*\|^2$ raises the terms of the form $\mathbb{E} \|x^{t-\tau} - x^*\|^2$, we have to do a summation of $\mathbb{E} \|x^{t+1} - x^*\|^2$ from $t = \tau$ to $t = T$. These terms greatly complicate the proof of Theorem 2 compared to the unbiased compressors. The results of Theorem 2 can be rewritten as an upper complexity bound on a number of iterations T of the Algorithm 1 by carefully tuning the step size γ .

Corollary 1 (Step tuning for Theorem 2).

• Under the conditions of Theorem 2 in the non-convex case, choosing γ as in Appendix E.2, in order to achieve the ε -approximate solution (in terms of $\mathbb{E} [\|\nabla f(x^T)\|^2] \leq \varepsilon^2$), it takes

$$\mathcal{O} \left(\frac{L\tau d^2}{m^2} F_\tau \left(\frac{\delta^2 + 1}{\varepsilon^2} + \frac{\sigma^2}{\varepsilon^4} \right) \right) \text{ iterations of Algorithm 1.}$$

• Under the conditions of Theorem 2 in the PL-condition (Assumption 3) case, choosing γ as in Appendix E.2 in order to achieve the ε -approximate solution (in terms of $\mathbb{E} [f(x^t) - f(x^*)] \leq \varepsilon$), it takes

$$\mathcal{O} \left(\frac{d^2 L\tau}{m^2 \mu} \left((\delta^2 + 1) \log \left(\frac{1}{\varepsilon} \right) + \frac{\sigma^2}{\mu \varepsilon} \right) \right) \text{ iterations of Algorithm 1.}$$

2.3 ACCELERATED METHOD

After giving the convergence result for the vanilla distributed SGD with Markovian compression operator, we now move on to the accelerated scheme. Since we do not assume boundedness of the gradient variance, the classical Nesterov acceleration Nesterov (2014) does not produce the expected effect, and therefore an additional momentum has to be introduced Nesterov (2012b); Vaswani et al. (2019). By applying a multi-step strategy partially similar to Beznosikov et al. (2023b), we obtain our Algorithm 2.

Algorithm 2 Accelerated Markovian QSGD (AMQSGD)

- 1: **Input:** starting point $x^0 \in \mathbb{R}^d$, step size $\gamma > 0$, momentums θ, η, β, p , number of iterations T
 - 2: **for** $t = 0$ **to** T **do**
 - 3: Update $x_g^t = \theta x_f^t + (1 - \theta)x^t$
 - 4: Broadcast x_g^t to all workers
 - 5: **for** $i = 1$ **to** n **in parallel do**
 - 6: Set $g_i^t = Q_t^i (\nabla f_i(x_g^t))$
 - 7: Send g_i^t to the server
 - 8: **end for**
 - 9: Aggregate $g^t = \frac{1}{n} \sum_{i=1}^n g_i^t$
 - 10: Update $x_f^{t+1} = x_g^t - p\gamma g^t$
 - 11: Update $x^{t+1} = \eta x_f^{t+1} + (p - \eta)x_f^t$
 - 12: + $(1 - p)(1 - \beta)x^t + (1 - p)\beta x_g^t$
 - 13: **end for**
-

Theorem 3 (Convergence of AMQSGD (Algorithm 2)). Consider Assumptions 1, 2, 4. Let the problem (1) be solved by Algorithm 2. Then for any $\gamma, \varepsilon > 0, T > \tau > \tau_{\text{mix}}(\varepsilon), \beta, \theta, \eta, p$ satisfying conditions, described in Appendix G.1, it holds that

$$F_{T+1} = \mathcal{O} \left(\exp \left[-(T - \tau) \sqrt{\frac{p^2 \mu \gamma}{3}} \right] F_\tau + \exp \left[-T \sqrt{\frac{p^2 \mu \gamma}{3}} \right] \Delta_\tau + \frac{\gamma}{\mu} \sigma^2 \right).$$

Here we use the notations: $F_t := \mathbb{E}[\|x^t - x^*\|^2 + 3/\mu(f(x_f^t) - f(x^*))]$ and $\Delta_\tau \leq \gamma^{1/2} \tau^{-4/3} \mu^{-1/3} \sum_{t=0}^{\tau} (\mathbb{E} \|\nabla f(x_g^t)\|^2 + \mathbb{E} \|x^t - x^*\|^2 + \mathbb{E}[f(x_f^t) - f(x^*)])$.

The above theorem shows that in the strongly convex case Accelerated Markov QSGD with constant step-size can attain sublinear convergence. In terms of dealing with Markovian stochasticity, its proof follows quite similar ideas as the proof of Theorem 2: here again we use the technique of *stepping back* for mixing time, which allows us to effectively deal with the bias of the gradient estimator. The full proof is provided in Appendix G.3. The results of Theorem 3 can be rewritten as an upper complexity bound on a number of iterations T of the Algorithm 2 by carefully tuning the step size γ .

Corollary 2 (Step tuning for Theorem 3). *Under the conditions of Theorem 3, choosing γ as in Appendix G.2 in order to achieve the ϵ -approximate solution (in terms of $\mathbb{E} [\|x^T - x^*\|^2] \leq \epsilon^2$), it takes*

$$\mathcal{O} \left(\frac{d^2 L^{\frac{2}{3}} \tau^{\frac{4}{3}}}{m^2 \mu^{\frac{2}{3}}} \left((\delta^2 + 1) \log \left(\frac{1}{\epsilon} \right) + \frac{\sigma^2}{\mu \epsilon} \right) \right) \text{ iterations of Algorithm 2.}$$

2.4 DISCUSSION

Our Example 1 and the numerical experiments in Section 3 show that the using of Markovian compressors could lead to a better performance quite well, however, the theoretical guarantees turn out to be poorer than in the unbiased case. In particular, if we use *Randm* in the QSGD algorithm, then we observe the following estimates Beznosikov et al. (2023a):

$$X_T = \mathcal{O} \left((1 - \mu\gamma)^T X_0 + \gamma \frac{d}{m} \frac{\sigma^2}{\mu n} \right),$$

where $X_t = \mathbb{E} [\|x^t - x^*\|^2]$ and $\gamma \lesssim \frac{1}{L(1+d/mn)}$. However, Theorem 2 gives us such estimates:

$$F_T = \mathcal{O} \left(\left(1 - \frac{\mu\gamma}{12} \right)^T F_\tau + \gamma \frac{d^2}{m^2} \frac{\tau L \sigma^2}{\mu} \right),$$

where $F_t := \mathbb{E} [f(x^T) - f(x^*)]$ and $\gamma \lesssim \frac{m^2}{L d^2 \tau (\delta^2 + 1)}$. It is important to note that not only has the theory for Markovian compressors not yet been studied, but also dealing with the Markovian stochasticity itself implies quite strict limitations. For instance,

d/m vs d²/m². We are forced to uniformly bound the noise of the compressor (linearity in the compression constant is prevented by this) due to the impossibility of using the expectation trick, in contrast to the unbiased case Beznosikov et al. (2023a), where the authors estimated the variance of the compressor noise. The assumption of uniformly bounded noise cannot be rejected by any authors who work with Markovian stochasticity Beznosikov et al. (2023b); Dorfman & Levy (2023); Doan et al. (2020a); Sun et al. (2018); Even (2023), therefore, there is no possibility to achieve linearity in the compression rate in our theoretical guaranties, according to the current theoretical advances.

Mixing time. Furthermore, it is imperative to emphasize that it follows from Theorems 2 and 3 that the convergence rate is improved as τ (and, consequently, K) diminishes. In other words, the distribution of the compressor's underlying Markov chain has to converge to a uniform distribution as fast as possible, but empirically one wants the choice of coordinates to depend on previous iterations rather than be random (e.g. for *Randm* compressor $\tau = 1, K = 0$). This causes a logical contradiction: while using a large K will theoretically give poorer convergence, in practice algorithms with non-zero values of K perform better (see Section 3). It is also worth mentioning that when Markovian stochasticity is employed, we can never avoid τ in our estimates, since it appears in the lower bounds on the convergence rate of methods that involve Markovian properties Bresler et al. (2020). Thus, our Algorithms 1 and 2 have a reasonably good polynomial dependence on mixing time (Theorem 2 shows an optimal estimation in terms of τ), considering the fact there are several works Doan et al. (2020b) whose bounds include terms that are even *exponential* in the mixing time.

L/μ. In spite of the difficulties listed above, we still can observe that the momentums implementation in Algorithm 2 gives an acceleration in terms of L/μ compared to vanilla QSGD (Algorithm 1). In

the classical version of accelerated Gradient Descent, one can achieve an acceleration of the form $\sqrt{L/\mu}$ Nesterov (1983), but our analysis allows only to achieve $(L/\mu)^{2/3}$ in Theorem 3. When Markovian stochasticity is employed, it is also possible to achieve estimation of the form $\sqrt{L/\mu}$ Beznosikov et al. (2023b), but it is obtained by using batches with size scaled as 2^j , where j is drawn from a truncated geometric distribution. Unfortunately, this specific batching technique cannot be applied in our paper, as we consider compressors that act as random sparsification (Definition 4), which necessitates that the gradient be compressed only once at each iteration.

Variance reduction. In our paper, we focus on the QSGD method and its accelerated version (Algorithms 1 and 2). However, in modern studies on distributed optimization, techniques of variance reduction are of a great interest (DIANA Mishchenko et al. (2019), MARINA Gorbunov et al. (2021a), DASHA Tyurin & Richtárik (2022)), because these methods converge linearly to the exact solution of the problem (1), while QSGD (Algorithms 1 and 2) converges only to the σ^2 -neighborhood of the solution. We implement Markovian compressors (Definitions 5 and 6) in these methods in our experiments, but we do not provide theoretical guarantees for such algorithms since we have just developed a theoretical baseline for the study of Markovian compressors. This represents a promising direction for future research.

Even though it is not entirely clear whether it is possible to achieve significant improvements in the theoretical results, due to the peculiarities of dealing with Markovian randomness, for now we could only highlight a significantly better performance of Algorithms 1 and 2 compared to a similar algorithms using a vanilla unbiased compressor Rand_m (see Section 3).

3 EXPERIMENTS

In order to justify the practical usage of the proposed methods and analyze their behavior, we conduct a series of experiments using Markovian compression on distributed optimization problems, specifically logistic regression and neural network-based image classification. We observe that Markovian compressors, when used with MQSGD and AMQSGD, as well as with classical SGD and DIANA Mishchenko et al. (2019), improve convergence on several benchmarks. Appendix H provides a description of the technical setup, extended experiments with hyperparameters analysis, and an application of Markovian compressors to model-parallel neural network training.

3.1 LOGISTIC REGRESSION

Firstly, we experiment on a classification task using a logistic regression model with L_2 regularization of the form:

$$\min_{w \in \mathbb{R}^d} \left\{ f(w) = \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_s w^T x_s)) + \lambda \|w\|^2 \right\},$$

with a regularization term $\lambda = 0.05$. The dataset is split among $n = 10$ clients. We use Mushrooms, A9A, and W8A datasets from LibSVM Chang & Lin (2011) and MNIST Deng (2012). Experiments are conducted using Python 3.10 and PyTorch, and a distributed environment is simulated. We experiment with MQSGD, AMQSGD, and DIANA optimizers, employing Rand -10% as a sparsification compressor. Markovian compressors were utilized independently on each client, with normalization activation function, and with all hyperparameters being fine-tuned.

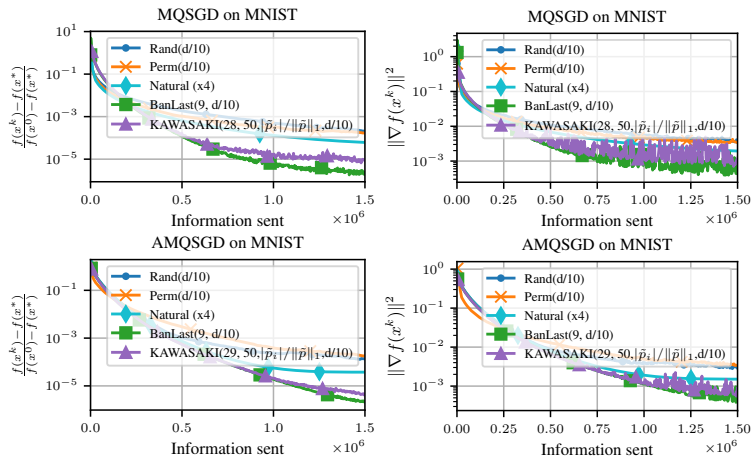


Figure 1: Logistic Regression on MNIST experiments results. All hyperparameters are fine-tuned, and best runs are selected.

Figure 1 shows the convergence of the Rand -10% baseline and Markovian compressors on the MQSGD and AMQSGD algorithms on MNIST dataset. Both Markovian compressors achieve faster convergence

than the baseline and more complex compressors like PermK Szlendak et al. (2021) and Natural compressors Horvath et al. (2022). In most of our results, BanLast and KAWASAKI show similar performance with fine-tuned hyperparameters. Experiments on other datasets, and tuning history size K tuning analysis appear in Appendix H.2. Additionally, as our compressors are fully compatible with classical compressors, we conduct experiments on combination with Natural compression in Appendix H.5.

3.2 NEURAL NETWORKS

We also apply Markovian compressors in more complex optimization tasks, such as image classification on CIFAR-10 Krizhevsky et al. (2009) dataset with ResNet-18 convolutional neural network He et al. (2016). Formally, we solve optimization problem:

$$\min_{w \in \mathbb{R}^d} \left\{ f(w) = \frac{1}{n} \sum_{i=1}^n l(\text{softmax}(f(x_i, w)), y_i) \right\},$$

where x_i is a training image, y_i is its respective class, and $l()$ is a cross-entropy loss function. Dataset is split equally between $n = 5$ clients. We use Rand-5% sparsification operator and SGD optimizer with cosine annealing LR schedule. Hyperparameters, such as the learning rate, batch size, and Markovian-specific ones are fine-tuned.

Figure 2 depicts the training loss and gradient norm, with the aggregate values shown in Table 1. As in the previous case, the application of the Markovian compressor favours faster convergence and better validation results. Note that for more complex optimization task, smoother history accumulation (as in KAWASAKI) is required.

Figure 3 presents comparison with Permutation and Natural compression, which confirm practical usefulness of Markovian compressors on more complex and non-convex optimization problems. Note that our compressors can be applied in combination with complex randomized compressor like Natural compression, making our method even more flexible.

4 CONCLUSION

In this paper, we propose a family of compression schemes, which takes into account previous iterations of algorithm and transform the input vector according to a Markov chain. We develop two sparsification methods BanLast (Definition 5) and KAWASAKI (Definition 6) based on this idea. These compressors are implemented in QSGD (Algorithm 1) and accelerated QSGD (Algorithm 2). We provide convergence rates under different assumptions on the objective function (Theorems 2 and 3). In experiments, we show that our compression methods outperform the baselines in the deep neural network optimisation problem. Future research may consider the implementation of our Markovian compressors in other optimization methods, e.g. using the variance reduction techniques.

Table 1: Numerical results of training ResNet-18 on CIFAR-10 with different compressors. Each cell represents mean \pm standard deviation over 5 runs.

	Rand-5%	Banlast	KAWASAKI
Train Loss	0.0743 \pm 0.003	0.0734 \pm 0.003	0.0305 \pm 0.001
Gradient Norm	1.403 \pm 0.029	1.383 \pm 0.035	0.745 \pm 0.015
Test Accuracy	87.9 \pm 0.179	88.0 \pm 0.122	89.05 \pm 0.294

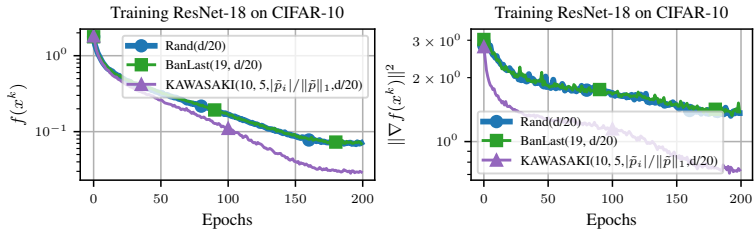


Figure 2: Image classification with ResNet-18 on CIFAR-10 experiments results. Best runs for each method are displayed.

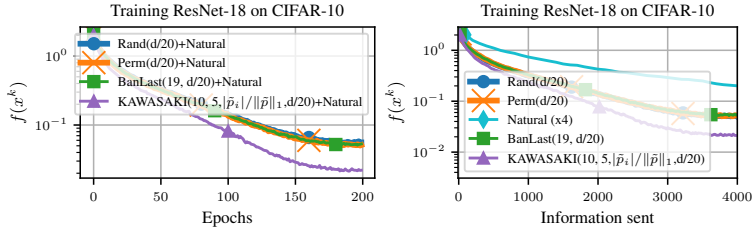


Figure 3: Comparison with other compressors on Resnet-18 training on CIFAR-10 dataset for Rand-5% sparsification on $N = 20$ clients. Natural compression factor is 4. Left figure is sequential combination with Natural compression. Right figure is comparison against PermK and Natural compressors independently, with information sent on x-axis.

540 REFERENCES

- 541 Alyazeed Albasyoni, Mher Safaryan, Laurent Condat, and Peter Richtárik. Optimal gradient com-
542 pression for distributed and federated learning, 2020.
- 543
- 544 Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. Qsgd: Communication-
545 efficient sgd via gradient quantization and encoding. *Advances in neural information processing*
546 *systems*, 30, 2017.
- 547 Dan Alistarh, Torsten Hoefler, Mikael Johansson, Sarit Khirirat, Nikola Konstantinov, and Cédric
548 Renggli. The convergence of sparsified gradient methods, 2018.
- 549
- 550 Yossi Arjevani and Ohad Shamir. Communication complexity of distributed convex learning and
551 optimization. *Advances in neural information processing systems*, 28, 2015.
- 552 Ghadir Ayache, Venkat Dassari, and Salim El Rouayheb. Walk for learning: A random walk approach
553 for federated learning from heterogeneous data, 2022.
- 554
- 555 Ron Bekkerman, Mikhail Bilenko, and John Langford. *Scaling up machine learning: Parallel and*
556 *distributed approaches*. Cambridge University Press, 2011.
- 557 Aleksandr Beznosikov, Pavel Dvurechenskii, Anastasiia Koloskova, Valentin Samokhin, Sebastian U
558 Stich, and Alexander Gasnikov. Decentralized local stochastic extra-gradient for variational
559 inequalities. *Advances in Neural Information Processing Systems*, 35:38116–38133, 2022.
- 560 Aleksandr Beznosikov, Samuel Horváth, Peter Richtárik, and Mher Safaryan. On biased compression
561 for distributed learning. *Journal of Machine Learning Research*, 24(276):1–50, 2023a.
- 562 Aleksandr Beznosikov, Sergey Samsonov, Marina Sheshukova, Alexander Gasnikov, Alexey Naumov,
563 and Eric Moulines. First order methods with markovian noise: from acceleration to variational
564 inequalities. *arXiv preprint arXiv:2305.15938*, 2023b.
- 565
- 566 Song Bian, Dacheng Li, Hongyi Wang, Eric P. Xing, and Shivaram Venkataraman. Does compressing
567 activations help model parallel training?, 2023.
- 568
- 569 Lukas Biewald. Experiment tracking with weights and biases, 2020. URL [https://](https://www.wandb.com/)
570 www.wandb.com/. Software available from wandb.com.
- 571 Guy Bresler, Prateek Jain, Dheeraj Nagaraj, Praneeth Netrapalli, and Xian Wu. Least squares
572 regression with markovian data: Fundamental limits and algorithms, 2020.
- 573 Chih-Chung Chang and Chih-Jen Lin. Libsvm: a library for support vector machines. *ACM*
574 *transactions on intelligent systems and technology (TIST)*, 2(3):1–27, 2011.
- 575
- 576 Mingzhe Chen, Deniz Gündüz, Kaibin Huang, Walid Saad, Mehdi Bennis, Aneta Vulgarakis Feljan,
577 and H Vincent Poor. Distributed learning in wireless networks: Recent progress and future
578 challenges. *IEEE Journal on Selected Areas in Communications*, 39(12):3579–3605, 2021.
- 579 Trishul Chilimbi, Yutaka Suzue, Johnson Apacible, and Karthik Kalyanaraman. Project adam:
580 Building an efficient and scalable deep learning training system. In *11th USENIX symposium on*
581 *operating systems design and implementation (OSDI 14)*, pp. 571–582, 2014.
- 582
- 583 Laurent Condat, Kai Yi, and Peter Richtárik. Ef-bv: A unified theory of error feedback and variance
584 reduction mechanisms for biased and unbiased compression in distributed optimization, 2023.
- 585 Li Deng. The mnist database of handwritten digit images for machine learning research. *IEEE Signal*
586 *Processing Magazine*, 29(6):141–142, 2012.
- 587
- 588 Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. Gpt3. int8 (): 8-bit matrix
589 multiplication for transformers at scale. *Advances in Neural Information Processing Systems*, 35:
590 30318–30332, 2022.
- 591 Michael Diskin, Alexey Bukhtiyarov, Max Ryabinin, Lucile Saulnier, Anton Sinitsin, Dmitry Popov,
592 Dmitry V Pyrkin, Maxim Kashirin, Alexander Borzunov, Albert Villanova del Moral, et al.
593 Distributed deep learning in open collaborations. *Advances in Neural Information Processing*
Systems, 34:7879–7897, 2021.

- 594 Thinh T. Doan, Lam M. Nguyen, Nhan H. Pham, and Justin Romberg. Convergence rates of
595 accelerated markov gradient descent with applications in reinforcement learning, 2020a.
596
- 597 Thinh T. Doan, Lam M. Nguyen, Nhan H. Pham, and Justin Romberg. Finite-time analysis of
598 stochastic gradient descent under markov randomness, 2020b.
- 599 Ron Dorfman and Kfir Y. Levy. Adapting to mixing time in stochastic optimization with markovian
600 data, 2023.
- 601 John C. Duchi, Alekh Agarwal, Mikael Johansson, and Michael I. Jordan. Ergodic mirror descent,
602 2012.
- 603
- 604 Mathieu Even. Stochastic gradient descent under markovian sampling schemes, 2023.
605
- 606 Ilyas Fatkhullin, Igor Sokolov, Eduard Gorbunov, Zhize Li, and Peter Richtárik. Ef21 with bells &
607 whistles: Practical algorithmic extensions of modern error feedback, 2021.
- 608 Venkata Gandikota, Daniel Kane, Raj Kumar Maity, and Arya Mazumdar. vqsgd: Vector quantized
609 stochastic gradient descent, 2020.
- 610
- 611 Eduard Gorbunov, Filip Hanzely, and Peter Richtárik. A unified theory of sgd: Variance reduction,
612 sampling, quantization and coordinate descent, 2019.
- 613
- 614 Eduard Gorbunov, Konstantin P Burlachenko, Zhize Li, and Peter Richtárik. Marina: Faster non-
615 convex distributed learning with compression. In *International Conference on Machine Learning*,
616 pp. 3788–3798. PMLR, 2021a.
- 617
- 618 Eduard Gorbunov, Filip Hanzely, and Peter Richtárik. Local sgd: Unified theory and new efficient
619 methods. In *International Conference on Artificial Intelligence and Statistics*, pp. 3556–3564.
620 PMLR, 2021b.
- 621
- 622 Kaja Grunkowska, Alexander Tyurin, and Peter Richtárik. Ef21-p and friends: Improved theoretical
623 communication complexity for distributed optimization with bidirectional compression, 2023.
- 624
- 625 Vipul Gupta, Avishek Ghosh, Michal Dereziński, Rajiv Khanna, Kannan Ramchandran, and Michael
626 Mahoney. Localnewton: Reducing communication bottleneck for distributed learning. *arXiv
627 preprint arXiv:2105.07320*, 2021.
- 628
- 629 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
630 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
631 pp. 770–778, 2016.
- 632
- 633 Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. Deberta: Decoding-enhanced bert
634 with disentangled attention, 2021. URL <https://arxiv.org/abs/2006.03654>.
- 635
- 636 Hadrien Hendrikx. A principled framework for the design and analysis of token algorithms, 2022.
- 637
- 638 Samuel Horvath, Chen-Yu Ho, Ludovít Horvath, Atal Narayan Sahu, Marco Canini, and Peter
639 Richtárik. Natural compression for distributed deep learning, 2022.
- 640
- 641 Samuel Horváth, Dmitry Kovalev, Konstantin Mishchenko, Sebastian Stich, and Peter Richtárik.
642 Stochastic distributed learning with gradient quantization and variance reduction, 2019.
- 643
- 644 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
645 and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021. URL <https://arxiv.org/abs/2106.09685>.
- 646
- 647 Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin
648 Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. Ad-
649 vances and open problems in federated learning. *Foundations and Trends® in Machine Learning*,
650 14(1–2):1–210, 2021.
- 651
- 652 Ahmed Khaled, Konstantin Mishchenko, and Peter Richtárik. Tighter theory for local sgd on identical
653 and heterogeneous data. In *International Conference on Artificial Intelligence and Statistics*, pp.
654 4519–4529. PMLR, 2020.

- 648 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
649 *arXiv:1412.6980*, 2014.
- 650 B Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A Al Sallab, Senthil Yogamani,
651 and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey. *IEEE*
652 *Transactions on Intelligent Transportation Systems*, 23(6):4909–4926, 2021.
- 653 Jakub Konečný, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh, and
654 Dave Bacon. Federated learning: Strategies for improving communication efficiency. *arXiv*
655 *preprint arXiv:1610.05492*, 2016.
- 656 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- 657 Li Li, Yuxi Fan, Mike Tse, and Kuo-Yi Lin. A review of applications in federated learning. *Computers*
658 *& Industrial Engineering*, 149:106854, 2020.
- 659 Chung-Yi Lin, Victoria Kostina, and Babak Hassibi. Differentially quantized gradient methods, 2022.
- 660 Xianghui Mao, Kun Yuan, Yubin Hu, Yuantao Gu, Ali H. Sayed, and Wotao Yin. Walkman: A
661 communication-efficient random-walk algorithm for decentralized optimization, 2019.
- 662 Prathamesh Mayekar and Himanshu Tyagi. Ratq: A universal fixed-length quantizer for stochastic
663 optimization, 2019.
- 664 Konstantin Mishchenko, Eduard Gorbunov, Martin Takáč, and Peter Richtárik. Distributed learning
665 with compressed gradient differences. *arXiv preprint arXiv:1901.09269*, 2019.
- 666 Yu. Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems.
667 *SIAM Journal on Optimization*, 22(2):341–362, 2012a. doi: 10.1137/100802001. URL <https://doi.org/10.1137/100802001>.
- 668 Yurii Nesterov. A method of solving a convex programming problem with convergence rate $o(1/k^{**$
669 $2)$. *Doklady Akademii Nauk SSSR*, 269(3):543, 1983.
- 670 Yurii Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems. *SIAM*
671 *J. Optim.*, 22:341–362, 2012b. URL [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:1424102)
672 1424102.
- 673 Yurii Nesterov. *Introductory Lectures on Convex Optimization: A Basic Course*. Springer Publishing
674 Company, Incorporated, 1 edition, 2014. ISBN 1461346916.
- 675 J v Neumann. Proof of the quasi-ergodic hypothesis. *Proceedings of the National Academy of*
676 *Sciences*, 18(1):70–82, 1932.
- 677 Peter Richtárik and Martin Takáč. Parallel coordinate descent methods for big data optimization.
678 *Mathematical Programming*, 156:433–484, 2016.
- 679 Peter Richtárik, Igor Sokolov, and Ilyas Fatkhullin. Ef21: A new, simpler, theoretically better,
680 and practically faster error feedback. *Advances in Neural Information Processing Systems*, 34:
681 4384–4396, 2021.
- 682 Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of mathematical*
683 *statistics*, pp. 400–407, 1951.
- 684 Frank Seide, Hao Fu, Jasha Droppo, Gang Li, and Dong Yu. 1-bit stochastic gradient descent and
685 its application to data-parallel distributed training of speech dnns. In *Interspeech*, 2014. URL
686 <https://api.semanticscholar.org/CorpusID:2189412>.
- 687 Ohad Shamir, Nati Srebro, and Tong Zhang. Communication-efficient distributed optimization
688 using an approximate newton-type method. In *International conference on machine learning*, pp.
689 1000–1008. PMLR, 2014.
- 690 Jaeyong Song, Jinkyu Yim, Jaewon Jung, Hongsun Jang, Hyung-Jin Kim, Youngsok Kim, and Jinho
691 Lee. Optimus-cc: Efficient large nlp model training with 3d parallelism aware communication
692 compression. In *Proceedings of the 28th ACM International Conference on Architectural Support*
693 *for Programming Languages and Operating Systems, Volume 2*, pp. 560–573, 2023.

- 702 Sebastian U Stich, Jean-Baptiste Cordonnier, and Martin Jaggi. Sparsified sgd with memory. *Advances*
703 *in Neural Information Processing Systems*, 31, 2018.
- 704
- 705 Tao Sun, Yuejiao Sun, and Wotao Yin. On markov chain gradient descent. *Advances in neural*
706 *information processing systems*, 31, 2018.
- 707
- 708 Tao Sun, Dongsheng Li, and Bao Wang. Adaptive random walk gradient descent for decentralized
709 optimization. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu,
710 and Sivan Sabato (eds.), *Proceedings of the 39th International Conference on Machine Learning*,
711 volume 162 of *Proceedings of Machine Learning Research*, pp. 20790–20809. PMLR, 17–23 Jul
2022. URL <https://proceedings.mlr.press/v162/sun22b.html>.
- 712
- 713 Rafał Szlendak, Alexander Tyurin, and Peter Richtárik. Permutation compressors for provably faster
714 distributed nonconvex optimization. *arXiv preprint arXiv:2110.03300*, 2021.
- 715
- 716 Alexander Tyurin and Peter Richtárik. Dasha: Distributed nonconvex optimization with communica-
717 tion compression, optimal oracle complexity, and no client synchronization, 2022.
- 718
- 719 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
720 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing*
systems, 30, 2017.
- 721
- 722 Sharan Vaswani, Francis Bach, and Mark Schmidt. Fast and faster convergence of sgd for over-
723 parameterized models and an accelerated perceptron, 2019.
- 724
- 725 Joost Verbraeken, Matthijs Wolting, Jonathan Katzy, Jeroen Kloppenburg, Tim Verbelen, and Jan S
726 Rellermeyer. A survey on distributed machine learning. *Acm computing surveys (csur)*, 53(2):
1–33, 2020.
- 727
- 728 Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. Glue:
729 A multi-task benchmark and analysis platform for natural language understanding, 2019. URL
<https://arxiv.org/abs/1804.07461>.
- 730
- 731 Blake E Woodworth, Kumar Kshitij Patel, and Nati Srebro. Minibatch vs local sgd for heterogeneous
732 distributed learning. *Advances in Neural Information Processing Systems*, 33:6281–6292, 2020.
- 733
- 734 Maziar Yazdani and Fariborz Jolai. Lion optimization algorithm (loa): a nature-inspired metaheuristic
735 algorithm. *Journal of computational design and engineering*, 3(1):24–36, 2016.
- 736
- 737
- 738
- 739
- 740
- 741
- 742
- 743
- 744
- 745
- 746
- 747
- 748
- 749
- 750
- 751
- 752
- 753
- 754
- 755

756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

Supplementary Material

CONTENTS

B	Mathematical calculations from Example 1	16
C	Proof of Theorem 1	17
D	Main lemmas	19
E	Extensions for Theorem 2	21
	E.1 Full version of Theorem 2	21
	E.2 Full version of Corollary 1	22
	E.3 Proof of Theorem 2, non-convex case	22
	E.4 Proof of Theorem 2, Under PL-condition	26
F	Convergence of Algorithm 1 without data similarity	29
G	Extensions for Theorem 3	32
	G.1 Full version of Theorem 3	32
	G.2 Full version of Corollary 2	33
	G.3 Proof of Theorem 6	33
H	Experiments	43
	H.1 Technical details	44
	H.2 Logistic Regression experiments	44
	H.3 Dependence on size history	45
	H.4 Comparison with Permutation & Natural Compression	45
	H.5 Combination with other compressors	46
	H.6 Neural Networks Experiments: Data Parallelism Case	47
	H.7 Neural Networks Experiments: Model Parallelism Case	48
	H.8 Fine-tuning DeBERTaV3-base on GLUE development set	49

810 A AUXILIARY LEMMAS AND FACTS

811 In this section we list auxiliary facts and our results that we use several times in our proofs.

812 A.1 CAUCHY-SCHWARZ INEQUALITY

813 For all $x, y \in \mathbb{R}^d$

$$814 \langle x, y \rangle \leq \|x\| \|y\|.$$

815 A.2 FENCHEL-YOUNG INEQUALITY

816 For all $x, y \in \mathbb{R}^d$ and $\beta > 0$

$$817 2 \langle x, y \rangle \leq \beta^{-1} \|x\|^2 + \beta \|y\|^2.$$

818 B MATHEMATICAL CALCULATIONS FROM EXAMPLE 1

819 By definition of the mathematical expectation of an integer positive random variable Z , we obtain
 820 that $\mathbb{E}[Z] = \sum_{s=1}^{\infty} s \cdot \mathbb{P}\{Z = s\}$. In our problem, Z is the number of an iteration where we first
 821 selected the desired coordinate. For $\text{Rand}m$ compressor, we have $\mathbb{P}\{Z = s\} = \frac{m}{d} \cdot \left(1 - \frac{m}{d}\right)^{s-1}$.
 822 The first term is the probability of picking the desired coordinate at iteration s and the second term
 823 is the probability of not picking the desired coordinate at iterations from 1 to $s - 1$. Using this, the
 824 mathematical expectation of the number of steps to quit the point x^t for $\text{Rand}m$ compressor is equal
 825 to

$$826 \sum_{s=1}^{\infty} s \left(1 - \frac{m}{d}\right)^{s-1} \frac{m}{d} = \frac{d}{m}. \quad (6)$$

827 Now we calculate the expectation for $\text{BanLast}(K, m)$ compressor (Definition 5). If $s > K$,
 828 similarly to the $\text{Rand}m$ case, we obtain that $\mathbb{P}\{Z = s\} = \frac{m}{d-Km} \left(1 - \frac{m}{d-Km}\right)^{s-1}$, because we
 829 cannot choose Km coordinates. If $s \leq K$, then the formula of $\mathbb{P}\{Z = s\}$ becomes a bit more
 830 complicated, because the probability of not picking the desired coordinate at iterations from 1 to
 831 $s - 1$ is different at each iteration and is equal to $\prod_{h=0}^{s-2} \left(1 - \frac{m}{d-hm}\right)$. If $s = 1$, then this probability
 832 is equal to one. Using this, we can calculate the mathematical expectation of the number of steps to
 833 leave the point x^t for $\text{BanLast}(K, m)$ compressor:

$$834 \begin{aligned} & \sum_{s=1}^K \frac{sm}{d - (s-1)m} \prod_{h=0}^{s-2} \left(1 - \frac{m}{d-hm}\right) + \sum_{s=K+1}^{\infty} s \left(1 - \frac{m}{d-Km}\right)^{s-1} \frac{m}{d-Km} \\ &= \sum_{s=1}^K \frac{sm}{d - (s-1)m} \prod_{h=0}^{s-2} \left(1 - \frac{m}{d-hm}\right) + \frac{d}{m} \left(1 - \frac{m}{d-Km}\right)^K \\ &= \sum_{s=1}^K \frac{s}{\alpha - (s-1)} \prod_{h=0}^{s-2} \left(1 - \frac{1}{\alpha-h}\right) + \alpha \left(1 - \frac{1}{\alpha-K}\right)^K, \end{aligned} \quad (7)$$

835 where we used the notation $\alpha = d/m$ to show that (7) depends only on d/m , but not on d and m
 836 separately. We can consider (7) as an optimization problem with respect to K . Since K is an integer
 837 and the objective function in (7) is complex, we numerically find the optimal K for different α . For
 838 the sake of clarity, we show the difference between formulas (6) and (7) on Figure 4(c).

We consider $\alpha \in [5.3, 6.7, 8.3, 10, 11.1, 12.5, 14.3, 16.7, 20]$ and find the optimal K by a complete brute force search – see Figure 4 (a). Then, we perform a linear approximation and obtain the formula $K^*(\alpha) \approx 0.7323\alpha$ – see Figure 4 (b). Since the correlation coefficient between the points and the approximated line is equal to 0.73, we can consider this formula to be accurate enough for practical applications.

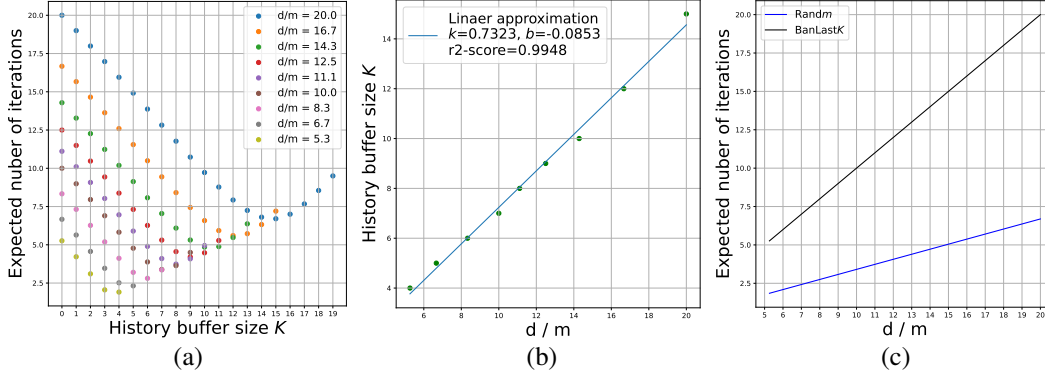


Figure 4: Theoretical estimate on dependence of history buffer size K on parameter $\alpha = d/m$: (a) represents expected number of iterations required to transfer all coordinates to server on history buffer size K for different α , (b) represents scaling of optimal history buffer size K^* on α . (c) represents comparison of expected number of iterations required to transfer all coordinates to server on problems parameter α for Randm and BanLastK.

C PROOF OF THEOREM 1

Lemma 1. *If P is a transition matrix of a finite homogeneous Markov chain, i.e.*

$$P := (p_{ij})_{i,j=1}^n,$$

where p_{ij} is probability of moving from i to j in one time step. And the matrix P is symmetric, i.e. $P^T = P$, then stationary distribution exists and it is uniformly distributed.

Proof of Lemma 1. Let us look at uniform distribution

$$\pi := \left(\frac{1}{n}, \frac{1}{n}, \dots, \frac{1}{n} \right).$$

We can easily obtain that π is a stationary distribution, using symmetry and stochastic property of matrix P :

$$\pi P = \frac{1}{n} \mathbf{1}^T P = \frac{1}{n} (P \mathbf{1})^T = \frac{1}{n} \mathbf{1}^T = \pi.$$

□

Proof of Theorem 1. We consider states of Markov chain as $s := \{\nu_1, \nu_2, \dots, \nu_K\}_{\nu_1, \dots, \nu_K \in M}$, where M is the set of all subsets of $\overline{1, d}$ of size m . We define $p(s, s', i)$ as the probability to move from state s to state s' for the number of steps i .

• For both compressors $\text{BanLast}(K, m)$ (Definition 5) and $\text{KAWASAKI}(K, b, \pi_\Delta, m)$ (Definition 6) corresponding Markov chain is finite and indecomposable.

The finiteness of the chain is apparent, as the number of states can be explicitly expressed as $|M| = \binom{d}{m}^K$. We show that both chains are indecomposable below. Then we deduce that the chain is ergodic based on the Ergodic Theorem Neumann (1932). Thus, we know that a stationary distribution exists. Then we show that the stationary distribution is uniform over the set of states using Lemma 1.

All that remains is to show that both chains are indecomposable and that transition matrices for both chains are symmetric.

We will start with $\text{BanLast}(K, m)$. Restriction on K, m and d is $d > (K + 1)m$. That makes obvious that any two states are communicated, i.e. for any s, s' there exists way from s to s' . Thus, the Markov chain is indecomposable.

For the compressor probability to move from s to s' in one time step can be explicitly expressed as:

$$p(s, s', 1) = \left(\frac{1}{C_{d-Km}^m} \right)^K,$$

where $C_{d-Km}^m = \frac{(d-Km)!}{m!(d-(K+1)m)!}$ is a binomial coefficient. And all these states are equal in probability. If $d = (K + 1)m$, then for s there will be only one set s' , such that $p(s, s', 1) > 0$, in this case chain will not be ergodic. If $d > (K + 1)m$, then there are more then one state s' , for witch $p(s, s', 1) > 0$, therefore chain will be ergodic.

• According to the Ergodic Theorem, $\rho = (1 - \delta)^{1/N_0}$ and $C = (1 - \delta)^{-1}$, where N_0 is the minimal number of iterations through which is strictly greater then zero and $\delta := \min_{s,s'} \{p(s, s', N_0)\} > 0$. For $\text{BanLast}(K, m)$ in case of $d > (2K + 1)m$ it holds that

$$N_0 = 2 \text{ and } \delta = p(s, s, 2) = \left(\frac{C_{d-2Km}^m}{C_{d-Km}^m} \right)^K \cdot \left(\frac{1}{C_{d-Km}^m} \right)^K,$$

because the smallest probability is to return to state s in two steps.

• For $\text{KAWASAKI}(K, b, \pi_\Delta, m)$ from any given state, there exists a path to any other state in just one iteration, because probabilities to choose any set of coordinates ν are non-zero. Thus, the corresponding markov chain is indecomposable.

We focus on the case where $K = 1$ and that generalize analysis to accommodate larger values of K . Let us look at probabilities to move from ν_i to ν_j and from ν_j to ν_i . We show that both these probabilities correspond to random choice of the same indexes with the same distribution vector p , defined in 6, i.e. the probabilities are equal. For this case let us define ν as operator

$$\Psi_i(\overline{1, d}) := \nu_i,$$

i.e. operator chooses indexes that are in ν_i from $\overline{1, d}$. And

$$\Phi(p, \Psi_i) := \mathbb{P}\{\text{choose } \nu_i \text{ with distribution vector } p\}.$$

According to 6, probability to move from ν_i to ν_j equals a probability to choose indexes ν_j with distribution

$$p_i = \pi_\Delta(\tilde{p}_i),$$

where

$$\tilde{p}_i^k = \begin{cases} 1/bd & \text{if } k \in \nu_i \\ 1/d & \text{if } k \notin \nu_i \end{cases},$$

i.e.

$$p_{ij} = \Phi(p_i, \Psi_j).$$

By the definition of Φ , for arbitrary permutation ϕ and index choice Ψ holds

$$\Phi(\phi(p), \Psi \circ \phi) = \Phi(p, \Psi).$$

Now we point out that for arbitrary ν_i and ν_j exists permutation ϕ_{ij} , such that

$$\Psi_j \circ \phi_{ij} = \Psi_i.$$

For such permutation holds $\phi_{ij}(\tilde{p}_i) = \tilde{p}_j$, i.e. the permutations moves indexes from ν_i to indexes from ν_j . Then we need to use the property of π_Δ to get the same equality for p_i, p_j :

$$\phi_{ij}(p_i) = \phi_{ij}(\pi_\Delta(\tilde{p}_i)) = \pi_\Delta \phi_{ij}(\tilde{p}_i) = \pi_\Delta(p_j).$$

This allows us to write

$$p_{ij} = \Phi(p_i, \Psi_j) = \Phi(\phi_{ij}(p_i), \Psi_j \circ \phi_{ij}) = \Phi(p_j, \Psi_i) = p_{ji}.$$

Thus we get equality of probabilities to move from ν_j to ν_i and to opposite way.

Now we can easily generalize the proof for arbitrary K . All that is required is to consider, instead of the sets of indices ν , combinations of sets of indices that were chosen for transmission over the previous K steps. In this way, the number of states is increased, but the logic of reasoning remains unchanged.

• As was mentioned above, for $\text{KAWASAKI}(K, b, \pi_\Delta, m)$ $N_0 = 1$. We now compute $\delta := p(s, s, 1)$, where $s = \{\nu, \dots, \nu\}$, where ν occurs K times. In this case probability to choose ν another K times is equal to $\mathbb{P}\{j \in \nu\}^{mK}$. And

$$\mathbb{P}\{j \in \nu\} = \min \left\{ \pi_\Delta \left[\tilde{p} := \left(\underbrace{\frac{1/d}{b^K}, \dots, \frac{1/d}{b^K}}_m, \underbrace{\frac{1/d}{1}, \dots, \frac{1/d}{1}}_{d-m} \right)^T \right] \right\}.$$

If we consider $(\pi_\Delta(\tilde{p}))_j = |\tilde{p}_j|/\|\tilde{p}\|_1$, then, since $\|\tilde{p}\|_1 = \frac{1}{db^K}(db^K - m(b^K - 1))$, it hold that $\delta = (db^K - m(b^K - 1))^{-mK}$. This finishes the proof. \square

D MAIN LEMMAS

Lemma 2. For any $i \in \overline{1, n}$, $\varepsilon > 0$, $\tau > \tau_{\text{mix}}(\varepsilon)$, $t > \tau$, for any $a^{t-\tau}, b^{t-\tau} \in \mathbb{R}^d$, such that if we fix all randomness up to step $t - \tau$, $a^{t-\tau}$ and $b^{t-\tau}$ become non-random, it holds that

$$\mathbb{E} [\langle Q_t^i(a^{t-\tau}) - a^{t-\tau}, b^{t-\tau} \rangle] \leq \frac{\varepsilon d}{m} \mathbb{E} [\|a^{t-\tau}\| \cdot \|b^{t-\tau}\|].$$

Proof. We begin by using tower property:

$$\mathbb{E} [\langle Q_t^i(a^{t-\tau}) - a^{t-\tau}, b^{t-\tau} \rangle] = \mathbb{E} [\langle \mathbb{E}_{t-\tau} [Q_t^i(a^{t-\tau}) - a^{t-\tau}], b^{t-\tau} \rangle], \quad (8)$$

where $\mathbb{E}_{t-\tau}[\cdot]$ is the conditional expectation with fixed randomness of all steps up to $t - \tau$. Since on a step t we compress vector $a^{t-\tau}$ according to distribution π_t^i by the formula $Q_t^i(a^{t-\tau}) = d/m a^{t-\tau} \odot \mathbb{1}(\nu_t^i)$, where ν_t^i is some set of m coordinates: $\nu_t^i \subset \overline{1, d}$ and $\mathbb{1}(\nu_t^i)$ is vector with 1 on coordinates ν_t^i on 0 otherwise. Using this we can obtain:

$$\mathbb{E}_{t-\tau} [Q_t^i(a^{t-\tau}) - a^{t-\tau}] = \sum_{\tilde{\nu}_i \in M} \left(\mathbb{P}_{t-\tau} \{ \nu_t^i = \tilde{\nu}_i \} - \frac{1}{C_d^m} \right) a^{t-\tau} \odot \mathbb{1}(\tilde{\nu}_i) \frac{d}{m},$$

where M is set of all subsets of $\overline{1, d}$ of size m . This equality follows from the fact that $\sum_{\tilde{\nu}_i \in M} a^{t-\tau} \odot \mathbb{1}(\tilde{\nu}_i) = C_{d-1}^{m-1} a^{t-\tau}$ and $C_{d-1}^{m-1}/C_d^m = m/d$. Now with the help of Cauchy–Schwarz inequality A.1 we can estimate (8):

$$(8) \leq \mathbb{E} \left[\sum_{\tilde{\nu}_i \in M} \left| \mathbb{P}_{t-\tau} \{ \nu_t^i = \tilde{\nu}_i \} - \frac{1}{C_d^m} \right| \|a^{t-\tau} \odot \mathbb{1}(\tilde{\nu}_i)\| \frac{d}{m} \|b^{t-\tau}\| \right]. \quad (9)$$

Since $t > \tau$ and $\tau > \tau_{\text{mix}}(\varepsilon)$ it holds that $|\mathbb{P}_{t-\tau} \{ \nu_t^i = \tilde{\nu}_i \} - 1/C_d^m| \leq \varepsilon \cdot 1/C_d^m$, because stationary distribution of our Markov chain is uniform. Using the fact that $\|a^{t-\tau} \odot \mathbb{1}(\tilde{\nu}_i)\| \leq \|a^{t-\tau}\|$ we can obtain:

$$(9) \leq \mathbb{E} \left[\sum_{\tilde{\nu}_i \in M} \varepsilon \frac{1}{C_d^m} \|a^{t-\tau}\| \frac{d}{m} \|b^{t-\tau}\| \right] = \frac{\varepsilon d}{m} \mathbb{E} [\|a^{t-\tau}\| \cdot \|b^{t-\tau}\|].$$

This finishes the proof. \square

Lemma 3. For any $i \in \overline{1, n}$, $\varepsilon > 0$, $\tau > \tau_{\text{mix}}(\varepsilon)$, $t > \tau$, for any $a^{t-\tau} \in \mathbb{R}^d$, such that if we fix all randomness up to step $t - \tau$, $a^{t-\tau}$ becomes non-random, it holds that

$$\mathbb{E} \left[\left\| \mathbb{E}_{t-\tau} [Q_t^i(a^{t-\tau})] - a^{t-\tau} \right\|^2 \right] \leq \frac{\varepsilon^2 d^2}{m^2} \mathbb{E} \left[\|a^{t-\tau}\|^2 \right].$$

Proof. Using same notation as in the proof of Lemma 3 we obtain

$$\begin{aligned} \mathbb{E} \left[\left\| \mathbb{E}_{t-\tau} [Q_t^i(a^{t-\tau})] - a^{t-\tau} \right\|^2 \right] &= \mathbb{E} \left[\left\| \sum_{\tilde{\nu}_i \in M} \left(\mathbb{P}_{t-\tau} \{ \nu_t^i = \tilde{\nu}_i \} - \frac{1}{C_d^m} \right) \frac{d}{m} a^{t-\tau} \odot \mathbb{1}(\tilde{\nu}_i) \right\|^2 \right] \\ &\leq \mathbb{E} \left[\frac{d^2}{m^2} C_d^m \sum_{\tilde{\nu}_i \in M} \left(\left| \mathbb{P}_{t-\tau} \{ \nu_t^i = \tilde{\nu}_i \} - \frac{1}{C_d^m} \right|^2 \|a^{t-\tau} \odot \mathbb{1}(\tilde{\nu}_i)\|^2 \right) \right]. \end{aligned}$$

Since $t > \tau$ and $\tau > \tau_{\text{mix}}(\varepsilon)$ it holds that $|\mathbb{P}_{t-\tau} \{ \nu_t^i = \tilde{\nu}_i \} - 1/C_d^m| \leq \varepsilon \cdot 1/C_d^m$, because stationary distribution of our Markov chain is uniform. Using the fact that $\|a^{t-\tau} \odot \mathbb{1}(\tilde{\nu}_i)\| \leq \|a^{t-\tau}\|$ we can obtain:

$$\mathbb{E} \left[\left\| \mathbb{E}_{t-\tau} [Q_t^i(a^{t-\tau})] - a^{t-\tau} \right\|^2 \right] \leq \frac{\varepsilon^2 d^2}{m^2} \mathbb{E} \left[\|a^{t-\tau}\|^2 \right].$$

This finishes the proof. \square

Lemma 4. For any $i \in \overline{1, n}$ and $a \in \mathbb{R}^d$ it holds that

$$\|Q^i(a)\|^2 \leq \frac{d^2}{m^2} \|a\|^2 \quad \text{and} \quad \|Q^i(a) - a\|^2 \leq 4 \frac{d^2}{m^2} \|a\|^2.$$

Proof. Consider the first inequality. Since $Q^i(a) = d/ma \odot \mathbb{1}(\nu^i)$, then $\|Q^i(a)\| \leq d/m \|a\|$, therefore

$$\|Q^i(a)\|^2 \leq \frac{d^2}{m^2} \|a\|^2.$$

Consider the second inequality. Using Fenchel-Young inequality A.2 with $\beta = 1$ we can estimate

$$\|Q^i(a) - a\|^2 \leq 2 \|Q^i(a)\|^2 + 2 \|a\|^2 \leq 2 \left(\frac{d^2}{m^2} + 1 \right) \|a\|^2 \leq 4 \frac{d^2}{m^2} \|a\|^2.$$

This finishes the proof. \square

Corollary 3. For any $i \in \overline{1, n}$, $\varepsilon > 0$, $\tau > \tau_{\text{mix}}(\varepsilon)$, $t > \tau$, for any $a^t, b^t \in \mathbb{R}^d$, such that if we fix all randomness up to step t , a^t and b^t become non-random. And for any $\hat{a}^{t-\tau}, \hat{b}^{t-\tau}$, such that if we fix all randomness up to step $t - \tau$, $\hat{a}^{t-\tau}$ and $\hat{b}^{t-\tau}$ become non-random, it holds that

$$\begin{aligned} 2 |\mathbb{E} [\langle Q_t^i(a^t) - a^t, b^t \rangle]| &\leq \frac{\varepsilon d}{m \beta_0} \mathbb{E} \left[\|\hat{a}^{t-\tau}\|^2 \right] + \frac{\varepsilon d \beta_0}{m} \mathbb{E} \left[\|\hat{b}^{t-\tau}\|^2 \right] + \frac{1}{\beta_2} \mathbb{E} \left[\|b^t\|^2 \right], \\ &+ \left(\frac{1}{\beta_1} + \frac{1}{\beta_3} \right) \mathbb{E} \left[\|b^t - \hat{b}^{t-\tau}\|^2 \right] + 4 \frac{d^2}{m^2} \beta_3 \mathbb{E} \left[\|a^t\|^2 \right] + 4 \frac{d^2 (\beta_1 + \beta_2)}{m^2} \mathbb{E} \left[\|a^t - \hat{a}^{t-\tau}\|^2 \right] \end{aligned}$$

where $\beta_0, \beta_1, \beta_2, \beta_3 > 0$.

1080 *Proof.* Using straightforward algebra we obtain

$$\begin{aligned}
1081 \mathbb{E} [\langle Q_t^i(a^t) - a^t, b^t \rangle] &= \mathbb{E} [\langle Q_t^i(\hat{a}^{t-\tau}) - \hat{a}^{t-\tau}, \hat{b}^{t-\tau} \rangle] \\
1082 &\quad - \mathbb{E} [\langle Q_t^i(a^t - \hat{a}^{t-\tau}) - a^t + \hat{a}^{t-\tau}, b^t - \hat{b}^{t-\tau} \rangle] \\
1083 &\quad + \mathbb{E} [\langle Q_t^i(a^t - \hat{a}^{t-\tau}) - a^t + \hat{a}^{t-\tau}, b^t \rangle] \\
1084 &\quad + \mathbb{E} [\langle Q_t^i(a^t) - a^t, b^t - \hat{b}^{t-\tau} \rangle].
\end{aligned}$$

1085 Using Lemma 2 with $a^{t-\tau} = \hat{a}^{t-\tau}$, $b^{t-\tau} = \hat{b}^{t-\tau}$ and Fenchel-Young inequality A.2 with $\beta_1, \beta_2, \beta_3 > 0$ we obtain:

$$\begin{aligned}
1086 2 |\mathbb{E} [\langle Q_t^i(a^t) - a^t, b^t \rangle]| &\leq 2 \frac{\varepsilon d}{m} \mathbb{E} [\|\hat{a}^{t-\tau}\| \cdot \|\hat{b}^{t-\tau}\|] \\
1087 &\quad + \beta_1 \mathbb{E} [\|Q_t^i(a^t - \hat{a}^{t-\tau}) - a^t + \hat{a}^{t-\tau}\|^2] + \frac{1}{\beta_1} \mathbb{E} [\|b^t - \hat{b}^{t-\tau}\|^2] \\
1088 &\quad + \beta_2 \mathbb{E} [\|Q_t^i(a^t - \hat{a}^{t-\tau}) - a^t + \hat{a}^{t-\tau}\|^2] + \frac{1}{\beta_2} \mathbb{E} [\|b^t\|^2] \\
1089 &\quad + \beta_3 \mathbb{E} [\|Q_t^i(a^t) - a^t\|^2] + \frac{1}{\beta_3} \mathbb{E} [\|b^t - \hat{b}^{t-\tau}\|^2].
\end{aligned}$$

1090 Using Lemma 4 and Fenchel-Young inequality A.2 with $\beta_0 > 0$ we obtain

$$\begin{aligned}
1091 2 |\mathbb{E} [\langle Q_t^i(a^t) - a^t, b^t \rangle]| &\leq \frac{\varepsilon d}{m\beta_0} \mathbb{E} [\|\hat{a}^{t-\tau}\|^2] + \frac{\varepsilon d\beta_0}{m} \mathbb{E} [\|\hat{b}^{t-\tau}\|^2] \\
1092 &\quad + 4 \frac{d^2}{m^2} (\beta_1 + \beta_2) \mathbb{E} [\|a^t - \hat{a}^{t-\tau}\|^2] + \left(\frac{1}{\beta_1} + \frac{1}{\beta_3}\right) \mathbb{E} [\|b^t - \hat{b}^{t-\tau}\|^2] \\
1093 &\quad + 4 \frac{d^2}{m^2} \beta_3 \mathbb{E} [\|a^t\|^2] + \frac{1}{\beta_2} \mathbb{E} [\|b^t\|^2].
\end{aligned}$$

1094 This finishes the proof. \square

1095 **Lemma 5.** Assume 4, then for any $x \in \mathbb{R}^d$ it holds that

$$\frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x)\|^2 \leq 2(\delta^2 + 1) \|\nabla f(x)\|^2 + 2\sigma^2.$$

1096 *Proof.* Using straightforward algebra and Fenchel-Young inequality A.2 with $\beta = 1$ we obtain

$$\begin{aligned}
1097 \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x)\|^2 &\leq \frac{2}{n} \sum_{i=1}^n \|\nabla f_i(x) - \nabla f(x)\|^2 + 2 \|\nabla f(x)\|^2 \\
1098 &\leq 2(\delta^2 + 1) \|\nabla f(x)\|^2 + 2\sigma^2.
\end{aligned}$$

1099 The last inequity follows from 4. This finishes the proof. \square

1100 E EXTENSIONS FOR THEOREM 2

1101 E.1 FULL VERSION OF THEOREM 2

1102 **Theorem 4** (Convergence of MQSGD (Algorithm 1), extension of 2). Consider Assumptions 1, 4 and 5. Let problem (1) be solved by Algorithm 1.

- 1134 • For any $\varepsilon > 0, \gamma > 0, \tau > \tau_{\text{mix}}(\varepsilon)$ and $T > \tau$ satisfying

$$1135 \quad \gamma \lesssim \frac{m^2}{d^2 L (\delta^2 + 1) \tau} \quad \text{and} \quad \varepsilon \lesssim \frac{m^2}{d^2 (\delta^2 + 1)},$$

1136 it holds that

$$1137 \quad \mathbb{E} \left[\|\nabla f(\hat{x}^T)\|^2 \right] = \mathcal{O} \left(\frac{F_\tau}{\gamma T} + \frac{\gamma L \tau d^2}{m^2} \sigma^2 \right),$$

1138 where \hat{x}^T is chosen uniformly from $\{x^t\}_{t=0}^T$.

- 1139 • If f additionally verifies the PL-condition (Assumption 3), then for any $\varepsilon > 0, \gamma > 0, \tau > \tau_{\text{mix}}(\varepsilon)$
1140 and $T > \tau$ satisfying

$$1141 \quad \gamma \lesssim \frac{m^2}{L d^2 \tau (\delta^2 + 1)} \quad \text{and} \quad \varepsilon = \sqrt{\gamma L \tau} \lesssim \frac{m}{d \sqrt{\delta^2 + 1}},$$

1142 it holds that

$$1143 \quad F_T = \mathcal{O} \left(\left(1 - \frac{\mu \gamma}{12} \right)^{T-\tau} F_\tau + \frac{\gamma d^2 L \tau}{\mu m^2} \sigma^2 \right).$$

1144 Here we use a notation $F_t := \mathbb{E} [f(x^t) - f(x^*)]$.

1145 E.2 FULL VERSION OF COROLLARY 1

1146 **Corollary 4** (Step tuning for Theorem 2, extension of Corollary 1).

- 1147 • Under the conditions of Theorem 2 in the non-convex case, choosing γ as

$$1148 \quad \gamma \lesssim \frac{m}{d \sqrt{L \tau}} \min \left\{ \frac{m}{d (\delta^2 + 1) \sqrt{L \tau}} ; \sqrt{\frac{F_\tau}{T \sigma^2}} \right\},$$

1149 in order to achieve ε -approximate solution (in terms of $\mathbb{E} [\|\nabla f(x^T)\|^2] \leq \varepsilon^2$) it takes

$$1150 \quad \mathcal{O} \left(\frac{L \tau d^2}{m^2} F_\tau \left(\frac{\delta^2 + 1}{\varepsilon^2} + \frac{\sigma^2}{\varepsilon^4} \right) \right) \text{ iterations of Algorithm 1.}$$

- 1151 • Under the conditions of Theorem 2 in the PL-condition (Assumption 3) case, choosing γ as

$$1152 \quad \gamma \lesssim \min \left\{ \frac{m^2}{L d^2 \tau (\delta^2 + 1)} ; \frac{\log \left(\max \left\{ 2 ; \frac{\mu^2 m^2 F_\tau T}{d^2 L \tau \sigma^2} \right\} \right)}{\mu T} \right\},$$

1153 in order to achieve ε -approximate solution (in terms of $\mathbb{E} [f(x^t) - f(x^*)] \leq \varepsilon$) it takes

$$1154 \quad \mathcal{O} \left(\frac{d^2 L \tau}{m^2 \mu} \left((\delta^2 + 1) \log \left(\frac{1}{\varepsilon} \right) + \frac{\sigma^2}{\mu \varepsilon} \right) \right) \text{ iterations of Algorithm 1.}$$

1155 E.3 PROOF OF THEOREM 2, NON-CONVEX CASE

1156 *Proof.* Denoting $F_t := \mathbb{E} [f(x^t) - f(x^*)]$, we have using L -smoothness:

$$1157 \quad F_{t+1} - F_t \leq -\gamma \mathbb{E} \left[\left\langle \frac{1}{n} \sum_{i=1}^n Q_i^i(\nabla f_i(x^t)), \nabla f(x^t) \right\rangle \right] + \frac{\gamma^2 L}{2} \mathbb{E} \left[\left\| \frac{1}{n} \sum_{i=1}^n Q_i^i(\nabla f_i(x^t)) \right\|^2 \right]. \quad (10)$$

1188 Consider $-\gamma \mathbb{E} \left[\left\langle \frac{1}{n} \sum_{i=1}^n Q_t^i(\nabla f_i(x^t)), \nabla f(x^t) \right\rangle \right]$. Using straightforward algebra: $\pm \nabla f_i(x^{t-\tau})$ and
 1189 $\pm \nabla f(x^{t-\tau})$ we can re-write this term:
 1190
 1191
 1192
 1193

$$\begin{aligned}
 & -\gamma \mathbb{E} \left[\left\langle \frac{1}{n} \sum_{i=1}^n Q_t^i(\nabla f_i(x^t)), \nabla f(x^t) \right\rangle \right] \\
 & = -\gamma \mathbb{E} \left[\underbrace{\left\langle \frac{1}{n} \sum_{i=1}^n Q_t^i(\nabla f_i(x^{t-\tau})), \nabla f(x^{t-\tau}) \right\rangle}_{\textcircled{1}} \right] \\
 & \quad - \gamma \mathbb{E} \left[\underbrace{\left\langle \frac{1}{n} \sum_{i=1}^n Q_t^i(\nabla f_i(x^t)), \nabla f(x^t) - \nabla f(x^{t-\tau}) \right\rangle}_{\textcircled{2}} \right] \\
 & \quad - \gamma \mathbb{E} \left[\underbrace{\left\langle \frac{1}{n} \sum_{i=1}^n Q_t^i(\nabla f_i(x^t) - \nabla f_i(x^{t-\tau})), \nabla f(x^{t-\tau}) \right\rangle}_{\textcircled{3}} \right].
 \end{aligned}$$

1209
 1210
 1211 Consider $\textcircled{1}$. Using straightforward algebra, tower property, Lemmas 3 and 5 we obtain
 1212
 1213
 1214
 1215

$$\begin{aligned}
 \textcircled{1} & = -\gamma \mathbb{E} \left[\left\langle \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{t-\tau} [Q_t^i(\nabla f_i(x^{t-\tau}))], \nabla f(x^{t-\tau}) \right\rangle \right] \\
 & = -\frac{\gamma}{2} \mathbb{E} \left[\left\| \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{t-\tau} [Q_t^i(\nabla f_i(x^{t-\tau}))] \right\|^2 \right] \\
 & \quad + \frac{\gamma}{2} \mathbb{E} \left[\left\| \nabla f(x^{t-\tau}) - \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{t-\tau} [Q_t^i(\nabla f_i(x^{t-\tau}))] \right\|^2 \right] - \frac{\gamma}{2} \mathbb{E} \left[\|\nabla f(x^{t-\tau})\|^2 \right] \quad (11) \\
 & \leq \frac{\gamma}{2} \varepsilon^2 \frac{d^2}{m^2} \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left[\|\nabla f_i(x^{t-\tau})\|^2 \right] - \frac{\gamma}{2} \mathbb{E} \left[\|\nabla f(x^{t-\tau})\|^2 \right] \\
 & \leq \gamma \left(\varepsilon^2 \frac{d^2}{m^2} (\delta^2 + 1) - \frac{1}{2} \right) \mathbb{E} \left[\|\nabla f(x^{t-\tau})\|^2 \right] + \gamma \varepsilon^2 \frac{d^2}{m^2} \sigma^2 \\
 & \leq -\frac{\gamma}{4} \mathbb{E} \left[\|\nabla f(x^{t-\tau})\|^2 \right] + \gamma \varepsilon^2 \frac{d^2}{m^2} \sigma^2.
 \end{aligned}$$

1233
 1234 The last inequality follows from the fact, that
 1235
 1236

$$\varepsilon \leq \frac{m}{2d\sqrt{\delta^2 + 1}}.$$

1237
 1238
 1239
 1240
 1241 Consider $\textcircled{2}$. Using Cauchy-Schwarz A.1 and Fenchel-Young A.2 with $\beta = 1$ inequalities we obtain

$$\begin{aligned}
& \textcircled{2} \leq \mathbb{E} \left[\left\| -\frac{\gamma}{n} \sum_{i=1}^n Q_t^i(\nabla f_i(x^t)) \right\| \|\nabla f(x^t) - \nabla f(x^{t-\tau})\| \right] \\
& \leq \gamma L \mathbb{E} \left[\left\| \frac{1}{n} \sum_{i=1}^n Q_t^i(\nabla f_i(x^t)) \right\| \|x^t - x^{t-\tau}\| \right] \\
& = \gamma^2 L \mathbb{E} \left[\left\| \frac{1}{n} \sum_{i=1}^n Q_t^i(\nabla f_i(x^t)) \right\| \left\| \sum_{s=t-\tau}^{t-1} \frac{1}{n} \sum_{i=1}^n Q_s^i(\nabla f_i(x^s)) \right\| \right] \\
& \leq \frac{\gamma^2 L}{2} \left(\tau \mathbb{E} \left[\left\| \frac{1}{n} \sum_{i=1}^n Q_t^i(\nabla f_i(x^t)) \right\|^2 \right] + \sum_{s=t-\tau}^{t-1} \mathbb{E} \left[\left\| \frac{1}{n} \sum_{i=1}^n Q_s^i(\nabla f_i(x^s)) \right\|^2 \right] \right).
\end{aligned} \tag{12}$$

Third equality holds since $x^t - x^{t-\tau} = \gamma \sum_{s=t-\tau}^{t-1} \frac{1}{n} \sum_{i=1}^n Q_s^i(\nabla f_i(x^s))$. Consider $\textcircled{3}$. Using Cauchy-Schwarz A.1 and Fenchel-Young A.2 with $\beta = m/d$ inequalities we obtain

$$\begin{aligned}
& \textcircled{3} \leq \mathbb{E} \left[\left\| -\frac{\gamma}{n} \sum_{i=1}^n Q_t^i(\nabla f_i(x^t) - \nabla f_i(x^{t-\tau})) \right\| \|\nabla f(x^{t-\tau})\| \right] \\
& \leq \gamma L \mathbb{E} \left[\left\| \frac{1}{n} \sum_{i=1}^n Q_t^i(\nabla f_i(x^t) - \nabla f_i(x^{t-\tau})) \right\| \|\nabla f(x^{t-\tau})\| \right] \\
& \leq \gamma^2 L \frac{d}{m} \mathbb{E} \left[\left\| \frac{1}{n} \sum_{i=1}^n Q_t^i(\nabla f_i(x^t)) \right\| \left\| \sum_{s=t-\tau}^{t-1} \frac{1}{n} \sum_{i=1}^n Q_s^i(\nabla f_i(x^s)) \right\| \right] \\
& \leq \frac{\gamma^2 L}{2} \left(\sum_{s=t-\tau}^{t-1} \mathbb{E} \left[\left\| \frac{1}{n} \sum_{i=1}^n Q_s^i(\nabla f_i(x^s)) \right\|^2 \right] + \frac{d^2 \tau}{m^2} \mathbb{E} [\|\nabla f(x^{t-\tau})\|^2] \right).
\end{aligned} \tag{13}$$

Wrapping (10) - (13) up we obtain

$$\begin{aligned}
F_{t+1} - F_t & \leq \frac{\gamma^2 L}{2} \mathbb{E} \left[\left\| \frac{1}{n} \sum_{i=1}^n Q_t^i(\nabla f_i(x^t)) \right\|^2 \right] - \frac{\gamma}{4} \mathbb{E} [\|\nabla f(x^{t-\tau})\|^2] + \gamma \varepsilon^2 \frac{d^2}{m^2} \sigma^2 \\
& \quad + \frac{\gamma^2 L}{2} \left(\tau \mathbb{E} \left[\left\| \frac{1}{n} \sum_{i=1}^n Q_t^i(\nabla f_i(x^t)) \right\|^2 \right] + \sum_{s=t-\tau}^{t-1} \mathbb{E} \left[\left\| \frac{1}{n} \sum_{i=1}^n Q_s^i(\nabla f_i(x^s)) \right\|^2 \right] \right) \\
& \quad + \frac{\gamma^2 L}{2} \left(\sum_{s=t-\tau}^{t-1} \mathbb{E} \left[\left\| \frac{1}{n} \sum_{i=1}^n Q_s^i(\nabla f_i(x^s)) \right\|^2 \right] + \frac{d^2 \tau}{m^2} \mathbb{E} [\|\nabla f(x^{t-\tau})\|^2] \right) \\
& \leq \gamma^2 L \tau \mathbb{E} \left[\left\| \frac{1}{n} \sum_{i=1}^n Q_t^i(\nabla f_i(x^t)) \right\|^2 \right] + \gamma \varepsilon^2 \frac{d^2}{m^2} \sigma^2 \\
& \quad + \gamma^2 L \sum_{s=t-\tau}^{t-1} \mathbb{E} \left[\left\| \frac{1}{n} \sum_{i=1}^n Q_s^i(\nabla f_i(x^s)) \right\|^2 \right] + \left(\frac{\gamma^2 L \tau d^2}{2m^2} - \frac{\gamma}{4} \right) \mathbb{E} [\|\nabla f(x^{t-\tau})\|^2].
\end{aligned}$$

Using Lemma 5 we obtain

$$\begin{aligned}
F_{t+1} - F_t &\leq \frac{2d^2\gamma^2L\tau}{m^2} \left((\delta^2 + 1)\mathbb{E} \left[\|\nabla f(x^t)\|^2 \right] + \sigma^2 \right) + \left(\frac{\gamma^2L\tau d^2}{2m^2} - \frac{\gamma}{4} \right) \mathbb{E} \left[\|\nabla f(x^{t-\tau})\|^2 \right] \\
&\quad + \frac{2d^2\gamma^2L}{m^2} \sum_{s=t-\tau}^{t-1} \left((\delta^2 + 1)\mathbb{E} \left[\|\nabla f(x^s)\|^2 \right] + \sigma^2 \right) + \gamma\varepsilon^2 \frac{d^2}{m^2} \sigma^2 \\
&= \frac{2d^2\gamma^2L(\delta^2 + 1)\tau}{m^2} \mathbb{E} \left[\|\nabla f(x^t)\|^2 \right] + \frac{2d^2\gamma^2L(\delta^2 + 1)}{m^2} \sum_{s=t-\tau}^{t-1} \mathbb{E} \left[\|\nabla f(x^s)\|^2 \right] \\
&\quad + \left(\frac{\gamma^2L\tau d^2}{2m^2} - \frac{\gamma}{4} \right) \mathbb{E} \left[\|\nabla f(x^{t-\tau})\|^2 \right] + \frac{\gamma d^2}{m^2} (4\gamma L\tau + \varepsilon^2) \sigma^2.
\end{aligned} \tag{14}$$

Summing (14) from $t = \tau$ to $t = T$ and using the fact that $\varepsilon^2 \leq \gamma L\tau$ and $1 + \delta^2 \geq 1$ we obtain

$$\begin{aligned}
\sum_{t=\tau}^T \frac{\gamma}{4} \mathbb{E} \left[\|\nabla f(x^{t-\tau})\|^2 \right] &\leq F_\tau + \frac{2d^2\gamma^2L(\delta^2 + 1)}{m^2} \left(\tau \sum_{t=\tau}^T \mathbb{E} \left[\|\nabla f(x^t)\|^2 \right] \right. \\
&\quad \left. + \sum_{t=\tau}^T \sum_{s=t-\tau}^{t-1} \mathbb{E} \left[\|\nabla f(x^s)\|^2 \right] + \tau \sum_{t=\tau}^T \mathbb{E} \left[\|\nabla f(x^{t-\tau})\|^2 \right] \right) + \sum_{t=\tau}^T 5 \frac{\gamma^2L\tau d^2}{m^2} \sigma^2.
\end{aligned}$$

Since $\sum_{t=\tau}^T \sum_{s=t-\tau}^{t-1} \mathbb{E} \left[\|\nabla f(x^s)\|^2 \right] \leq \tau \sum_{t=0}^T \mathbb{E} \left[\|\nabla f(x^t)\|^2 \right]$, we get

$$\gamma \sum_{t=0}^{T-\tau} \mathbb{E} \left[\|\nabla f(x^t)\|^2 \right] \leq 4F_\tau + \frac{24d^2\gamma^2L(\delta^2 + 1)\tau}{m^2} \sum_{t=0}^T \mathbb{E} \left[\|\nabla f(x^t)\|^2 \right] + 20 \sum_{t=\tau}^T \frac{\gamma^2L\tau d^2}{m^2} \sigma^2.$$

Taking

$$\gamma \leq \frac{m^2}{48d^2L(\delta^2 + 1)\tau},$$

we obtain

$$\gamma \sum_{t=0}^{T-\tau} \mathbb{E} \left[\|\nabla f(x^t)\|^2 \right] \leq 8F_\tau + \frac{48d^2\gamma^2L(\delta^2 + 1)\tau}{m^2} \sum_{t=T-\tau}^T \mathbb{E} \left[\|\nabla f(x^t)\|^2 \right] + 40 \sum_{t=\tau}^T \frac{\gamma^2L\tau d^2}{m^2} \sigma^2. \tag{15}$$

We now prove that for any $t \geq 0$, we have

$$\sup_{t \leq s \leq t+\tau} \left\{ \mathbb{E} \left[\|\nabla f(x^s)\|^2 \right] \right\} \leq 4\mathbb{E} \left[\|\nabla f(x^t)\|^2 \right] + 8L^2\gamma^2\tau^2 \frac{d^2}{m^2} \sigma^2.$$

For $t \leq s \leq t + \tau$ it holds that

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365

$$\begin{aligned}
\mathbb{E} \left[\|\nabla f(x^s)\|^2 \right] &\leq 2\mathbb{E} \left[\|\nabla f(x^t)\|^2 \right] + \mathbb{E} \left[\|\nabla f(x^s) - \nabla f(x^t)\|^2 \right] \\
&\leq 2\mathbb{E} \left[\|\nabla f(x^t)\|^2 \right] + 2L^2\gamma^2\mathbb{E} \left[\left\| \sum_{r=t}^{s-1} \frac{1}{n} \sum_{i=1}^n Q_r^i(\nabla f_i(x^r)) \right\|^2 \right] \\
&\leq 2\mathbb{E} \left[\|\nabla f(x^t)\|^2 \right] + 2L^2\gamma^2\tau \frac{d^2}{m^2} \sum_{r=t}^{s-1} \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left[\|\nabla f_i(x^r)\|^2 \right] \\
&\leq 2\mathbb{E} \left[\|\nabla f(x^t)\|^2 \right] + 4L^2\gamma^2\tau \frac{d^2}{m^2} \sum_{r=t}^{s-1} \left((\delta^2 + 1)\mathbb{E} \left[\|\nabla f(x^r)\|^2 \right] + \sigma^2 \right) \\
&\leq 2\mathbb{E} \left[\|\nabla f(x^t)\|^2 \right] + 4L^2\gamma^2\tau^2 \frac{d^2}{m^2} \left((\delta^2 + 1) \sup_{t \leq s \leq t+\tau} \left\{ \mathbb{E} \left[\|\nabla f(x^s)\|^2 \right] \right\} + \sigma^2 \right).
\end{aligned}$$

1366 Since

$$1367 \quad \gamma \leq \frac{m}{\sqrt{8dL\sqrt{\delta^2 + 1}\tau}},$$

1368 it holds that

$$1369 \quad \sup_{t \leq s \leq t+\tau} \left\{ \mathbb{E} \left[\|\nabla f(x^s)\|^2 \right] \right\} \leq 4\mathbb{E} \left[\|\nabla f(x^t)\|^2 \right] + 8L^2\gamma^2\tau^2 \frac{d^2}{m^2}\sigma^2.$$

1370 Using this (15) takes form

$$\begin{aligned}
1371 \quad \gamma \sum_{t=0}^{T-\tau} \mathbb{E} \left[\|\nabla f(x^t)\|^2 \right] &\leq 8F_\tau + \frac{192d^2\gamma^2L(\delta^2 + 1)\tau}{m^2} \sum_{t=T-2\tau}^{T-\tau} \mathbb{E} \left[\|\nabla f(x^t)\|^2 \right] \\
1372 &+ 384L^3\gamma^4\tau^3 \frac{d^4}{m^4}(\delta^2 + 1)\sigma^2 + 40 \sum_{t=\tau}^T \frac{\gamma^2L\tau d^2}{m^2}\sigma^2.
\end{aligned}$$

1373 Taking

$$1374 \quad \gamma \leq \frac{m}{384dL\sqrt{\delta^2 + 1}\tau},$$

1375 and dividing both sides of the inequality by $T - \tau$, we obtain

$$1376 \quad \frac{1}{T - \tau} \sum_{t=0}^{T-\tau} \mathbb{E} \left[\|\nabla f(x^t)\|^2 \right] \leq 16 \frac{F_\tau}{\gamma(T - \tau)} + 80 \frac{\gamma^2L\tau d^2}{m^2}\sigma^2.$$

1377 Therefore, if \hat{x}^T is chosen uniformly from $\{x^t\}_{t=0}^{T-1}$, then it holds that

$$1378 \quad \mathbb{E} \left[\|\nabla f(\hat{x}^T)\|^2 \right] \leq 16 \frac{F_\tau}{\gamma T} + 80 \frac{\gamma^2L\tau d^2}{m^2}\sigma^2.$$

1379 This finishes the proof. □

1400
1401
1402
1403

E.4 PROOF OF THEOREM 2, UNDER PL-CONDITION

Proof. We start from (14):

$$\begin{aligned}
F_{t+1} - F_t &= \frac{2d^2\gamma^2L(\delta^2 + 1)\tau}{m^2} \mathbb{E} \left[\|\nabla f(x^t)\|^2 \right] + \frac{2d^2\gamma^2L(\delta^2 + 1)}{m^2} \sum_{s=t-\tau}^{t-1} \mathbb{E} \left[\|\nabla f(x^s)\|^2 \right] \\
&\quad + \left(\frac{\gamma^2L\tau d^2}{2m^2} - \frac{\gamma}{4} \right) \mathbb{E} \left[\|\nabla f(x^{t-\tau})\|^2 \right] + \frac{\gamma d^2}{m^2} (4\gamma L\tau + \varepsilon^2) \sigma^2.
\end{aligned}$$

If f satisfies PL-inequality (Assumption 3), then $-\mathbb{E} \left[\|\nabla f(x^{t-\tau})\|^2 \right] \leq -2\mu F_{t-\tau}$, so that, for some $0 < \alpha < 1$ we obtain

$$\begin{aligned}
F_{t+1} - F_t &= \frac{2d^2\gamma^2L(\delta^2 + 1)\tau}{m^2} \mathbb{E} \left[\|\nabla f(x^t)\|^2 \right] + \frac{2d^2\gamma^2L(\delta^2 + 1)}{m^2} \sum_{s=t-\tau}^{t-1} \mathbb{E} \left[\|\nabla f(x^s)\|^2 \right] \\
&\quad + \left(\frac{\gamma^2L\tau d^2}{2m^2} - \frac{(1-\alpha)\gamma}{4} \right) \mathbb{E} \left[\|\nabla f(x^{t-\tau})\|^2 \right] \\
&\quad - \frac{\alpha\gamma\mu}{2} F_{t-\tau} + \frac{\gamma d^2}{m^2} (4\gamma L\tau + \varepsilon^2) \sigma^2.
\end{aligned} \tag{16}$$

For $t \geq 0$, let $p_t = p^t$ and $p = (1 - \alpha\mu\gamma/4)^{-1}$. We multiply the above expression by p_t and sum for $t < T$, hoping for cancellations. Using PL-condition (Assumption 3), for $T \geq \tau$ we obtain

$$\begin{aligned}
\sum_{t=\tau}^{T-1} p_{t+1} \left(F_t - F_{t+1} - \frac{\alpha\gamma\mu}{4} F_{t-\tau} \right) &= \sum_{t=\tau}^{T-1} p_{t+1} \left[\left(1 - \frac{\alpha\gamma\mu}{4} \right) F_t - F_{t+1} + \frac{\alpha\gamma\mu}{4} (F_t - F_{t-\tau}) \right] \\
&= \sum_{t=\tau}^{T-1} p_t F_t - \sum_{t=\tau+1}^T p_t F_t + \frac{\alpha\gamma\mu}{4} \sum_{t=\tau}^{T-1} p_{t+1} (F_t - F_{t-\tau}) \\
&\leq p_\tau F_\tau - p_T F_T + \frac{\alpha\gamma\mu}{4} \sum_{t=\tau}^{T-1} p_{t+1} F_t \\
&\quad - \frac{\alpha\gamma\mu p_\tau}{4} \sum_{t=0}^{T-1-\tau} p_{t+1} F_t \\
&\leq p_\tau F_\tau - p_T F_T + \frac{\alpha\gamma\mu}{4} \sum_{t=T-\tau}^{T-1} p_{t+1} F_t \\
&\leq p_\tau F_\tau - p_T F_T + \frac{\alpha\gamma}{8} \sum_{t=T-\tau}^{T-1} p_{t+1} \mathbb{E} \left[\|\nabla f(x^t)\|^2 \right].
\end{aligned}$$

For any $t \geq 0$ we use a notation $b_t := \mathbb{E} \left[\|\nabla f(x^t)\|^2 \right]$. We now handle b_t terms from (16).

$$-\sum_{t=\tau}^{T-1} \frac{(1-\alpha)\gamma}{4} p_{t+1} b_{t-\tau} + \gamma^2 L \frac{d^2}{m^2} \sum_{t=\tau}^{T-1} p_{t+1} \left(2\tau(\delta^2 + 1)b_t + 2(\delta^2 + 1) \sum_{s=t-\tau}^{t-1} b_s + \frac{\tau}{2} b_{t-\tau} \right). \tag{17}$$

If $p_t = p^t$, $p = (1 - \alpha\mu\gamma/2)^{-1}$ and $\gamma = \gamma_1/\tau$, then, using the fact that $(1 - a/x)^{-x} \leq 2e^a \leq 2e$ if $x \geq 2$ and $0 \leq a \leq 1$, we can get that $1 \geq p_\tau = (1 - \mu\gamma_1/(2\tau))^{-\tau} \leq 2e^{\mu\gamma_1/2} \leq 2e \leq 6$. Then

$$\sum_{t=\tau}^T p_{t+1} \sum_{s=t-\tau}^{t-1} b_s \leq p^\tau \sum_{t=\tau}^T \sum_{s=t-\tau}^{t-1} p_{s+1} b_s \leq 6\tau \sum_{t=0}^T p_{t+1} b_t.$$

1458 Now we can estimate (17):
 1459
 1460

$$\begin{aligned}
 1461 \quad (17) &\leq - \sum_{t=0}^{T-\tau-1} \frac{(1-\alpha)\gamma}{4} p_{t+1} b_t + \gamma^2 L \frac{d^2 \tau}{m^2} \left(2(\delta^2 + 1) \sum_{t=\tau}^{T-1} b_t + 12(\delta^2 + 1) \sum_{t=0}^{T-1} b_t + 3 \sum_{t=0}^{T-\tau} b_t \right) \\
 1462 &\leq - \sum_{t=0}^{T-\tau-1} p_{t+1} \gamma b_t \left(\frac{1-\alpha}{4} - 17\gamma L \frac{d^2 \tau (\delta^2 + 1)}{m^2} \right) + 14\gamma^2 L \frac{d^2 \tau (\delta^2 + 1)}{m^2} \sum_{t=T-\tau}^{T-1} p_{t+1} b_t.
 \end{aligned}$$

1467 Taking

$$1468 \quad \gamma \leq \frac{m^2(1-\alpha)}{136Ld^2\tau(\delta^2+1)\beta},$$

1469 where $\beta \geq 1$, we obtain
 1470
 1471

$$1472 \quad (17) \leq -\frac{(1-\alpha)\gamma}{8} \sum_{t=0}^{T-\tau-1} p_{t+1} b_t + \frac{(1-\alpha)\gamma}{4\beta} \sum_{t=T-\tau}^{T-1} p_{t+1} b_t.$$

1473 Now we can estimate (16):
 1474
 1475

$$\begin{aligned}
 1476 \quad 0 &\leq p_\tau F_\tau - p_T F_T + \left(\frac{\alpha\gamma}{8} + \frac{(1-\alpha)\gamma}{4\beta} \right) \sum_{t=T-\tau}^{T-1} p_{t+1} b_t - \frac{(1-\alpha)\gamma}{8} \sum_{t=0}^{T-\tau-1} p_{t+1} b_t \\
 1477 &\quad + \sum_{t=\tau}^{T-1} p_{t+1} \frac{\gamma d^2}{m^2} (4\gamma L\tau + \varepsilon^2) \sigma^2.
 \end{aligned} \tag{18}$$

1478 Using that we proved in E.3 we have $b_t \leq 4b_{t-\tau} + 8L^2\gamma^2\tau^2 \frac{d^2}{m^2} \sigma^2$. Then, we can obtain
 1479
 1480

$$\begin{aligned}
 1481 \quad \gamma \left(\frac{\alpha}{8} + \frac{1-\alpha}{4\beta} \right) \sum_{t=T-\tau}^{T-1} p_{t+1} b_t &\leq 24\gamma \left(\frac{\alpha}{8} + \frac{1-\alpha}{4\beta} \right) \sum_{t=T-2\tau}^{T-\tau-1} p_{t+1} b_t \\
 1482 &\quad + 48L^2\gamma^3\tau^3 \frac{d^2}{m^2} \left(\frac{\alpha}{8} + \frac{1-\alpha}{4\beta} \right) \sigma^2.
 \end{aligned}$$

1483 Taking $\alpha = 1/6$ and $\beta = 4$, we obtain
 1484
 1485

$$1486 \quad \frac{\alpha}{8} + \frac{1-\alpha}{4\beta} = \frac{1-\alpha}{8},$$

1487 and (18) takes form
 1488
 1489

$$1490 \quad 0 \leq p_\tau F_\tau - p_T F_T + 48L^2\gamma^3\tau^3 \frac{d^2}{m^2} \sigma^2 + \sum_{t=\tau}^{T-1} p_{t+1} \frac{\gamma d^2}{m^2} (4\gamma L\tau + \varepsilon^2) \sigma^2. \tag{19}$$

1491 Using the fact that
 1492
 1493

$$1494 \quad \sum_{t=\tau}^T \left(1 - \frac{\alpha\mu\gamma}{2}\right)^{T-t} = \sum_{t=0}^{T-\tau} \left(1 - \frac{\alpha\mu\gamma}{2}\right)^t \leq \sum_{t=0}^{+\infty} \left(1 - \frac{\alpha\mu\gamma}{2}\right)^t = \frac{2}{\alpha\mu\gamma},$$

1495 and taking
 1496
 1497

1512

1513

1514

1515

$$\gamma \leq \frac{m^2}{625Ld^2\tau(\delta^2 + 1)} \quad \text{and} \quad \varepsilon = \sqrt{\gamma L\tau} \leq \frac{m}{25d\sqrt{\delta^2 + 1}},$$

1516

by dividing (19) by p_τ , we obtain

1517

1518

1519

1520

1521

$$\mathbb{E} [f(x^T) - f(x^*)] \leq \left(1 - \frac{\mu\gamma}{12}\right)^{T-\tau} \mathbb{E} [f(x^\tau) - f(x^*)] + 636 \frac{\gamma d^2 L\tau}{\mu m^2} \sigma^2.$$

1522

This finishes the proof. □

1523

1524

1525

1526

F CONVERGENCE OF ALGORITHM 1 WITHOUT DATA SIMILARITY

1527

1528

1529

1530

Theorem 5 (Convergence of GD Algorithm 1 without data similarity). *Consider Assumptions 1 and 2. Let problem (1) be solved by Algorithm 1. Then for any $\varepsilon > 0$, $\gamma > 0$, $\tau > \tau_{\text{mix}}(\varepsilon)$ and $T > \tau$ satisfying*

1531

1532

1533

1534

$$\gamma \leq \frac{m^2 \sqrt{\mu}}{24d^2 L^{3/2} \tau} \quad \text{and} \quad \varepsilon \leq \frac{m\sqrt{\mu}}{24d} \min \left\{ \frac{1}{L^{3/2}}; \sqrt{\mu} \right\},$$

1535

it holds that

1536

1537

1538

1539

$$\mathbb{E} [\|x^{T+1} - x^*\|^2] \leq \left(1 - \frac{\mu\gamma}{2}\right)^{T-\tau} \mathbb{E} [\|x^\tau - x^*\|^2] + \left(1 - \frac{\mu\gamma}{2}\right)^T \Delta_\tau + 26 \frac{\gamma d^2 \tau}{\mu m^2} \sigma_*^2,$$

1540

where

1541

1542

1543

1544

$$\Delta_\tau = \mathcal{O} \left(\frac{\gamma^2 d^2}{m^2} \sqrt{\frac{\mu}{L}} \sum_{t=0}^{\tau} \left[\tau \mathbb{E} [\|x^t - x^*\|^2] + 4L \mathbb{E} [f(x^t) - f(x^*)] \right] \right).$$

1545

Proof of Theorem 5. We start by writing out step of the Algorithm 1:

1546

1547

1548

1549

1550

1551

1552

1553

1554

Consider $\mathbb{E} [\langle Q_t^i (\nabla f_i(x^t)) - \nabla f_i(x^t), x^t - x^* \rangle]$. Using Corollary 3 with $a^t = \nabla f_i(x^t)$, $b^t = x^t - x^*$, $\hat{a}^{t-\tau} = \nabla f_i(x^{t-\tau})$ and $\hat{b}^{t-\tau} = x^{t-\tau} - x^*$ we obtain

1558

1559

1560

1561

1562

1563

1564

1565

$$\begin{aligned} 2 \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n |\langle Q_t^i (\nabla f_i(x^t)) - \nabla f_i(x^t), x^t - x^* \rangle| \right] &\leq \frac{\varepsilon d}{m\beta_0} \frac{1}{n} \sum_{i=1}^n \mathbb{E} [\|\nabla f_i(x^{t-\tau})\|^2] \\ &+ \frac{\varepsilon d\beta_0}{m} \mathbb{E} [\|x^{t-\tau} - x^*\|^2] + 4 \frac{d^2 L^2}{m^2} (\beta_1 + \beta_2) \mathbb{E} [\|x^t - x^\tau\|^2] + \left(\frac{1}{\beta_1} + \frac{1}{\beta_3} \right) \mathbb{E} [\|x^t - x^\tau\|^2] \\ &+ 4 \frac{d^2}{m^2} \beta_3 \frac{1}{n} \sum_{i=1}^d \mathbb{E} [\|\nabla f_i(x^t)\|^2] + \frac{1}{\beta_2} \mathbb{E} [\|x^t - x^*\|^2]. \end{aligned}$$

(21)

Using the fact that f_i are L -smooth, we can obtain:

$$\begin{aligned}
\frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^t)\|^2 &= \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^t) - \nabla f_i(x^*) + \nabla f_i(x^*)\|^2 \\
&\leq \frac{2}{n} \sum_{i=1}^n \|\nabla f_i(x^t) - \nabla f_i(x^*)\|^2 + \frac{2}{n} \sum_{i=1}^n \|\nabla f_i(x^*)\|^2 \\
&\leq \frac{4L}{n} \sum_{i=1}^n (f_i(x^t) - f_i(x^*) - \langle \nabla f_i(x^*), x^t - x^* \rangle) + 2\sigma_*^2 \\
&= 4L(f(x^t) - f(x^*)) + 2\sigma_*^2,
\end{aligned} \tag{22}$$

where we use a notation $\sigma_*^2 := \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^*)\|^2$. Now we can estimate (21):

$$\begin{aligned}
(21) &\leq \frac{2\epsilon d}{m\beta_0} (2L\mathbb{E}[f(x^{t-\tau}) - f(x^*)] + \sigma_*^2) + \frac{\epsilon d\beta_0}{m} \mathbb{E}[\|x^{t-\tau} - x^*\|^2] \\
&\quad + \left(4\frac{d^2L^2}{m^2}(\beta_1 + \beta_2) + \frac{1}{\beta_1} + \frac{1}{\beta_3}\right) \mathbb{E}\left[\left\|-\gamma \sum_{s=t-\tau}^{t-1} \frac{1}{n} \sum_{i=1}^n Q_s^i(\nabla f_i(x^s))\right\|^2\right] \\
&\quad + 8\frac{d^2}{m^2}\beta_3(2L\mathbb{E}[f(x^t) - f(x^*)] + \sigma_*^2) + \frac{1}{\beta_2} \mathbb{E}[\|x^t - x^*\|^2].
\end{aligned}$$

Now we can estimate (20). Using Lemma 4 and Assumption 2 we can obtain

$$\begin{aligned}
\mathbb{E}[\|x^{t+1} - x^*\|^2] &\leq \left(1 - \mu\gamma + \frac{\gamma}{\beta_2}\right) \mathbb{E}[\|x^t - x^*\|^2] + \frac{\epsilon d\beta_0\gamma}{m} \mathbb{E}[\|x^{t-\tau} - x^*\|^2] \\
&\quad + 4L\mathbb{E}\left[\frac{\epsilon d\gamma}{m\beta_0}(f(x^{t-\tau}) - f(x^*)) + 4\frac{d^2\beta_3\gamma}{m^2}(f(x^t) - f(x^*))\right. \\
&\quad + \left.4\left(\frac{d^2L^2}{m^2}(\beta_1 + \beta_2) + \frac{1}{\beta_1} + \frac{1}{\beta_3}\right) \frac{\gamma^3\tau d^2}{m^2} \sum_{s=t-\tau}^{t-1} (f(x^s) - f(x^*))\right. \\
&\quad + \left.\frac{\gamma^2 d^2}{m^2}(f(x^t) - f(x^*)) - \frac{\gamma}{2L}(f(x^t) - f(x^*))\right] \\
&\quad + 2\left[\frac{\epsilon d\gamma}{m\beta_0} + 4\frac{d^2\beta_3\gamma}{m^2} + \left(4\frac{d^2L^2}{m^2}(\beta_1 + \beta_2) + \frac{1}{\beta_1} + \frac{1}{\beta_3}\right) \frac{\gamma^3\tau^2 d^2}{m^2} + \frac{\gamma^2 d^2}{m^2}\right] \sigma_*^2.
\end{aligned} \tag{23}$$

Taking $\beta_0 = \beta_1 = 1, \beta_3 = \gamma, \beta_2 = 4/\mu$ and using fact, that $\epsilon \leq \gamma\tau d/m$ inequality (23) takes form

$$\begin{aligned}
\mathbb{E}[\|x^{t+1} - x^*\|^2] &\leq \left(1 - \frac{3}{4\mu\gamma}\right) \mathbb{E}[\|x^t - x^*\|^2] + \frac{\epsilon d\beta_0\gamma}{m} \mathbb{E}[\|x^{t-\tau} - x^*\|^2] \\
&\quad + 4L\mathbb{E}\left[\frac{\epsilon d\gamma}{m\beta_0}(f(x^{t-\tau}) - f(x^*)) + 5\frac{d^2\gamma^2}{m^2}(f(x^t) - f(x^*))\right. \\
&\quad + \left.20\frac{d^4L^2}{m^4} \frac{\gamma^3\tau}{\mu} \sum_{s=t-\tau}^{t-1} (f(x^s) - f(x^*)) - \frac{\gamma}{2L}(f(x^t) - f(x^*))\right] \\
&\quad + 4\frac{d^2\gamma^2\tau}{m^2} \left[3 + 10\frac{d^2L^2}{m^2} \frac{\gamma}{\mu}\right] \sigma_*^2.
\end{aligned} \tag{24}$$

Let us perform the summation from $t = \tau$ to $t = T > \tau$ of equations (24) with coefficients p_k :

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636

$$\begin{aligned}
\sum_{t=\tau}^T p_t \mathbb{E} [\|x^{t+1} - x^*\|^2] &\leq \sum_{t=\tau}^T p_t (1 - \frac{3\mu\gamma}{4}) \mathbb{E} [\|x^t - x^*\|^2] \\
&+ \sum_{t=\tau}^T p_t \frac{\gamma\epsilon d}{m} \mathbb{E} [\|x^{t-\tau} - x^*\|^2] \\
&+ \sum_{t=\tau}^T p_t 4L \left(\frac{\gamma\epsilon d}{m} + 5 \frac{\gamma^2 d^2 \tau}{m^2} - \frac{\gamma}{2L} \right) \mathbb{E} [f(x^t) - f(x^*)] \quad (25) \\
&+ 20 \sum_{t=\tau}^T p_t 4L \frac{d^4 L^2 \gamma^3 \tau}{m^4 \mu} \sum_{s=t-\tau}^{t-1} \mathbb{E} [f(x^s) - f(x^*)] \\
&+ \sum_{t=\tau}^T p_t 4 \frac{d^2 \gamma^2 \tau}{m^2} \left[3 + 10 \frac{d^2 L^2 \gamma}{m^2 \mu} \right] \sigma_*^2.
\end{aligned}$$

1637 If $p_t = p^t$, $p = (1 - \mu\gamma/2)^{-1}$ and $\gamma = \gamma_1/\tau$, then, using the fact that $(1 - a/x)^{-x} \leq 2e^a \leq 2e$ if
1638 $x \geq 2$ and $0 \leq a \leq 1$, we can get that $p_\tau = (1 - \mu\gamma_1/(2\tau))^{-\tau} \leq 2e^{\mu\gamma_1/2} \leq 2e \leq 6$.

1639
1640
1641
1642

$$\sum_{t=\tau}^T p_t \sum_{s=t-\tau}^{t-1} a_s \leq p^\tau \sum_{t=\tau}^T \sum_{s=t-\tau}^{t-1} p_s a_s \leq 6\tau \sum_{t=0}^T p_t a_t.$$

1643 Using this we can estimate (25):

1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656

$$\begin{aligned}
\sum_{t=\tau}^T p_t \mathbb{E} [\|x^{t+1} - x^*\|^2] &\leq \sum_{t=\tau}^T p_t \left(1 - \mu\gamma + 6 \frac{\gamma\epsilon d}{m} \right) \mathbb{E} [\|x^t - x^*\|^2] \\
&+ \sum_{t=\tau}^T 4p_t L \left(\frac{\gamma\epsilon d}{m} + 5 \frac{\gamma^2 d^2 \tau}{m^2} + 120 \frac{d^4 L^2 \gamma^3 \tau^2}{m^4 \mu} - \frac{\gamma}{2L} \right) \mathbb{E} [f(x^t) - f(x^*)] \quad (26) \\
&+ 4 \sum_{t=\tau}^T p_t \left[3 + 10 \frac{d^2 L^2 \gamma}{m^2 \mu} \right] \sigma_*^2 + \sum_{t=0}^{\tau} p_{t+\tau} \frac{\gamma\epsilon d}{m} \mathbb{E} [\|x^t - x^*\|^2] \\
&+ 80 \sum_{t=0}^{\tau} p_{t+\tau} L \frac{d^4 L^2 \gamma^3 \tau}{m^4 \mu} \mathbb{E} [f(x^t) - f(x^*)].
\end{aligned}$$

1657 Taking

1658
1659
1660
1661

$$\gamma \leq \frac{m^2 \sqrt{\mu}}{24d^2 L^{3/2} \tau} \quad \text{and} \quad \epsilon = \min \left\{ \frac{\gamma d \tau}{m}; \frac{\mu m}{24d} \right\} \leq \frac{m \sqrt{\mu}}{24d} \min \left\{ \frac{1}{L^{3/2}}; \sqrt{\mu} \right\}.$$

1662 We get

1663
1664
1665
1666

$$\frac{\gamma\epsilon d}{m} + 5 \frac{\gamma^2 d^2 \tau}{m^2} + 120 \frac{d^4 L^2 \gamma^3 \tau^2}{m^4 \mu} - \frac{\gamma}{2L} \leq 0 \quad \text{and} \quad 1 - \frac{3\mu\gamma}{4} + 6 \frac{\gamma\epsilon d}{m} = 1 - \frac{\mu\gamma}{2}.$$

1667 Assume a notation

1668
1669
1670
1671
1672
1673

$$\begin{aligned}
\Delta_\tau &:= \sum_{t=0}^{\tau} p_{t+\tau} \frac{\gamma\epsilon d}{m} \mathbb{E} [\|x^t - x^*\|^2] + 80 \sum_{t=0}^{\tau} p_{t+\tau} L \frac{d^4 L^2 \gamma^3 \tau}{m^4 \mu} \mathbb{E} [f(x^t) - f(x^*)] \\
&\leq 120 \frac{\gamma^2 d^2}{m^2} \sqrt{\frac{\mu}{L}} \sum_{t=0}^{\tau} \left(\tau \mathbb{E} [\|x^t - x^*\|^2] + 4L \mathbb{E} [f(x^t) - f(x^*)] \right).
\end{aligned}$$

Using the notation of Δ_τ , (26) takes form

$$\sum_{t=\tau}^T p_t \mathbb{E} [\|x^{t+1} - x^*\|^2] \leq \sum_{t=\tau}^T p_t \left(1 - \frac{\mu\gamma}{2}\right) \mathbb{E} [\|x^t - x^*\|^2] + \sum_{t=\tau}^T 13 p_t \frac{\gamma^2 d^2 \tau}{m^2} \sigma_*^2 + \Delta_\tau.$$

Using $p_t = p^t$ and $p = (1 - \mu\gamma/2)^{-1}$ we can obtain:

$$\begin{aligned} \sum_{t=\tau}^T \left(1 - \frac{\mu\gamma}{2}\right)^{-t} \mathbb{E} [\|x^{t+1} - x^*\|^2] &\leq \sum_{t=\tau}^T \left(1 - \frac{\mu\gamma}{2}\right)^{-t+1} \mathbb{E} [\|x^t - x^*\|^2] \\ &+ \sum_{t=\tau}^T 13 \left(1 - \frac{\mu\gamma}{2}\right)^{-t} \frac{\gamma^2 d^2 \tau}{m^2} \sigma_*^2 + \Delta_\tau. \end{aligned}$$

The summed terms on the left and right sides are reduced, therefore this expression takes the form:

$$\begin{aligned} \left(1 - \frac{\mu\gamma}{2}\right)^{-T} \mathbb{E} [\|x^{T+1} - x^*\|^2] &\leq \left(1 - \frac{\mu\gamma}{2}\right)^{-\tau} \mathbb{E} [\|x^\tau - x^*\|^2] \\ &+ \sum_{t=\tau}^T 13 \left(1 - \frac{\mu\gamma}{2}\right)^{-t} \frac{\gamma^2 d^2 \tau}{m^2} \sigma_*^2 + \Delta_\tau. \end{aligned}$$

We can re-arrange this inequality:

$$\begin{aligned} \mathbb{E} [\|x^{T+1} - x^*\|^2] &\leq \left(1 - \frac{\mu\gamma}{2}\right)^{T-\tau} \mathbb{E} [\|x^\tau - x^*\|^2] \\ &+ \sum_{t=\tau}^T 13 \left(1 - \frac{\mu\gamma}{2}\right)^{T-t} \frac{\gamma^2 d^2 \tau}{m^2} \sigma_*^2 + \left(1 - \frac{\mu\gamma}{2}\right)^T \Delta_\tau. \end{aligned}$$

Using the fact that

$$\sum_{t=\tau}^T \left(1 - \frac{\mu\gamma}{2}\right)^{T-t} = \sum_{t=0}^{T-\tau} \left(1 - \frac{\mu\gamma}{2}\right)^t \leq \sum_{t=0}^{+\infty} \left(1 - \frac{\mu\gamma}{2}\right)^t = \frac{2}{\mu\gamma}.$$

We can estimate:

$$\mathbb{E} [\|x^{T+1} - x^*\|^2] \leq \left(1 - \frac{\mu\gamma}{2}\right)^{T-\tau} \mathbb{E} [\|x^\tau - x^*\|^2] + \left(1 - \frac{\mu\gamma}{2}\right)^T \Delta_\tau + 26 \frac{\gamma d^2 \tau}{\mu m^2} \sigma_*^2.$$

This finishes the proof. □

G EXTENSIONS FOR THEOREM 3

G.1 FULL VERSION OF THEOREM 3

Theorem 6 (Convergence of AMQSGD Algorithm 2, full version). *Consider Assumptions 1, 2 and 4. Let problem (1) be solved by Algorithm 2. Then for any $\gamma > 0, \varepsilon > 0, \tau > \tau_{\text{mix}}(\varepsilon), T > \tau$ and β, θ, η, p satisfying*

$$\gamma \lesssim \frac{\mu^{\frac{1}{3}} m^{\frac{1}{2}}}{\tau L^{\frac{4}{3}} d^{\frac{1}{2}}}, \quad p \lesssim \frac{m^2}{\tau^2 d^2 (\delta^2 + 1)}, \quad \varepsilon \lesssim \min \left\{ \frac{m^{\frac{7}{4}}}{d^{\frac{7}{4}} \tau^{\frac{5}{4}} L (\delta^2 + 1)}; \frac{m^{\frac{15}{4}}}{d^{\frac{15}{4}} \tau^{\frac{13}{4}} (\delta^2 + 1)^2} \right\}$$

$$\beta = \sqrt{\frac{2p^2 \mu \gamma}{3}}, \quad \eta = \sqrt{\frac{3}{2\mu\gamma}}, \quad \theta = \frac{p\eta^{-1} - 1}{\beta p \eta^{-1} - 1}$$

it holds that

$$F_{T+1} = \mathcal{O} \left(\exp \left[-(T - \tau) \sqrt{\frac{p^2 \mu \gamma}{3}} \right] F_\tau + \exp \left[-T \sqrt{\frac{p^2 \mu \gamma}{3}} \right] \Delta_\tau + \frac{\gamma}{\mu} \sigma^2 \right).$$

Here we use notations: $F_t := \mathbb{E}[\|x^t - x^*\|^2 + \frac{3}{\mu}(f(x_f^t) - f(x^*))]$ and $\Delta_\tau \leq \frac{\sqrt{\gamma}}{\tau^{\frac{4}{3}} \mu^{\frac{1}{3}}} \sum_{t=0}^{\tau} (\mathbb{E} \|\nabla f(x_g^t)\| + \mathbb{E} \|x^t - x^*\|^2 + \mathbb{E}[f(x_f^t) - f(x^*)])$.

G.2 FULL VERSION OF COROLLARY 2

Corollary 5 (Step tuning for Theorem 3, full version of Corollary 2). *Under the conditions of Theorem 3, choosing γ as*

$$\gamma \lesssim \min \left\{ \frac{\mu^{\frac{1}{3}}}{L^{\frac{4}{3}} \tau^{\frac{8}{3}}}; \frac{\log \left(\max \left\{ 2; \frac{\mu^{\frac{2}{3}} (F_\tau + \Delta_\tau) T}{\tau^{\frac{4}{3}} L^{\frac{2}{3}} \sigma^2} \right\} \right)}{\mu p^2 T^2} \right\},$$

in order to achieve ϵ -approximate solution (in terms of $\mathbb{E}[\|x^T - x^*\|^2] \leq \epsilon^2$) it takes

$$\mathcal{O} \left(\frac{d^2 L^{\frac{2}{3}} \tau^{\frac{4}{3}}}{m^2 \mu^{\frac{2}{3}}} \left((\delta^2 + 1) \log \left(\frac{1}{\epsilon} \right) + \frac{\sigma^2}{\mu \epsilon} \right) \right) \text{ iterations.}$$

G.3 PROOF OF THEOREM 6

Lemma 6. *Consider Algorithm 2 with $\theta = (p\eta^{-1} - 1)/(\beta\eta^{-1} - 1) < 1$. Then for any $y^t = \kappa x_f^t + (1 - \kappa)x^t \in \text{conv} \{x_f^t, x^t\}$ for any $s < t$ exist constants $\alpha_f^s, \alpha^s \geq 0$ and $c_r \geq 0$ such that*

$$y^t = \tilde{y}^s - p\gamma \sum_{r=s}^{t-1} c_r g^r = \alpha_f^s x_f^s + \alpha^s x^s - p\gamma \sum_{r=s}^{t-1} c_r g^r.$$

And $\alpha_f^s + \alpha^s = 1$ for any $s < t$. If $(1 - \kappa)\eta \leq 1$, then $c_r \leq t - s + 2$, otherwise we can only use the estimate $c_r \leq \eta$.

Proof. We start by writing out lines 3 and 10 of Algorithm 2:

$$x_f^s = x_g^{s-1} - p\gamma g^{s-1} = \theta x_f^{s-1} + (1 - \theta)x^{s-1} - p\gamma g^{s-1}. \quad (27)$$

Now let us handle expression $\eta x_g^k + (p - \eta)x_f^k + (1 - p)(1 - \beta)x^k + (1 - p)\beta x_g^k - x^*$ for a while. Taking into account the choice of θ such that $\theta = (p\eta^{-1} - 1)/(\beta p \eta^{-1} - 1)$ (in particular, $(p\eta^{-1} - 1) = (\beta p \eta^{-1} - 1)\theta$ and $\eta(1 - \beta p \eta^{-1})(1 - \theta) = p(1 - \beta)$), we get

$$\eta x_g^k + (p - \eta)x_f^k + (1 - p)(1 - \beta)x^k + (1 - p)\beta x_g^k$$

$$\begin{aligned}
&= (\eta + (1-p)\beta)x_g^k + (p-\eta)x_f^k + (1-p)(1-\beta)x^k \\
&= (\eta + (1-p)\beta)x_g^k + \eta(p\eta^{-1}-1)x_f^k + (1-p)(1-\beta)x^k \\
&= (\eta + (1-p)\beta)x_g^k + \eta(\beta p\eta^{-1}-1)\theta x_f^k + (1-p)(1-\beta)x^k \\
&= (\eta + (1-p)\beta)x_g^k + \eta(\beta p\eta^{-1}-1)(x_g^k - (1-\theta)x^k) + (1-p)(1-\beta)x^k \\
&= (\eta + (1-p)\beta)x_g^k + \eta(\beta p\eta^{-1}-1)(x_g^k - (1-\theta)x^k) + (1-p)(1-\beta)x^k \\
&= \beta x_g^k - \eta(\beta p\eta^{-1}-1)(1-\theta)x^k + (1-p)(1-\beta)x^k \\
&= \beta x_g^k + p(1-\beta)x^k + (1-p)(1-\beta)x^k \\
&= \beta x_g^k + (1-\beta)x^k.
\end{aligned}$$

Now we write out line 11 of Algorithm 2:

$$\begin{aligned}
x^s &= \beta x_g^{s-1} + (1-\beta)x^{s-1} - \eta x_g^{s-1} + \eta x_f^s = \beta x_g^{s-1} + (1-\beta)x^{s-1} - \eta p\gamma g^{s-1} \\
&= \beta(\theta x_f^{s-1} + (1-\theta)x^{s-1}) + (1-\beta)x^{s-1} - \eta p\gamma g^{s-1} \\
&= \beta\theta x_f^{s-1} + (1-\beta\theta)x^{s-1} - \eta p\gamma g^{s-1}.
\end{aligned} \tag{28}$$

Now we use induction. $x_f^t = \theta x_f^{s-1} + (1-\theta)x^{s-1} - p\gamma g^{s-1}$, then $\alpha_f^{t-1} = \theta \geq 0$, $\alpha^{t-1} = 1 - \theta \geq 0$, $c_r = 1 \leq \eta$ and $\alpha_f^{t-1} + \alpha^{t-1} = 1$, therefore base step is fulfilled. If $x_f^t = \alpha_f^s x_f^s + \alpha^s x^s - p\gamma \sum_{r=s}^{t-1} c_r g^r$ for some $s < t$, when with help of (27) and (28) we can write out

$$\begin{aligned}
x_f^t &= \alpha_f^s \left(\theta x_f^{s-1} + (1-\theta)x^{s-1} - p\gamma g^{s-1} \right) \\
&\quad + \alpha^s \left(\beta\theta x_f^{s-1} + (1-\beta\theta)x^{s-1} - \eta p\gamma g^{s-1} \right) - p\gamma \sum_{r=s}^{t-1} c_r g^r.
\end{aligned}$$

Therefore $\alpha_f^{s-1} = \alpha_f^s \theta + \alpha^s \beta\theta \geq 0$, $\alpha^{s-1} = \alpha_f^s (1-\theta) + \alpha^s (1-\beta\theta) \geq 0$ and $c_{s-1} = \alpha_f^s + \eta \alpha^s \leq \eta$. Then, the step of the induction is fulfilled, since $\alpha_f^{s-1} + \alpha^{s-1} = 1$. Therefore results of this Lemma are true for $y^t = x_f^t \in \text{conv} \{x_f^t, x^t\}$.

Consider $y^t = x^t \in \text{conv} \{x_f^t, x^t\}$. Form (28) follows that $\alpha_f^{t-1} = \beta\theta$ and $\alpha^{t-1} = 1 - \beta\theta$, therefore base step is fulfilled. The step of the induction will be the same as in $y^t = x_f^t$. Therefore results of this Lemma are true for $y^t = x^t$. Then, they are true for any $y^t \in \text{conv} \{x_f^t, x^t\}$.

If $y^t = \kappa x_f^t + (1-\kappa)x^t$, then $\alpha^s(y) = \kappa \alpha^s(x_f^t) + (1-\kappa)\alpha^s(x^t)$. Since $(1-\theta)\eta \leq 1$, then $\alpha^{t-1}(x_f^t)\eta \leq 1 = t - (t-1)$. Therefore $\alpha^s(x_f^t)\eta \leq t - s$ by induction, since $\alpha^{s-1}(x_f^t)\eta = \alpha_f^s(x_f^t)(1-\theta)\eta + (1-\beta\theta)\alpha^s(x_f^t)\eta \leq \alpha_f^s(x_f^t) + (1-\beta\theta)(t-s) \leq t - s + 1$.

Then, if $(1-\kappa)\eta \leq 1$, then $\alpha^s(y^t)\eta = \kappa \alpha^s(x_f^t)\eta + (1-\kappa)\eta \alpha^s(x^t) \leq \kappa(t-s) + \alpha^s(x^t) \leq t - s + 1$. Now we consider $c_s(y^t)$. $c_s(y^t) = \alpha_f^s(y^t) + \alpha^s(y^t)\eta \leq \alpha_f^s(y^t) + t - s + 1 \leq t - s + 2$.

□

Lemma 7. Assume 1, 2 and 4. Then for iterates of Algorithm 2 with $\theta = (p\eta^{-1}-1)/(\beta p\eta^{-1}-1)$, $\theta > 0, \eta \geq 1$, it holds that

$$\begin{aligned}
&\mathbb{E} \|x^{t+1} - x^*\|^2 \\
&\leq (1-\beta)\left(1 + \frac{\beta}{4}\right) \mathbb{E} \|x^t - x^*\|^2 + \beta\left(1 + \frac{\beta}{4}\right) \mathbb{E} \|x_g^t - x^*\|^2 + (\beta^2 - \beta) \mathbb{E} \|x^t - x_g^t\|^2 \\
&\quad + 10 \frac{d^2}{m^2} (\delta^2 + 1) p^2 \gamma^2 \eta^2 \mathbb{E} \|\nabla f(x_g^t)\|^2 + p^2 \gamma^2 \eta^2 \tau \left(32 \frac{\tau^2 d^2 L^2 p^2 \gamma^2}{m^2 \beta} + \frac{5}{4} \right) \sum_{r=t-\tau}^{t-1} \|g^r\|^2
\end{aligned}$$

$$\begin{aligned}
& + 3\varepsilon p\gamma\eta L \frac{d}{m} \sqrt{\delta^2 + 1} \mathbb{E} \left[\|x^{t-\tau} - x^*\|^2 \right] + 3\varepsilon p\gamma\eta L \frac{d}{m} \sqrt{\delta^2 + 1} \mathbb{E} \left[\|x_f^{t-\tau} - x^*\|^2 \right] \\
& - 2\gamma\eta^2 \mathbb{E} \langle \nabla f(x_g^t), x_g^t + (p\eta^{-1} - 1)x_f^t - p\eta^{-1}x^* \rangle + 2p\gamma\eta \left(\frac{\varepsilon d}{m\sqrt{\delta^2 + 1}L} + 4p\gamma\eta \frac{d^2}{m^2} \right) \sigma^2.
\end{aligned} \tag{29}$$

Proof. Using lines 10 and 11 of Algorithm 2, we get

$$\begin{aligned}
\mathbb{E} \|x^{t+1} - x^*\|^2 &= \mathbb{E} \left\| \eta x_f^{t+1} + (p - \eta)x_f^t + (1 - p)(1 - \beta)x^t + (1 - p)\beta x_g^t - x^* \right\|^2 \\
&= \mathbb{E} \left\| \eta x_g^t - p\gamma\eta g^t + (p - \eta)x_f^t + (1 - p)(1 - \beta)x^t + (1 - p)\beta x_g^t - x^* \right\|^2 \\
&= \mathbb{E} \left\| \eta x_g^t + (p - \eta)x_f^t + (1 - p)(1 - \beta)x^t + (1 - p)\beta x_g^t - x^* \right\|^2 + p^2\gamma^2\eta^2 \mathbb{E} \|g^t\|^2 \\
&\quad - 2p\gamma\eta \mathbb{E} \langle g^t, \eta x_g^t + (p - \eta)x_f^t + (1 - p)(1 - \beta)x^t + (1 - p)\beta x_g^t - x^* \rangle \\
&= \underbrace{\mathbb{E} \left\| \eta x_g^t + (p - \eta)x_f^t + (1 - p)(1 - \beta)x^t + (1 - p)\beta x_g^t - x^* \right\|^2}_{\textcircled{1}} + \underbrace{p^2\gamma^2\eta^2 \mathbb{E} \|g^t\|^2}_{\textcircled{2}} \\
&\quad - \underbrace{2p\gamma\eta \mathbb{E} \langle g^t - \nabla f(x_g^t), \eta x_g^t + (p - \eta)x_f^t + (1 - p)(1 - \beta)x^t + (1 - p)\beta x_g^t - x^* \rangle}_{\textcircled{3}} \\
&\quad - \underbrace{2p\gamma\eta \mathbb{E} \langle \nabla f(x_g^t), \eta x_g^t + (p - \eta)x_f^t + (1 - p)(1 - \beta)x^t + (1 - p)\beta x_g^t - x^* \rangle}_{\textcircled{4}}.
\end{aligned}$$

Consider ①. From Lemma 6, we know that

$$\eta x_g^t + (p - \eta)x_f^t + (1 - p)(1 - \beta)x^t + (1 - p)\beta x_g^t = \beta x_g^t + (1 - \beta)x^t.$$

It implies

$$\begin{aligned}
& \left\| \eta x_g^t + (p - \eta)x_f^t + (1 - p)(1 - \beta)x^t + (1 - p)\beta x_g^t - x^* \right\|^2 \\
&= \left\| \beta x_g^t + (1 - \beta)x^t - x^* \right\|^2 \\
&= \left\| \beta(x_g^t - x^t) + x^t - x^* \right\|^2 \\
&= \|x^t - x^*\|^2 + 2\beta \langle x^t - x^*, x_g^t - x^t \rangle + \beta^2 \|x_g^t - x^t\|^2 \\
&= \|x^t - x^*\|^2 + \beta(\|x_g^t - x^*\|^2 - \|x^t - x^*\|^2 - \|x_g^t - x^t\|^2) + \beta^2 \|x_g^t - x^t\|^2 \\
&= (1 - \beta) \|x^t - x^*\|^2 + \beta \|x_g^t - x^*\|^2 + (\beta^2 - \beta) \|x^t - x_g^t\|^2.
\end{aligned} \tag{30}$$

Consider ②. Using convexity of squared Euclidean norm and Lemma 4, one can obtain

$$\begin{aligned}
p^2\gamma^2\eta^2 \mathbb{E} \|g^t\|^2 &= p^2\gamma^2\eta^2 \mathbb{E} \left\| \frac{1}{n} \sum_{i=1}^n Q_t^i(\nabla f_i(x_g^t)) \right\|^2 \\
&\leq p^2\gamma^2\eta^2 \frac{1}{n} \sum_{i=1}^n \mathbb{E} \|Q_t^i(\nabla f_i(x_g^t))\|^2 \\
&\stackrel{(4)}{\leq} p^2\gamma^2\eta^2 \frac{d^2}{m^2} \frac{1}{n} \sum_{i=1}^n \mathbb{E} \|\nabla f_i(x_g^t)\|^2 \\
&\stackrel{(5)}{\leq} 2p^2\gamma^2\eta^2 \frac{d^2}{m^2} (\delta^2 + 1) \mathbb{E} \|\nabla f(x_g^t)\|^2 + 2p^2\gamma^2\eta^2 \frac{d^2}{m^2} \sigma^2,
\end{aligned} \tag{31}$$

where in the last inequality we used Lemma 5.

Consider ③. We first use Lemma 6 twice

$$x_g^t = \theta x_f^t + (1 - \theta)x^t = \alpha_f^{t-\tau} x_f^{t-\tau} + \alpha^{t-\tau} x^{t-\tau} - p\gamma \sum_{r=t-\tau}^{t-1} c_r g^r$$

$$\begin{aligned}
1890 \quad & \eta x_g^t + (p - \eta)x_f^t + (1 - p)(1 - \beta)x^t + (1 - p)\beta x_g^t = \beta x_g^t + (1 - \beta)x^t \\
1891 \quad & = \beta \theta x_f^t + (1 - \beta \theta)x^t \\
1892 \quad & \\
1893 \quad & = \hat{\alpha}_f^{t-\tau} x_f^{t-\tau} + \hat{\alpha}^{t-\tau} x^{t-\tau} - p\gamma \sum_{r=t-\tau}^{t-1} \hat{c}_r g^r. \\
1894 \quad & \\
1895 \quad &
\end{aligned}$$

1896 Next, we apply Corollary 3 with $\hat{a}^{t-\tau} = \nabla f_i(\tilde{x}_g^{t-\tau})$, where $\tilde{x}_g^{t-\tau} = \alpha_f^{t-\tau} x_f^{t-\tau} + \alpha^{t-\tau} x^{t-\tau}$, and
1897 $\hat{b}^{t-\tau} = \hat{\alpha}_f^{t-\tau} x_f^{t-\tau} + \hat{\alpha}^{t-\tau} x^{t-\tau} - x^*$, leading us to

$$\begin{aligned}
1899 \quad & -2p\gamma\eta \mathbb{E} \left\langle g^t - \nabla f(x_g^t), \eta x_g^t + (p - \eta)x_f^t + (1 - p)(1 - \beta)x^t + (1 - p)\beta x_g^t - x^* \right\rangle \\
1900 \quad & = -2p\gamma\eta \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left\langle Q_t^i(\nabla f_i(x_g^t)) - \nabla f_i(x_g^t), \eta x_g^t + (p - \eta)x_f^t + (1 - p)(1 - \beta)x^t \right. \\
1901 \quad & \left. + (1 - p)\beta x_g^t - x^* \right\rangle \\
1902 \quad & \leq \frac{\varepsilon d}{m\beta_0} p\gamma\eta \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left[\|\nabla f_i(\tilde{x}_g^{t-\tau})\|^2 \right] + \frac{\varepsilon d\beta_0}{m} p\gamma\eta \mathbb{E} \left[\left\| \hat{\alpha}_f^{t-\tau} x_f^{t-\tau} + \hat{\alpha}^{t-\tau} x^{t-\tau} - x^* \right\|^2 \right] \\
1903 \quad & + 4 \frac{d^2}{m^2} p\gamma\eta (\beta_1 + \beta_2) \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left[\|\nabla f_i(x_g^t) - \nabla f_i(\tilde{x}_g^{t-\tau})\|^2 \right] \\
1904 \quad & + p\gamma\eta \left(\frac{1}{\beta_1} + \frac{1}{\beta_3} \right) \mathbb{E} \left[\left\| -p\gamma \sum_{r=t-\tau}^{t-1} \hat{c}_r g^r \right\|^2 \right] \\
1905 \quad & + 4 \frac{d^2}{m^2} p\gamma\eta \beta_3 \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left[\|\nabla f_i(x_g^t)\|^2 \right] + \frac{p\gamma\eta}{\beta_2} \mathbb{E} \left[\|\beta x_g^t + (1 - \beta)x^t - x^*\|^2 \right]. \\
1906 \quad & \\
1907 \quad & \\
1908 \quad & \\
1909 \quad & \\
1910 \quad & \\
1911 \quad & \\
1912 \quad & \\
1913 \quad & \\
1914 \quad & \\
1915 \quad & \\
1916 \quad & \\
1917 \quad &
\end{aligned}$$

1918 Using Assumption 1 and Lemma 5 with $c_r \leq \tau \leq 2\tau$ and $\hat{c}_r \leq \eta$ one might obtain

$$\begin{aligned}
1919 \quad & -2p\gamma\eta \mathbb{E} \left\langle g^t - \nabla f(x_g^t), \eta x_g^t + (p - \eta)x_f^t + (1 - p)(1 - \beta)x^t + (1 - p)\beta x_g^t - x^* \right\rangle \\
1920 \quad & \leq \frac{2\varepsilon d}{m\beta_0} p\gamma\eta (\delta^2 + 1) \mathbb{E} \left[\|\nabla f(\tilde{x}_g^{t-\tau})\|^2 \right] + \frac{\varepsilon d\beta_0}{m} p\gamma\eta \mathbb{E} \left[\left\| \hat{\alpha}_f^{t-\tau} x_f^{t-\tau} + \hat{\alpha}^{t-\tau} x^{t-\tau} - x^* \right\|^2 \right] \\
1921 \quad & + 4 \frac{d^2 L^2}{m^2} p\gamma\eta (\beta_1 + \beta_2) \mathbb{E} \left[\left\| -p\gamma \sum_{r=t-\tau}^{t-1} c_r g^r \right\|^2 \right] + p\gamma\eta \left(\frac{1}{\beta_1} + \frac{1}{\beta_3} \right) \mathbb{E} \left[\left\| -p\gamma \sum_{r=t-\tau}^{t-1} \hat{c}_r g^r \right\|^2 \right] \\
1922 \quad & + 8 \frac{d^2}{m^2} (\delta^2 + 1) p\gamma\eta \beta_3 \mathbb{E} \left[\|\nabla f(x_g^t)\|^2 \right] + \frac{p\gamma\eta}{\beta_2} \mathbb{E} \left[\|\beta x_g^t + (1 - \beta)x^t - x^*\|^2 \right] \\
1923 \quad & + 2p\gamma\eta \left(\frac{\varepsilon d}{m\beta_0} + 4 \frac{d^2 \beta_3}{m^2} \right) \sigma^2 \tag{32} \\
1924 \quad & \leq \frac{\varepsilon d}{m} p\gamma\eta \left(2(\delta^2 + 1)L^2 \alpha_f^{t-\tau} \frac{1}{\beta_0} + \beta_0 \hat{\alpha}_f^{t-\tau} \right) \mathbb{E} \left[\|x_f^{t-\tau} - x^*\|^2 \right] \\
1925 \quad & + \frac{\varepsilon d}{m} p\gamma\eta \left(2(\delta^2 + 1)L^2 \alpha^{t-\tau} \frac{1}{\beta_0} + \beta_0 \hat{\alpha}^{t-\tau} \right) \mathbb{E} \left[\|x^{t-\tau} - x^*\|^2 \right] \\
1926 \quad & + p^3 \gamma^3 \eta \tau \left(4 \frac{\tau^2 d^2 L^2}{m^2} (\beta_1 + \beta_2) + \eta^2 \left(\frac{1}{\beta_1} + \frac{1}{\beta_3} \right) \right) \sum_{r=t-\tau}^{t-1} \|g^r\|^2 \\
1927 \quad & + 8 \frac{d^2}{m^2} (\delta^2 + 1) p\gamma\eta \beta_3 \mathbb{E} \left[\|\nabla f(x_g^t)\|^2 \right] \\
1928 \quad & + \frac{p\gamma\eta}{\beta_2} \beta \mathbb{E} \left[\|x_g^t - x^*\|^2 \right] + \frac{p\gamma\eta}{\beta_2} (1 - \beta) \mathbb{E} \left[\|x^t - x^*\|^2 \right] + 2p\gamma\eta \left(\frac{\varepsilon d}{m\beta_0} + 4 \frac{d^2 \beta_3}{m^2} \right) \sigma^2. \\
1929 \quad & \\
1930 \quad & \\
1931 \quad & \\
1932 \quad & \\
1933 \quad & \\
1934 \quad & \\
1935 \quad & \\
1936 \quad & \\
1937 \quad & \\
1938 \quad & \\
1939 \quad & \\
1940 \quad & \\
1941 \quad & \\
1942 \quad & \\
1943 \quad &
\end{aligned}$$

1944 Consider ④. Taking into account line 4 and the choice of θ such that $\theta = (p\eta^{-1} - 1)/(\beta p\eta^{-1} - 1)$,
 1945 one can note

$$\begin{aligned}
 1946 & \eta x_g^k + (p - \eta)x_f^k + (1 - p)(1 - \beta)x^k + (1 - p)\beta x_g^k - x^* \\
 1947 & = (\eta + (1 - p)\beta)x_g^k + (p - \eta)x_f^k + (1 - p)(1 - \beta)x^k - x^* \\
 1948 & = \eta p^{-1} \left((p + (1 - p)p^{-1}\eta\beta)x_g^k + (p\eta^{-1} - 1)px_f^k + (1 - p)(1 - \beta)p\eta^{-1}x^k - \eta^{-1}px^* \right) \\
 1949 & = \eta p^{-1} \left((p + (1 - p)p^{-1}\eta\beta)x_g^k + (p\eta^{-1} - 1)px_f^k + (1 - p)(1 - \beta p\eta^{-1})(1 - \theta)x^k - \eta^{-1}px^* \right) \\
 1950 & = \eta p^{-1} \left((p + (1 - p)p^{-1}\eta\beta)x_g^k + (p\eta^{-1} - 1)px_f^k + (1 - p)(1 - \beta p\eta^{-1})(x_g^k - \theta x_f^k) - \eta^{-1}px^* \right) \\
 1951 & = \eta p^{-1} \left(x_g^k + (p\eta^{-1} - 1)px_f^k - (1 - p)(1 - \beta p\eta^{-1})\theta x_f^k - \eta^{-1}px^* \right) \\
 1952 & = \eta p^{-1} \left(x_g^k + (p\eta^{-1} - 1)px_f^k - (1 - p)(p\eta^{-1} - 1)x_f^k - \eta^{-1}px^* \right) \\
 1953 & = \eta p^{-1} \left(x_g^k + (p\eta^{-1} - 1)x_f^k - \eta^{-1}px^* \right). \tag{33}
 \end{aligned}$$

1954 Using that, we get

$$\begin{aligned}
 1955 & -2p\gamma\eta \mathbb{E} \left\langle \nabla f(x_g^t), \eta x_g^t + (p - \eta)x_f^t + (1 - p)(1 - \beta)x^t + (1 - p)\beta x_g^t - x^* \right\rangle \\
 1956 & = -2\gamma\eta^2 \mathbb{E} \left\langle \nabla f(x_g^t), x_g^t + (p\eta^{-1} - 1)x_f^t - p\eta^{-1}x^* \right\rangle. \tag{34}
 \end{aligned}$$

1957 Summing (30), (31), (32) and (34) with $\beta_0 = \sqrt{\delta^2 + 1}L$, $\beta_1 = \beta_2 = \frac{4p\gamma\eta}{\beta}$ and $\beta_3 = p\gamma\eta$ we finish
 1958 the proof. \square

1959 **Lemma 8.** Assume 1, 2 and 4. Then for iterates of Algorithm 2 and for any $u \in \mathbb{R}^d$ it holds that

$$\begin{aligned}
 1960 & \mathbb{E} \left[f(x_f^{t+1}) \right] \leq \mathbb{E} [f(u)] - \mathbb{E} \left[\langle \nabla f(x_g^t), u - x_g^t \rangle \right] - \frac{\mu}{2} \|u - x_g^t\| - \frac{p\gamma}{2} \mathbb{E} \left[\|\nabla f(x_g^t)\|^2 \right] \\
 1961 & + 2\varepsilon\gamma \mathbb{E} \left[\|\nabla f(\tilde{x}_g^{t-\tau})\|^2 \right] + 20 \frac{L^2 d^3 \gamma^3 p^2 \tau^3 (\delta^2 + 1)}{m^3} \sum_{s=t-\tau}^{t-1} \mathbb{E} \left[\|\nabla f(x_g^s)\|^2 \right] + 23 \frac{L^2 d^3 \gamma^3 p^2 \tau^4}{m^3} \sigma^2,
 \end{aligned}$$

1962 where

$$\gamma \leq \frac{1}{L} \quad \text{and} \quad p \leq \frac{m^2}{12(\delta^2 + 1)d^2}.$$

1963 *Proof.* Using 1 with $x = x_f^{t+1}$, $y = x_g^t$ and line 3 of Algorithm 2 we get

$$\begin{aligned}
 1964 & \mathbb{E} \left[f(x_f^{t+1}) \right] \leq \mathbb{E} [f(x_g^t)] + \mathbb{E} \left[\langle \nabla f(x_g^t), x_f^{t+1} - x_g^t \rangle \right] + \frac{L}{2} \mathbb{E} \left[\|x_f^{t+1} - x_g^t\|^2 \right] \\
 1965 & = \mathbb{E} [f(x_g^t)] - p\gamma \mathbb{E} \left[\langle \nabla f(x_g^t), g^t \rangle \right] + \frac{Lp^2\gamma^2}{2} \mathbb{E} \left[\|g^t\|^2 \right] \\
 1966 & = \mathbb{E} [f(x_g^t)] - p\gamma \mathbb{E} \left[\langle \nabla f(x_g^t), \nabla f(x_g^t) \rangle \right] - p\gamma \mathbb{E} \left[\langle \nabla f(x_g^t), g^k - \nabla f(x_g^t) \rangle \right] \\
 1967 & + \frac{Lp^2\gamma^2}{2} \mathbb{E} \left[\|g^t\|^2 \right]. \tag{35}
 \end{aligned}$$

1968 Consider $\mathbb{E} \left[\langle \nabla f(x_g^t), g^k - \nabla f(x_g^t) \rangle \right]$. Using Corollary 3 with $a^t = \nabla f_i(x_g^t)$, $b^t =$
 1969 $\nabla f(x_g^t)$, $\hat{a}^{t-\tau} = \nabla f_i(\tilde{x}_g^{t-\tau})$, $\hat{b}^{t-\tau} = \nabla f(\tilde{x}_g^{t-\tau})$, where $x_g^t \in \text{conv} \{x_f^t, x^t\} = \tilde{x}_g^{t-\tau} -$
 1970 $p\gamma \sum_{s=t-\tau}^{t-1} c_s g^s$ from Lemma 6. Using Assumption 1 we obtain

$$\begin{aligned}
 1971 & 2 \left| \mathbb{E} \left[\langle \nabla f(x_g^t), g^k - \nabla f(x_g^t) \rangle \right] \right| \leq \frac{\varepsilon d}{m\beta_0} \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \|\nabla f_i(\tilde{x}_g^{t-\tau})\|^2 \right] + \frac{\varepsilon d\beta_0}{m} \mathbb{E} \left[\|\nabla f(\tilde{x}_g^{t-\tau})\|^2 \right] \\
 1972 & + 4 \frac{d^2 L^2}{m^2} (\beta_1 + \beta_2) \mathbb{E} \left[\|x_g^t - \tilde{x}_g^{t-\tau}\|^2 \right] + L^2 \left(\frac{1}{\beta_1} + \frac{1}{\beta_3} \right) \mathbb{E} \left[\|x_g^t - \tilde{x}_g^{t-\tau}\|^2 \right]
 \end{aligned}$$

$$+ 4 \frac{d^2}{m^2} \beta_3 \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x_g^t)\|^2 \right] + \frac{1}{\beta_2} \mathbb{E} \left[\|\nabla f(x_g^t)\|^2 \right].$$

Taking $\beta_0 = \sqrt{\delta^2 + 1}$, $\beta_1 = m/d$, $\beta_2 = m/(dp)$, $\beta_3 = pm/d$ and using results from Lemma 5 we obtain

$$\begin{aligned} 2 \left| \mathbb{E} [\langle \nabla f(x_g^t), g^k - \nabla f(x_g^t) \rangle] \right| &\leq \frac{2\varepsilon d}{m} \left(\sqrt{\delta^2 + 1} \mathbb{E} \left[\|\nabla f(\tilde{x}_g^{t-\tau})\|^2 \right] + \frac{\sigma^2}{\sqrt{\delta^2 + 1}} \right) \\ &+ \frac{dp}{m} \mathbb{E} \left[\|\nabla f(x_g^t)\|^2 \right] + 10 \frac{L^2 d}{pm} \mathbb{E} \left[\left\| -p\gamma \sum_{s=t-\tau}^{t-1} c_s \frac{1}{n} \sum_{i=1}^n Q_s^i(\nabla f_i(x_g^s)) \right\|^2 \right] \\ &+ \frac{8dp}{m} \left((\delta^2 + 1) \mathbb{E} \left[\|\nabla f(x_g^t)\|^2 \right] + \sigma^2 \right) + \frac{\varepsilon d \sqrt{\delta^2 + 1}}{m} \mathbb{E} \left[\|\nabla f(\tilde{x}_g^{t-\tau})\|^2 \right]. \end{aligned}$$

Using Lemma 4 and 5, convexity of the squared norm and the fact that $c_s \leq t - s + 2 \leq \tau + 2 \leq 2\tau$ we obtain

$$\begin{aligned} 2 \left| \mathbb{E} [\langle \nabla f(x_g^t), g^k - \nabla f(x_g^t) \rangle] \right| &\leq \frac{3\varepsilon d \sqrt{\delta^2 + 1}}{m} \mathbb{E} \left[\|\nabla f(\tilde{x}_g^{t-\tau})\|^2 \right] + \\ &+ 40 \frac{L^2 d^3 \gamma^2 p \tau^3}{m^3} \sum_{s=t-\tau}^{t-1} \mathbb{E} \left[(\delta^2 + 1) \|\nabla f(x_g^s)\|^2 + \sigma^2 \right] \\ &+ \frac{9dp(\delta^2 + 1)}{m} \mathbb{E} \left[\|\nabla f(x_g^t)\|^2 \right] + \frac{2d}{m} \left(\frac{\varepsilon}{\sqrt{\delta^2 + 1}} + p \right) \sigma^2. \end{aligned}$$

Using the fact that $L^2 \gamma^2 d^2 / m^2 \tau^4 \eta^2 \geq 1$ and $\varepsilon \leq \sqrt{\delta^2 + 1} p$ we obtain

$$\begin{aligned} 2 \left| \mathbb{E} [\langle \nabla f(x_g^t), g^k - \nabla f(x_g^t) \rangle] \right| &\leq \frac{3\varepsilon d \sqrt{\delta^2 + 1}}{m} \mathbb{E} \left[\|\nabla f(\tilde{x}_g^{t-\tau})\|^2 \right] + 44 \frac{L^2 d^3 \gamma^2 p \eta^2 \tau^4}{m^3} \sigma^2 \\ &+ 40 \frac{L^2 d^3 \gamma^2 p \tau^3 (\delta^2 + 1)}{m^3} \sum_{s=t-\tau}^{t-1} \mathbb{E} \left[\|\nabla f(x_g^s)\|^2 \right] + \frac{9dp(\delta^2 + 1)}{m} \mathbb{E} \left[\|\nabla f(x_g^t)\|^2 \right]. \end{aligned}$$

Using this result, Lemmas 4 and 5 we can estimate (35):

$$\begin{aligned} \mathbb{E} \left[f(x_f^{t+1}) \right] &= \mathbb{E} \left[f(x_g^t) \right] - p\gamma \mathbb{E} \left[\|\nabla f(x_g^t)\|^2 \right] \\ &- p\gamma \mathbb{E} \left[\langle \nabla f(x_g^t), g^k - \nabla f(x_g^t) \rangle \right] + \frac{L}{2} \mathbb{E} \left[\|g^t\|^2 \right] \\ &\leq \mathbb{E} \left[f(x_g^t) \right] - p\gamma \mathbb{E} \left[\|\nabla f(x_g^t)\|^2 \right] + \frac{2\varepsilon p \gamma d \sqrt{\delta^2 + 1}}{m} \mathbb{E} \left[\|\nabla f(\tilde{x}_g^{t-\tau})\|^2 \right] + \\ &+ 20 \frac{L^2 d^3 \gamma^3 p^2 \tau^3 (\delta^2 + 1)}{m^3} \sum_{s=t-\tau}^{t-1} \mathbb{E} \left[\|\nabla f(x_g^s)\|^2 \right] + \frac{5d\gamma p^2 (\delta^2 + 1)}{m} \mathbb{E} \left[\|\nabla f(x_g^t)\|^2 \right] \\ &+ 22 \frac{L^2 d^3 \gamma^3 p^2 \tau^4}{m^3} \sigma^2 + \frac{L p^2 \gamma^2 d^2}{m^2} (\delta^2 + 1) \mathbb{E} \left[\|\nabla f(x_g^t)\|^2 \right] + \frac{L p^2 \gamma^2 d^2}{m^2} \sigma^2. \end{aligned}$$

Taking

$$\gamma \leq \frac{1}{L} \quad \text{and} \quad p \leq \frac{m^2}{12(\delta^2 + 1)d^2},$$

we obtain

$$\begin{aligned} \mathbb{E} [f(x_f^{t+1})] &\leq \mathbb{E} [f(x_g^t)] - \frac{p\gamma}{2} \mathbb{E} [\|\nabla f(x_g^t)\|^2] + 2\varepsilon\gamma \mathbb{E} [\|\nabla f(\tilde{x}_g^{t-\tau})\|^2] + \\ &+ 20 \frac{L^2 d^3 \gamma^3 p^2 \tau^3 (\delta^2 + 1)}{m^3} \sum_{s=t-\tau}^{t-1} \mathbb{E} [\|\nabla f(x_g^s)\|^2] + 23 \frac{L^2 d^3 \gamma^3 p^2 \tau^4}{m^3} \sigma^2. \end{aligned}$$

Using 2 with $x = u$ and $y = x_g^t$, one can conclude that for any $u \in \mathbb{R}^d$ it holds

$$\begin{aligned} \mathbb{E} [f(x_f^{t+1})] &\leq \mathbb{E} [f(u)] - \mathbb{E} [\langle \nabla f(x_g^t), u - x_g^t \rangle] - \frac{\mu}{2} \|u - x_g^t\| \\ &- \frac{p\gamma}{2} \mathbb{E} [\|\nabla f(x_g^t)\|^2] + 2\varepsilon\gamma \mathbb{E} [\|\nabla f(\tilde{x}_g^{t-\tau})\|^2] + \\ &+ 20 \frac{L^2 d^3 \gamma^3 p^2 \tau^3 (\delta^2 + 1)}{m^3} \sum_{s=t-\tau}^{t-1} \mathbb{E} [\|\nabla f(x_g^s)\|^2] + 23 \frac{L^2 d^3 \gamma^3 p^2 \tau^4}{m^3} \sigma^2. \end{aligned}$$

This finishes the proof. \square

Theorem 7 (Theorem 3). *Consider Assumptions 1, 2 and 4. Let problem (1) be solved by Algorithm 2. Then for any $\gamma > 0, \varepsilon > 0, \tau > \tau_{\text{mix}}(\varepsilon), T > \tau$ and β, θ, η, p satisfying*

$$\begin{aligned} \gamma &\leq \frac{\mu^{\frac{1}{3}} m^{\frac{1}{2}}}{2\tau L^{\frac{1}{3}} d^{\frac{1}{2}}}, \quad \varepsilon \leq \min \left\{ \frac{m^{\frac{7}{4}}}{6d^{\frac{7}{4}} \tau^{\frac{5}{4}} L(\delta^2 + 1)}; \frac{m^{\frac{5}{4}}}{\sqrt{2}\tau^{\frac{3}{4}} \mu^{\frac{1}{3}} L^{\frac{2}{3}} d^{\frac{5}{4}}}; \frac{m^{\frac{15}{4}}}{6d^{\frac{15}{4}} \tau^{\frac{13}{4}} (\delta^2 + 1)^2} \right\}, \\ p &\leq \frac{m^2}{13d^2(\delta^2 + 1)\tau^2}, \quad \beta = \sqrt{\frac{2p^2\mu\gamma}{3}}, \quad \eta = \sqrt{\frac{3}{2\mu\gamma}}, \quad \theta = \frac{p\eta^{-1} - 1}{\beta p\eta^{-1} - 1}. \end{aligned}$$

it holds that

$$\begin{aligned} \mathbb{E} [\|x^{T+1} - x^*\|^2 + \frac{3}{\mu} (f(x_f^{T+1}) - f(x^*))] &\leq \exp \left(-(T - \tau) \sqrt{\frac{2p^2\mu\gamma}{3}} \right) F_\tau \\ &+ \exp \left(-T \sqrt{\frac{2p^2\mu\gamma}{3}} \right) \Delta_\tau + \frac{45\gamma}{\mu} \sigma^2, \end{aligned}$$

where $F_\tau := \mathbb{E} [\|x^\tau - x^*\|^2 + \frac{3}{\mu} (f(x_f^\tau) - f(x^*))]$ and $\Delta_\tau \leq \frac{\sqrt{\gamma}}{\tau^{\frac{1}{3}} \mu^{\frac{1}{3}}} \sum_{t=0}^{\tau} (\mathbb{E} \|\nabla f(x_g^t)\| + \mathbb{E} \|x^t - x^*\|^2 + \mathbb{E} [f(x_f^t) - f(x^*)])$.

Proof. We start by using Lemma 8 with $u = x^*$ and $u = x_f^t$

$$\begin{aligned} \mathbb{E} [f(x_f^{t+1})] &\leq \mathbb{E} [f(x^*)] - \mathbb{E} [\langle \nabla f(x_g^t), x^* - x_g^t \rangle] - \frac{\mu}{2} \|x^* - x_g^t\| - \frac{p\gamma}{2} \mathbb{E} [\|\nabla f(x_g^t)\|^2] \\ &+ 2\varepsilon\gamma \mathbb{E} [\|\nabla f(\tilde{x}_g^{t-\tau})\|^2] + 20 \frac{L^2 d^3 \gamma^3 p^2 \tau^3 (\delta^2 + 1)}{m^3} \sum_{s=t-\tau}^{t-1} \mathbb{E} [\|\nabla f(x_g^s)\|^2] + 23 \frac{L^2 d^3 \gamma^3 p^2 \tau^4}{m^3} \sigma^2, \\ \mathbb{E} [f(x_f^{t+1})] &\leq \mathbb{E} [f(x_f^t)] - \mathbb{E} [\langle \nabla f(x_g^t), x_f^t - x_g^t \rangle] - \frac{\mu}{2} \|x_f^t - x_g^t\| - \frac{p\gamma}{2} \mathbb{E} [\|\nabla f(x_g^t)\|^2] \\ &+ 2\varepsilon\gamma \mathbb{E} [\|\nabla f(\tilde{x}_g^{t-\tau})\|^2] + 20 \frac{L^2 d^3 \gamma^3 p^2 \tau^3 (\delta^2 + 1)}{m^3} \sum_{s=t-\tau}^{t-1} \mathbb{E} [\|\nabla f(x_g^s)\|^2] + 23 \frac{L^2 d^3 \gamma^3 p^2 \tau^4}{m^3} \sigma^2. \end{aligned}$$

Summing the first inequality with coefficient $2p\gamma\eta$, the second with coefficient $2p\gamma\eta(\eta - p)$ and (29), we get

$$\mathbb{E} [\|x^{t+1} - x^*\|^2 + 2\gamma\eta^2 f(x_f^{t+1})]$$

$$\begin{aligned}
&\leq (1 - \beta)(1 + \frac{\beta}{4}) \mathbb{E} \|x^t - x^*\|^2 + \beta(1 + \frac{\beta}{4}) \mathbb{E} \|x_g^t - x^*\|^2 + (\beta^2 - \beta) \mathbb{E} \|x^t - x_g^t\|^2 \\
&\quad + 10 \frac{d^2}{m^2} (\delta^2 + 1) p^2 \gamma^2 \eta^2 \mathbb{E} \|\nabla f(x_g^t)\|^2 + p^2 \gamma^2 \eta^2 \tau \left(32 \frac{\tau^2 d^2 L^2 p^2 \gamma^2}{m^2 \beta} + \frac{5}{4} \right) \sum_{r=t-\tau}^{t-1} \|g^r\|^2 \\
&\quad + 3\epsilon p \gamma \eta L \frac{d}{m} \sqrt{\delta^2 + 1} \mathbb{E} [\|x^{t-\tau} - x^*\|^2] + 3\epsilon p \gamma \eta L \frac{d}{m} \sqrt{\delta^2 + 1} \mathbb{E} [\|x_f^{t-\tau} - x^*\|^2] \\
&\quad - 2\gamma \eta^2 \mathbb{E} \langle \nabla f(x_g^t), x_g^t + (p\eta^{-1} - 1)x_f^t - p\eta^{-1}x^* \rangle + 2p\gamma \eta \left(\frac{\epsilon d}{m\sqrt{\delta^2 + 1}L} + 4p\gamma \eta \frac{d^2}{m^2} \right) \sigma^2 \\
&\quad + 2p\gamma \eta \left(\mathbb{E} [f(x^*)] - \mathbb{E} [\langle \nabla f(x_g^t), x^* - x_g^t \rangle] - \frac{\mu}{2} \|x^* - x_g^t\| - \frac{p\gamma}{2} \mathbb{E} [\|\nabla f(x_g^t)\|^2] \right) \\
&\quad + 2\epsilon \gamma \mathbb{E} [\|\nabla f(\tilde{x}_g^{t-\tau})\|^2] + 20 \frac{L^2 d^3 \gamma^3 p^2 \tau^3 (\delta^2 + 1)}{m^3} \sum_{s=t-\tau}^{t-1} \mathbb{E} [\|\nabla f(x_g^s)\|^2] \\
&\quad + 23 \frac{L^2 d^3 \gamma^3 p^2 \tau^4}{m^3} \sigma^2) \\
&\quad + 2\gamma \eta (\eta - p) \left(\mathbb{E} [f(x_f^t)] - \mathbb{E} [\langle \nabla f(x_g^t), x_f^t - x_g^t \rangle] - \frac{\mu}{2} \|x_f^t - x_g^t\| - \frac{p\gamma}{2} \mathbb{E} [\|\nabla f(x_g^t)\|^2] \right) \\
&\quad + 2\epsilon \gamma \mathbb{E} [\|\nabla f(\tilde{x}_g^{t-\tau})\|^2] + 20 \frac{L^2 d^3 \gamma^3 p^2 \tau^3 (\delta^2 + 1)}{m^3} \sum_{s=t-\tau}^{t-1} \mathbb{E} [\|\nabla f(x_g^s)\|^2] \\
&\quad + 23 \frac{L^2 d^3 \gamma^3 p^2 \tau^4}{m^3} \sigma^2) \\
&\leq (1 - \beta)(1 + \frac{\beta}{4}) \mathbb{E} \|x^t - x^*\|^2 + (\beta + \frac{\beta^2}{4} - p\gamma \eta \mu) \mathbb{E} \|x_g^t - x^*\|^2 + (\beta^2 - \beta) \mathbb{E} \|x^t - x_g^t\|^2 \\
&\quad + p^2 \gamma^2 \eta^2 \left(10 \frac{d^2}{m^2} (\delta^2 + 1) - \frac{1}{p} \right) \mathbb{E} \|\nabla f(x_g^t)\|^2 + 2p\gamma \eta \mathbb{E} f(x^*) + 2\gamma \eta (\eta - p) \mathbb{E} f(x_f^t) \\
&\quad + p^2 \gamma^2 \eta^2 \tau (\delta^2 + 1) \frac{d^2}{m^2} \left(32 \frac{\tau^2 d^2 L^2 p^2 \gamma^2}{m^2 \beta} + \frac{5}{4} \right) \sum_{r=t-\tau}^{t-1} \mathbb{E} \|\nabla f(x_g^r)\|^2 \\
&\quad + \epsilon \gamma \eta L (3p \frac{d}{m} \sqrt{\delta^2 + 1} + 2\gamma \eta L) \mathbb{E} [\|x^{t-\tau} - x^*\|^2] \\
&\quad + \epsilon \gamma \eta L (3p \frac{d}{m} \sqrt{\delta^2 + 1} + 2\gamma \eta L) \mathbb{E} [\|x_f^{t-\tau} - x^*\|^2] \\
&\quad + 2p\gamma \eta \left(\frac{\epsilon d}{m\sqrt{\delta^2 + 1}L} + 4p\gamma \eta \frac{d^2}{m^2} \right. \\
&\quad \quad \left. + 23p\gamma^3 \eta \tau^4 \frac{d^3}{m^3} L^2 + p\gamma \eta \tau^2 \frac{d^2}{m^2} \left(16 \frac{\tau^2 d^2 L^2 p^2 \gamma^2}{m^2 \beta} + \frac{5}{8} \right) \right) \sigma^2,
\end{aligned}$$

where in the last inequality we used Lemma 5 and Assumption 1. Since $\beta < 1$, the choice of $p\gamma \eta \mu = \frac{3\beta}{2}$ gives

$$\begin{aligned}
(1 - \beta)(1 + \frac{\beta}{4}) &\leq 1 - \frac{3\beta}{4}, \\
\beta + \frac{\beta^2}{4} - p\gamma \eta \mu &\leq \frac{3\beta}{2} - p\gamma \eta \mu \leq 0, \\
\beta^2 - \beta &\leq 0.
\end{aligned}$$

This lead us to

$$\begin{aligned}
& \mathbb{E}[\|x^{t+1} - x^*\|^2 + 2\gamma\eta^2(f(x_f^{t+1}) - f(x^*))] \\
& \leq (1 - \frac{3\beta}{4}) \mathbb{E}\|x^t - x^*\|^2 + 2p\gamma\eta^2(1 - \frac{p}{\eta}) \mathbb{E}[f(x_f^t) - f(x^*)] \\
& \quad + p^2\gamma^2\eta^2 \left(10 \frac{d^2}{m^2} (\delta^2 + 1) - \frac{1}{p} \right) \mathbb{E}\|\nabla f(x_g^t)\| \\
& \quad + p^2\gamma^2\eta^2\tau(\delta^2 + 1) \frac{d^2}{m^2} \left(32 \frac{\tau^2 d^2 L^2 p^2 \gamma^2}{m^2 \beta} + \frac{5}{4} \right) \sum_{r=t-\tau}^{t-1} \mathbb{E}\|\nabla f(x_g^r)\| \quad (36) \\
& \quad + \varepsilon\gamma\eta L \left(3p \frac{d}{m} \sqrt{\delta^2 + 1} + 2\gamma\eta L \right) \mathbb{E}\left[\|x^{t-\tau} - x^*\|^2\right] \\
& \quad + \varepsilon\gamma\eta L \left(3p \frac{d}{m} \sqrt{\delta^2 + 1} + 2\gamma\eta L \right) \frac{2}{\mu} \mathbb{E}[f(x_f^{t-\tau}) - f(x^*)] \\
& \quad + 2p\gamma\eta \left(\frac{\varepsilon d}{m\sqrt{\delta^2 + 1}L} + 4p\gamma\eta \frac{d^2}{m^2} \right. \\
& \quad \left. + 23p\gamma^3\eta\tau^4 \frac{d^3}{m^3} L^2 + p\gamma\eta\tau^2 \frac{d^2}{m^2} \left(16 \frac{\tau^2 d^2 L^2 p^2 \gamma^2}{m^2 \beta} + \frac{5}{8} \right) \right) \sigma^2,
\end{aligned}$$

where we also used Assumption 2 and subtracted $2\gamma\eta^2 f(x^*)$ from both sides. Next, we perform the summation from $t = \tau$ to $t = T > \tau$ of equations (36) with coefficients p_t :

$$\begin{aligned}
& \sum_{t=\tau}^T p_t \mathbb{E}[\|x^{t+1} - x^*\|^2 + 2\gamma\eta^2(f(x_f^{t+1}) - f(x^*))] \\
& \leq \sum_{t=\tau}^T p_t \left(1 - \frac{3\beta}{4} \right) \mathbb{E}\|x^t - x^*\|^2 \\
& \quad + \sum_{t=\tau}^T p_t 2p\gamma\eta^2 \left(1 - \frac{p}{\eta} \right) \mathbb{E}[f(x_f^t) - f(x^*)] + \sum_{t=\tau}^T p_t p^2\gamma^2\eta^2 \left(10 \frac{d^2}{m^2} (\delta^2 + 1) - \frac{1}{p} \right) \mathbb{E}\|\nabla f(x_g^t)\| \\
& \quad + \sum_{t=\tau}^T p_t p^2\gamma^2\eta^2\tau(\delta^2 + 1) \frac{d^2}{m^2} \left(32 \frac{\tau^2 d^2 L^2 p^2 \gamma^2}{m^2 \beta} + \frac{5}{4} \right) \sum_{r=t-\tau}^{t-1} \mathbb{E}\|\nabla f(x_g^r)\| \\
& \quad + \sum_{t=\tau}^T p_t \varepsilon\gamma\eta L \left(3p \frac{d}{m} \sqrt{\delta^2 + 1} + 2\gamma\eta L \right) \mathbb{E}\left[\|x^{t-\tau} - x^*\|^2\right] \\
& \quad + \sum_{t=\tau}^T p_t \varepsilon\gamma\eta L \left(3p \frac{d}{m} \sqrt{\delta^2 + 1} + 2\gamma\eta L \right) \frac{2}{\mu} \mathbb{E}[f(x_f^{t-\tau}) - f(x^*)] \\
& \quad + \sum_{t=\tau}^T p_t 2p\gamma\eta \left(\frac{\varepsilon d}{m\sqrt{\delta^2 + 1}L} + 4p\gamma\eta \frac{d^2}{m^2} \right. \\
& \quad \left. + 23p\gamma^3\eta\tau^4 \frac{d^3}{m^3} L^2 + p\gamma\eta\tau^2 \frac{d^2}{m^2} \left(16 \frac{\tau^2 d^2 L^2 p^2 \gamma^2}{m^2 \beta} + \frac{5}{8} \right) \right) \sigma^2.
\end{aligned}$$

Similar as in Theorem 5 we take $p_t = p^t$, $p = (1 - \frac{\beta}{2})^{-1}$, it implies $p_\tau \leq 6$ and therefore

$$\begin{aligned}
& \sum_{t=\tau}^T p_t \mathbb{E}[\|x^{t+1} - x^*\|^2 + 2\gamma\eta^2(f(x_f^{t+1}) - f(x^*))] \\
& \leq \sum_{t=\tau}^T p_t \left(1 - \frac{3\beta}{4} + 6\varepsilon\gamma\eta L \left(3p \frac{d}{m} \sqrt{\delta^2 + 1} + 2\gamma\eta L \right) \right) \mathbb{E}\|x^t - x^*\|^2
\end{aligned}$$

$$\begin{aligned}
& + \sum_{t=\tau}^T p_t \left(2p\gamma\eta^2 \left(1 - \frac{p}{\eta}\right) + 12 \frac{\varepsilon\gamma\eta L}{\mu} \left(3p \frac{d}{m} \sqrt{\delta^2 + 1} + 2\gamma\eta L\right) \right) \mathbb{E}[f(x_f^t) - f(x^*)] \\
& + \sum_{t=\tau}^T p_t p^2 \gamma^2 \eta^2 \left(10 \frac{d^2}{m^2} (\delta^2 + 1) - \frac{1}{p} + \tau^2 (\delta^2 + 1) \frac{d^2}{m^2} \left(32 \frac{\tau^2 d^2 L^2 p^2 \gamma^2}{m^2 \beta} + \frac{5}{4} \right) \right) \mathbb{E} \|\nabla f(x_g^t)\| \\
& + \sum_{t=0}^{\tau} p_{t+\tau} 8p^2 \gamma^4 \eta^2 (\delta^2 + 1) \frac{d^3}{m^3} \tau^3 L^2 \left(\frac{2p^2 d}{m\beta} + 5 \right) \sum_{r=t-\tau}^{t-1} \mathbb{E} \|\nabla f(x_g^r)\| \\
& + \sum_{t=0}^{\tau} p_{t+\tau} \varepsilon \gamma \eta L \left(3p \frac{d}{m} \sqrt{\delta^2 + 1} + 2\gamma\eta L \right) \mathbb{E} [\|x^t - x^*\|^2] \\
& + \sum_{t=0}^{\tau} p_{t+\tau} \varepsilon \gamma \eta L \left(3p \frac{d}{m} \sqrt{\delta^2 + 1} + 2\gamma\eta L \right) \frac{2}{\mu} \mathbb{E}[f(x_f^t) - f(x^*)] \\
& + \sum_{t=\tau}^T p_t 2p\gamma\eta \left(\frac{\varepsilon d}{m\sqrt{\delta^2 + 1}L} + 4p\gamma\eta \frac{d^2}{m^2} + 23p\gamma^3\eta\tau^4 \frac{d^3}{m^3} L^2 \right. \\
& \quad \left. + p\gamma\eta\tau^2 \frac{d^2}{m^2} \left(16 \frac{\tau^2 d^2 L^2 p^2 \gamma^2}{m^2 \beta} + \frac{5}{8} \right) \right) \sigma^2.
\end{aligned}$$

Taking

$$\begin{aligned}
& \gamma \leq \frac{\mu^{\frac{1}{3}} m^{\frac{1}{2}}}{2\tau L^{\frac{4}{3}} d^{\frac{1}{2}}}, \quad p \leq \frac{m^2}{13d^2(\delta^2 + 1)\tau^2}, \\
& \varepsilon \leq \min \left\{ \frac{m^{\frac{7}{4}}}{6d^{\frac{7}{4}}\tau^{\frac{5}{4}}L(\delta^2 + 1)}; \frac{m^{\frac{5}{4}}}{\sqrt{2}\tau^{\frac{3}{4}}\mu^{\frac{1}{3}}L^{\frac{2}{3}}d^{\frac{5}{4}}}; \frac{m^{\frac{15}{4}}}{6d^{\frac{15}{4}}\tau^{\frac{13}{4}}(\delta^2 + 1)^2} \right\},
\end{aligned}$$

we get

$$\begin{aligned}
& 10 \frac{d^2}{m^2} (\delta^2 + 1) - \frac{1}{p} + \tau^2 (\delta^2 + 1) \frac{d^2}{m^2} \left(32 \frac{\tau^2 d^2 L^2 p^2 \gamma^2}{m^2 \beta} + \frac{5}{4} \right) \leq 0, \\
& 6\varepsilon\gamma\eta L \left(3p \frac{d}{m} \sqrt{\delta^2 + 1} + 2\gamma\eta L \right) \leq \frac{\beta}{4}, \\
& 12 \frac{\varepsilon\gamma\eta L}{\mu} \left(3p \frac{d}{m} \sqrt{\delta^2 + 1} + 2\gamma\eta L \right) \leq 2p\gamma\eta^2 \frac{p}{2\eta},
\end{aligned}$$

and therefore with $\beta = \frac{p}{\eta}$

$$\begin{aligned}
& \sum_{t=\tau}^T p_t \mathbb{E} [\|x^{t+1} - x^*\|^2 + 2\gamma\eta^2 (f(x_f^{t+1}) - f(x^*))] \\
& \leq \sum_{t=\tau}^T p_t \left(1 - \frac{\beta}{2} \right) \mathbb{E} [\|x^t - x^*\|^2 + 2\gamma\eta^2 (f(x_f^t) - f(x^*))] \\
& \quad + \sum_{t=0}^{\tau} p_{t+\tau} 8p^2 \gamma^4 \eta^2 (\delta^2 + 1) \frac{d^3}{m^3} \tau^3 L^2 \left(\frac{2p^2 d}{m\beta} + 5 \right) \sum_{r=t-\tau}^{t-1} \mathbb{E} \|\nabla f(x_g^r)\| \\
& \quad + \sum_{t=0}^{\tau} p_{t+\tau} \varepsilon \gamma \eta L \left(3p \frac{d}{m} \sqrt{\delta^2 + 1} + 2\gamma\eta L \right) \mathbb{E} [\|x^t - x^*\|^2] \\
& \quad + \sum_{t=0}^{\tau} p_{t+\tau} \varepsilon \gamma \eta L \left(3p \frac{d}{m} \sqrt{\delta^2 + 1} + 2\gamma\eta L \right) \frac{2}{\mu} \mathbb{E}[f(x_f^t) - f(x^*)] \\
& \quad + \sum_{t=\tau}^T p_t 2p\gamma\eta \left(\frac{\varepsilon d}{m\sqrt{\delta^2 + 1}L} + 4p\gamma\eta \frac{d^2}{m^2} \right)
\end{aligned}$$

2268
2269
2270

$$+ 23p\gamma^3\eta\tau^4 \frac{d^3}{m^3} L^2 + p\gamma\eta\tau \frac{d^2}{m^2} \left(16 \frac{\tau^2 d^2 L^2 p^2 \gamma^2}{m^2 \beta} + \frac{5}{8} \right) \sigma^2.$$

2271
2272

Assume the following notation

2273
2274
2275
2276
2277
2278
2279
2280

$$\begin{aligned} \Delta_\tau &:= \sum_{t=0}^{\tau} p_{t+\tau} 8p^2 \gamma^4 \eta^2 (\delta^2 + 1) \frac{d^3}{m^3} \tau^3 L^2 \left(\frac{2p^2 d}{m\beta} + 5 \right) \sum_{r=t-\tau}^{t-1} \mathbb{E} \|\nabla f(x_g^r)\| \\ &\quad + \sum_{t=0}^{\tau} p_{t+\tau} \varepsilon \gamma \eta L \left(3p \frac{d}{m} \sqrt{\delta^2 + 1} + 2\gamma \eta L \right) \mathbb{E} [\|x^t - x^*\|^2] \\ &\quad + \sum_{t=0}^{\tau} p_{t+\tau} \varepsilon \gamma \eta L \left(3p \frac{d}{m} \sqrt{\delta^2 + 1} + 2\gamma \eta L \right) \frac{2}{\mu} \mathbb{E}[f(x_f^t) - f(x^*)] \\ &\leq \frac{\sqrt{\gamma}}{\tau^{\frac{4}{3}} \mu^{\frac{1}{3}}} \sum_{t=0}^{\tau} \left(\mathbb{E} \|\nabla f(x_g^t)\| + \mathbb{E} \|x^t - x^*\|^2 + \mathbb{E}[f(x_f^t) - f(x^*)] \right) \end{aligned}$$

2281
2282
2283

Now we substitute p_t , this lead us to

2284
2285
2286
2287
2288

$$\begin{aligned} &\sum_{t=\tau}^T \left(1 - \frac{\beta}{2} \right)^{-t} \mathbb{E} [\|x^{t+1} - x^*\|^2 + 2\gamma \eta^2 (f(x_f^{t+1}) - f(x^*))] \\ &\leq \sum_{t=\tau}^T \left(1 - \frac{\beta}{2} \right)^{-t+1} \mathbb{E} [\|x^t - x^*\|^2 + 2\gamma \eta^2 (f(x_f^t) - f(x^*))] + \Delta_\tau \\ &\quad + \sum_{t=\tau}^T \left(1 - \frac{\beta}{2} \right)^{-t} 2p\gamma\eta \left(\frac{\varepsilon d}{m\sqrt{\delta^2 + 1}L} + 4p\gamma\eta \frac{d^2}{m^2} \right. \\ &\quad \left. + 23p\gamma^3\eta\tau^4 \frac{d^3}{m^3} L^2 + p\gamma\eta\tau \frac{d^2}{m^2} \left(16 \frac{\tau^2 d^2 L^2 p^2 \gamma^2}{m^2 \beta} + \frac{5}{8} \right) \right) \sigma^2. \end{aligned}$$

2289
2290
2291

This implies

2292
2293
2294

$$\begin{aligned} \left(1 - \frac{\beta}{2} \right)^{-T} \mathbb{E} [\|x^{T+1} - x^*\|^2 + 2\gamma \eta^2 (f(x_f^{T+1}) - f(x^*))] &\leq \left(1 - \frac{\beta}{2} \right)^{\tau} \mathbb{E} [\|x^{\tau} - x^*\|^2 \\ &\quad + 2\gamma \eta^2 (f(x_f^{\tau}) - f(x^*))] + \Delta_\tau \\ &\quad + \sum_{t=\tau}^T \left(1 - \frac{\beta}{2} \right)^{-t} 2p\gamma\eta \left(\frac{\varepsilon d}{m\sqrt{\delta^2 + 1}L} + 4p\gamma\eta \frac{d^2}{m^2} \right. \\ &\quad \left. + 23p\gamma^3\eta\tau^4 \frac{d^3}{m^3} L^2 + p\gamma\eta\tau \frac{d^2}{m^2} \left(16 \frac{\tau^2 d^2 L^2 p^2 \gamma^2}{m^2 \beta} + \frac{5}{8} \right) \right) \sigma^2. \end{aligned}$$

2301
2302
2303
2304

2305
2306
2307

Rearranging this inequality and taking $\varepsilon \leq \frac{\sqrt{\gamma}m}{\sqrt{\mu}d}$ we obtain

2310
2311
2312
2313

$$\begin{aligned} &\mathbb{E} [\|x^{T+1} - x^*\|^2 + 2\gamma \eta^2 (f(x_f^{T+1}) - f(x^*))] \\ &\leq \left(1 - \frac{\beta}{2} \right)^{T-\tau} \mathbb{E} [\|x^{\tau} - x^*\|^2 + 2\gamma \eta^2 (f(x_f^{\tau}) - f(x^*))] + \left(1 - \frac{\beta}{2} \right)^T \Delta_\tau + 6\sqrt{\frac{\gamma}{\mu}} \sigma^2. \end{aligned}$$

2314
2315

This finishes the proof. \square

2316
2317
2318

H EXPERIMENTS

2319
2320
2321

This section provides description of the experiment setup, presents and analyses results of logistic regression experiments on LIBSVM datasets, studies dependence of history size over convergence. Moreover, experiments with neural networks optimization for data-parallelism and model-parallelism are presented and discussed.

H.1 TECHNICAL DETAILS

Our implementation of compression operators and algorithms is written in Python 3.10, with the use of PyTorch optimization library. We implement a simulation of distributed optimization system on a single machine, which is equivalent in terms of convergence analysis. Our server is AMD Ryzen Threadripper 2950X 16-Core Processor @ 2.2 GHz CPU and x2 NVIDIA GeForce GTX 1080 Ti GPU. We use Weights&Biases Biewald (2020) for experiments tracking and hyperparameters tuning.

H.2 LOGISTIC REGRESSION EXPERIMENTS

We conduct experiments on classification with logistic regression on four datasets: Mushrooms, A9A, W8A, MNIST. We apply the following optimization algorithms: proposed MQSGD and its accelerated version AMQSGD, and also use Markovian compressors with popular DIANA Mishchenko et al. (2019) algorithm. In all of our experiments, we do not utilize the steps of the optimizer, but rather the information that is transmitted by each worker at the current timestamp t . This implies that there are n workers, with each worker sending m coordinates at each iteration of the optimization step. Consequently, the x -axis displays numbers of the form $mn \cdot 1, mn \cdot 2, \dots, mn \cdot t, \dots, mn \cdot T$. This allows us to understand the performance of compressors with varying values of m and n .

We use convex logistic regression loss with a regularization term $\lambda = 0.05$. Each dataset is split horizontally (by rows) equally between $N = 10$ clients. The feature dimension is denoted as d in the figures, varying from hundreds to almost a thousand between datasets. The underlying sparsification compressors in Rand-10% for all logistic regression experiments. Learning rate initial value and decay rate are fine-tuned for each problem and compressor. Additionally, Markovian-specific parameters such as history size K , forgetting rate b are also fine-tuned. Table 2 provides hyperparameters grid for the tuning. We obtain optimal solution x^* for each problem with `scipy.optimize` method in order to use this value for the graphics.

Table 2: Hyperparameters values used for tuning in the experiments.

Hyperparameters	Values List
Learning rate	[0.01, 0.03, 0.05, 0.1, 0.3, 0.5, 1]
Learning rate decay rate	[0.5, 0.8, 1]
History size K	[1 ... 40]
Forgetting rate B	[1, 10, 15, 20, 30, 50]

Figures 5, 6 and 7 present relative distance to the optimum and gradient norm for the best runs on MQSGD, AMQSGD and DIANA, respectively. We observe that Markovian compressors consistently outperform the Rand-10% baseline in all scenarios, as the diverging trend can be seen. Only in some experiments with DIANA (MNIST) the advantage is negligible although present. We also observe that simpler and computational-effective BanLast compressor is often enough to achieve substantial convergence improvement. Notably, fine-tuned hyperparameters are similar across datasets and algorithms: for example, BanLast tends to perform best with largest possible values of history size K , and KAWASAKI forgetting rate b is large. Notice that BanLast compressor with largest K turns into round-robin compressor with (almost) no stochasticity in coordinates choice.

2376
2377
2378
2379
2380
2381
2382
2383
2384
2385

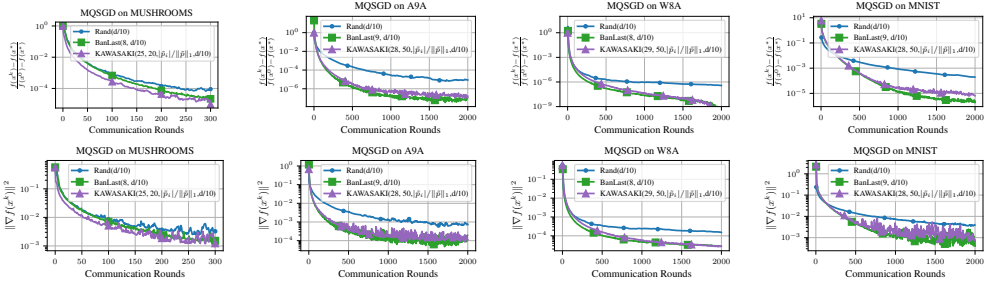


Figure 5: MQSGD LIBSVM logistic regression experiments. Best run after hyperparameters tuning is displayed for each method.

2386
2387
2388
2389
2390
2391
2392
2393
2394
2395
2396
2397
2398
2399

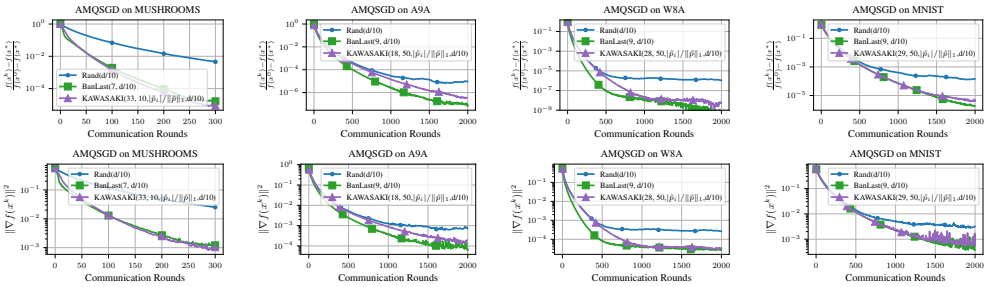


Figure 6: AMQSGD LIBSVM logistic regression experiments. Best run after hyperparameters tuning is displayed for each method.

2400
2401
2402
2403
2404
2405
2406
2407
2408
2409
2410
2411
2412

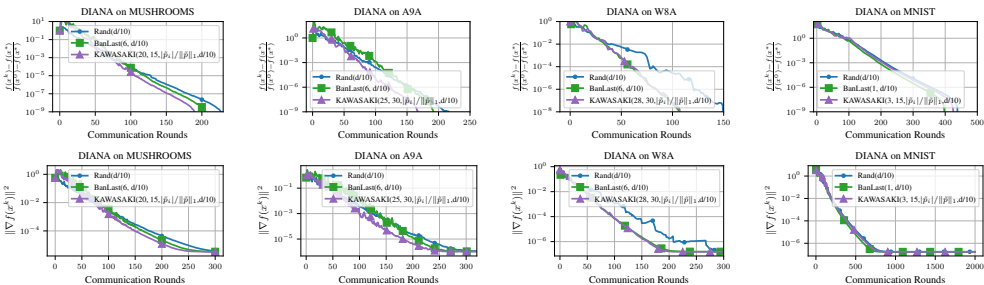


Figure 7: DIANA LIBSVM logistic regression experiments. Best run after hyperparameters tuning is displayed for each method.

2413
2414
2415
2416
2417

H.3 DEPENDENCE ON SIZE HISTORY

2418
2419
2420
2421
2422
2423
2424
2425
2426

As a part of hyperparameter tuning, we additionally analyze how history size K affects the convergence of Markovian compression-based methods. Figure 8 presents dependence of distance to optimum metric on history size for logistic regression experiments. We observe that BanLast performs better around larger values of $K = 8$ or $K = 9$. In such case for Rand10% used along with BanLast(9), the compression procedure resembles a permutation: for each 10 iterations, no indices are repeated, and the transmission cycle repeats after that. KAWASAKI history size seems to have periodical spikes and drops, achieving minimum at around $K = 25$. However, statistics for DIANA differ drastically, indicating that history size should be adjusted for each problem independently.

2427
2428
2429

H.4 COMPARISON WITH PERMUTATION & NATURAL COMPRESSION

In this section, we provide empirical comparison of the proposed compressors with other complex compression schemes.

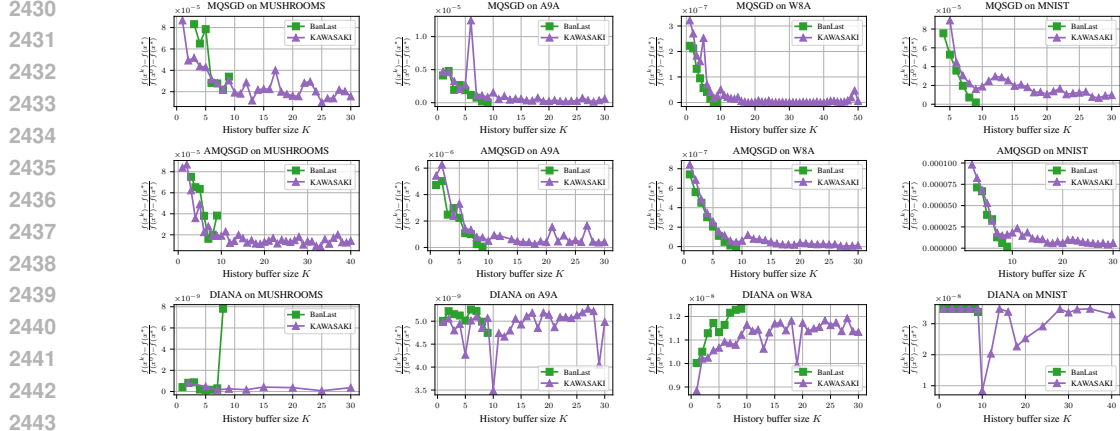


Figure 8: Convergence of Markovian-based algorithms on history size K

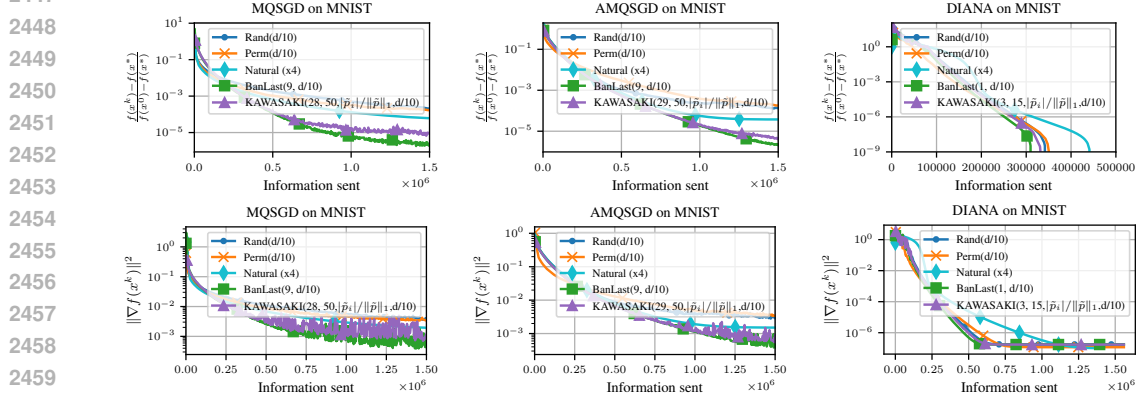


Figure 9: Comparison with PermK compressor and Natural compression. PermK compression factor is 10, Natural compression factor is 4. Logistic regression with L2 regularization on MNIST dataset for MQSGD, AMQSGD and DIANA algorithms on $N = 5$ clients. Best run is shown after fine-tuning learning rate, its decay, and Markovian compression parameters. X axis represent amount of information communicated.

Markovian compressors proposed in the paper compress vector coordinates dependently over optimization epochs. A similar idea of distributed compression is proposed in PermK Szlendak et al. (2021), where coordinates are arranged between workers at each iteration. Another compressor in the consideration is Natural compression Horvath et al. (2022), an unbiased randomized compressor.

Results of comparison of these compressors on MNIST dataset are presented in Figure 9. The results justify that Markovian compressors tend to converge faster than the competitors, allowing larger learning rates.

H.5 COMBINATION WITH OTHER COMPRESSORS

Although markovian compressors are initially targeted to work with sparsification-based compressors, refining coordinates selection probabilities, they are fully compatible with other compressors afterwards. To illustrate this, and to conduct additional comparison with PermK compressor, we setup experiments combined with Natural Compression. Precisely, we compare RandK+Natural, PermK+Natural, BanLast+Natural and KAWASAKI+Natural compressors on logistic regression on MNIST dataset.

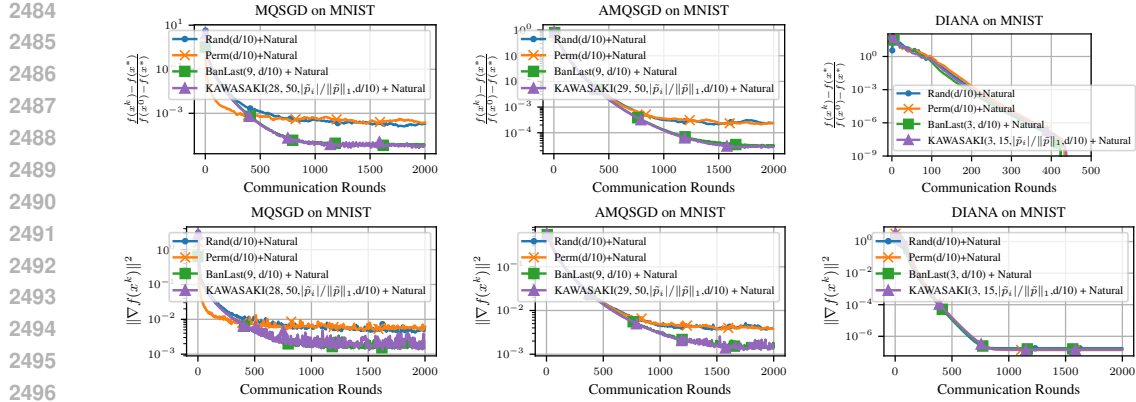


Figure 10: Experiments with Natural compression, MNIST logistic regression experiments. Best run after hyperparameters tuning is displayed for each method.

Figure 10 shows results of combination of mentioned sparsification compressors with natural compression.

H.6 NEURAL NETWORKS EXPERIMENTS: DATA PARALLELISM CASE

To adopt Markovian compression to a more complex task, we perform image classification on CIFAR-10 Krizhevsky et al. (2009) with Resnet-18 He et al. (2016) convolutional neural network. We split the training set of size 50,000 equally between $N = 5$ clients. We use SGD optimizer with momentum 0.9 and weight decay $5 \cdot 10^{-4}$. Hyperparameters such as batch size and learning rate are fine-tuned. Markovian compressors hyperparameters, such as history size K and forgetting rate b are fine-tuned, while activation function is set to ordinary normalization. Experiments are conducted with several sparsification compressors, such as Rand-5%, Rand-7%, and Rand-10%, with number of epochs adjusted for each case.

Figures 11, 12 and 13 present train loss, gradient norm and test accuracy for each baseline method and Markovian compressors for Rand-5%, Rand-7% and Rand-10% scenarios, respectively. Summary on best test accuracy is presented in Table 3, and extended numerical results for Rand-5% compressor were presented in main experiments Table 1. We observe that in such complex, batched optimization problem only KAWASAKI obtains a substantial convergence improvement, as opposed to simpler logistic regression. Nevertheless, BanLast still performs the best when used with large history size, while both history size and forgetting rate are low for KAWASAKI. In terms of achieved test set accuracy, methods differ significantly only on higher compression rates like Rand-5%. This may imply that Markovian compression tolerates stronger compression, which is useful in practice. To summarize, Markovian compressors can be successfully applied in neural networks training, with KAWASAKI compressor significantly improving convergence.

Finally, we also conduct the comparison with Permutatio and Natural compression, both independently and in combination. Figure 14 shows learning curves for training with $N = 20$ clients. KAWASAKI compressor appears to have best convergence in both independently and in combination with Natural compression against Permutation compressor.

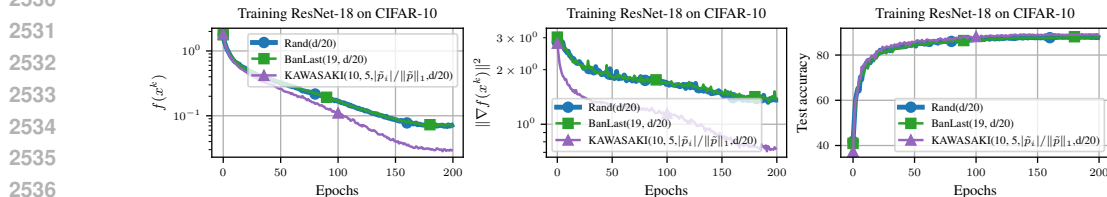


Figure 11: Resnet-18 on CIFAR-10 training results for Rand-5% sparsification.

2538
2539
2540
2541
2542
2543
2544

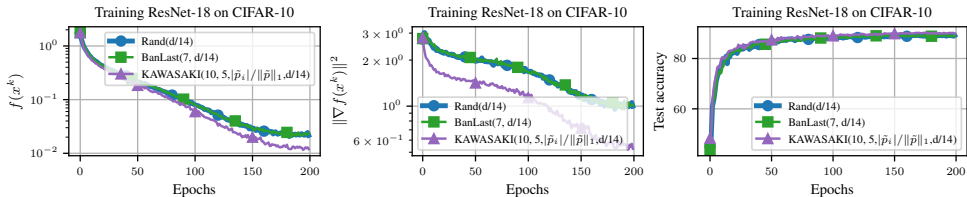


Figure 12: Resnet-18 on CIFAR-10 training results for Rand-7% sparsification.

2545
2546
2547
2548
2549
2550
2551
2552
2553
2554

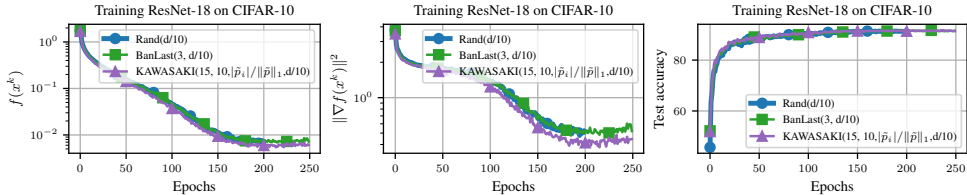


Figure 13: Resnet-18 on CIFAR-10 training results for Rand-10% sparsification.

2555
2556
2557
2558
2559
2560
2561
2562
2563
2564
2565
2566
2567

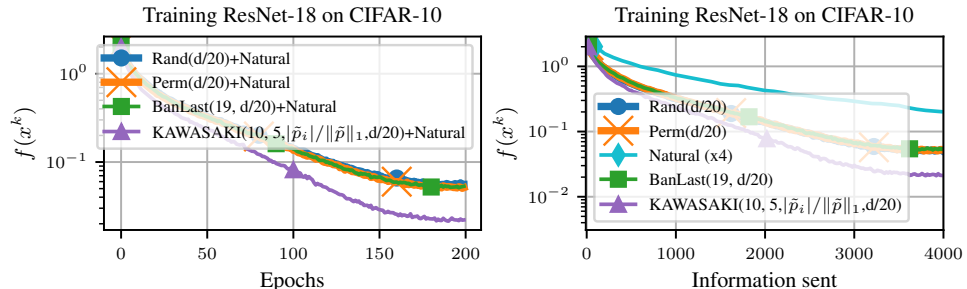


Figure 14: Comparison with other compressors on Resnet-18 training on CIFAR-10 dataset for Rand-5% sparsification on $N = 20$ clients. Natural compression factor is 4. Left figure is sequential combination with Natural compression. Right figure is comparison against PermK and Natural compressors independently, with information sent on x-axis.

2571
2572
2573
2574
2575
2576
2577
2578
2579

Table 3: Best test accuracy % of training ResNet-18 on CIFAR-10 with different compressors

	Rand-K%	Banlast	KAWASAKI
	88.03	88.1	89.27
	89.31	89.38	90.28
	91.46	91.72	91.78

H.7 NEURAL NETWORKS EXPERIMENTS: MODEL PARALLELISM CASE

2582
2583
2584
2585
2586
2587
2588

As opposed to data-parallel setting, model parallelism is paradigm which splits the model (typically a deep neural network) to a pipeline of layers between workers. Such distributed scenario is especially relevant for large language models (LLM), which consist of billions of trainable parameters. As communication is a typical bottleneck in such systems Diskin et al. (2021), various compression techniques are applied to layer activations and their respective gradients that are transferred between adjacent pipeline workers. Such techniques include quantization and sparsification Dettmers et al. (2022); Bian et al. (2023), as well as low-rank compression Song et al. (2023) techniques.

2589
2590
2591

We perform training of Resnet-18 He et al. (2016) convolutional neural network on CIFAR-10 dataset Krizhevsky et al. (2009). We split the ResNet onto 4 workers by resnet blocks, simulated on a single device with compression of activations and their respective gradients in the places of communication. We apply Markovian compressors only to gradients in model-parallel setup, using

same RandK compression for both activations and gradients independently for each compression block.

Table 4: Best test accuracy % for model parallelism experiments with Resnet-18 classification of CIFAR-10

Compressor	Compression ON	Compression OFF
No compression	92.8	92.8
Rand10%	84.6	86.1
BanLastK+Rand10%	85.2	86.4
KAWASAKI(simplex projection)+Rand10%	84.5	85.0
KAWASAKI(normalize)+Rand10%	85.2	86.8
KAWASAKI(softmax)+Rand10%	85.3	87.3

Table 4 presents best test set accuracy achieved for training with different compressors. While compression indeed decreases accuracy for Rand-10%, application of Markov compressors, especially KAWASAKI with normalization and softmax activation functions, favours the final test accuracy on a whole one percent. Note that compression is not applied during inference, only on training phase. This case illustrates potential of Markov compressors beyond data-parallelism setup considered in theory. In practical training of large neural networks, where both data-parallelism and model-parallelism are often applied simultaneously, Markov compressors could also be useful, as per shown efficiency on both these setups in separate.

H.8 FINE-TUNING DEBERTAV3-BASE ON GLUE DEVELOPMENT SET

In this series of experiments, we examine a distributed approach to fine-tuning language models using LoRA (Hu et al., 2021). This method is based on freezing the model weights that are pre-trained on a large dataset, and add a low rank adapter with matrices $A \in \mathbb{R}^{n \times r}$ and $B \in \mathbb{R}^{r \times m}$ to some selected layers $W_{old} \in \mathbb{R}^{n \times m}$ of this model, such that $W_{new} = W_{old} + A \cdot B$. Since in practice the parameter r is chosen to be much smaller than n and m , the new model has much fewer trainable parameters and can be efficiently trained on downstream tasks.

In our experiments, we apply LoRA adapters with fixed rank $r = 8$ to the attention layers of the DeBERTaV3-base model (He et al., 2021). The downstream task is the classical GLUE benchmark for natural language understanding (Wang et al., 2019). We consider only random sparsification compressors (Definition 4) with 25% compression rate, due to the large computational cost of this experiment. Figure 15 shows learning curves for training with $N = 10$ clients. Our Markovian compressors appears to have best convergence against independent Rand m compressor.

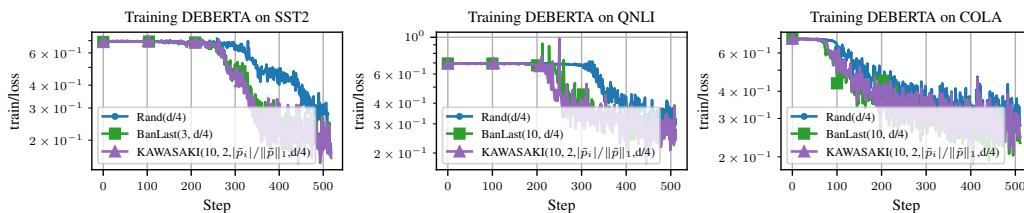


Figure 15: Comparison with other compressors on fine-tuning task on GLUE benchmark on $N = 10$ clients. We performed experiments on SST2, QNLI and COLA tasks, they are arranged from left to right.