

Singapore Management University

## Institutional Knowledge at Singapore Management University

---

Research Collection School Of Computing and Information Systems

School of Computing and Information Systems

---

7-2021

### MARINA: Faster non-convex distributed learning with compression

Eduard GORBUNOV

Konstantin BURLACHENKO

Zhize LI

Singapore Management University, zhizeli@smu.edu.sg

Peter RICHTARIK

Follow this and additional works at: [https://ink.library.smu.edu.sg/sis\\_research](https://ink.library.smu.edu.sg/sis_research)



Part of the [Databases and Information Systems Commons](#)

---

#### Citation

GORBUNOV, Eduard; BURLACHENKO, Konstantin; LI, Zhize; and RICHTARIK, Peter. MARINA: Faster non-convex distributed learning with compression. (2021). *Proceedings of the 38th International Conference on Machine Learning (ICML 2021), Virtual Conference, July 18-24*. 1-41.

Available at: [https://ink.library.smu.edu.sg/sis\\_research/8682](https://ink.library.smu.edu.sg/sis_research/8682)

This Conference Proceeding Article is brought to you for free and open access by the School of Computing and Information Systems at Institutional Knowledge at Singapore Management University. It has been accepted for inclusion in Research Collection School Of Computing and Information Systems by an authorized administrator of Institutional Knowledge at Singapore Management University. For more information, please email [cherylds@smu.edu.sg](mailto:cherylds@smu.edu.sg).

---

# MARINA: Faster Non-Convex Distributed Learning with Compression

---

Eduard Gorbunov<sup>1,2,3</sup> Konstantin Burlachenko<sup>3</sup> Zhize Li<sup>3</sup> Peter Richtárik<sup>3</sup>

## Abstract

We develop and analyze MARINA: a new communication efficient method for non-convex distributed learning over heterogeneous datasets. MARINA employs a novel communication compression strategy based on the compression of gradient differences that is reminiscent of but different from the strategy employed in the DIANA method of Mishchenko et al. (2019). Unlike virtually all competing distributed first-order methods, including DIANA, ours is based on a carefully designed *biased* gradient estimator, which is the key to its superior theoretical and practical performance. The communication complexity bounds we prove for MARINA are evidently better than those of all previous first-order methods. Further, we develop and analyze two variants of MARINA: VR-MARINA and PP-MARINA. The first method is designed for the case when the local loss functions owned by clients are either of a finite sum or of an expectation form, and the second method allows for a partial participation of clients – a feature important in federated learning. All our methods are superior to previous state-of-the-art methods in terms of oracle/communication complexity. Finally, we provide a convergence analysis of all methods for problems satisfying the Polyak-Łojasiewicz condition.

## 1. Introduction

Non-convex optimization problems appear in various applications of machine learning, such as training deep neural networks (Goodfellow et al., 2016) and matrix completion and recovery (Ma et al., 2018; Bhojanapalli et al., 2016). Because of their practical importance, these problems gained much attention in recent years, which led to a rapid develop-

ment of new efficient methods for non-convex optimization problems (Danilova et al., 2020), and especially the training of deep learning models (Sun, 2019).

Training deep neural networks is notoriously computationally challenging and time-consuming. In the quest to improve the generalization performance of modern deep learning models, practitioners resort to using increasingly larger datasets in the training process, and to support such workloads, it is imperative to use advanced parallel and distributed hardware, systems, and algorithms. Distributed computing is often necessitated by the desire to train models from data naturally distributed across several edge devices, as is the case in federated learning (Konečný et al., 2016; McMahan et al., 2017). However, even when this is not the case, distributed methods are often very efficient at reducing the training time (Goyal et al., 2017; You et al., 2020). Due to these and other reasons, distributed optimization has gained immense popularity in recent years.

However, distributed methods almost invariably suffer from the so-called *communication bottleneck*: the communication cost of information necessary for the workers to jointly solve the problem at hand is often very high, and depending on the particular compute architecture, workload, and algorithm used, it can be orders of magnitude higher than the computation cost. A popular technique for resolving this issue is *communication compression* (Seide et al., 2014; Konečný et al., 2016; Suresh et al., 2017), which is based on applying a lossy transformation/compression to the models, gradients, or tensors to be sent over the network to save on communication. Since applying a lossy compression generally decreases the utility of the exchanged messages, such an approach will typically lead to an increase in the number of communications, and the overall usefulness of this technique manifests itself in situations where the communication savings are larger compared to the increased need for the number of communication rounds (Horváth et al., 2019).

The optimization and machine learning communities have exerted considerable effort in recent years to design distributed methods supporting compressed communication. From many methods proposed, we emphasize VR-DIANA (Horváth et al., 2019), FedCOMGATE (Haddadpour et al., 2020), and FedSTEPH (Das et al., 2020) because these pa-

---

<sup>1</sup>Moscow Institute of Physics and Technology, Moscow, Russia <sup>2</sup>Yandex, Moscow, Russia <sup>3</sup>King Abdullah University of Science and Technology, Thuwal, Saudi Arabia. Correspondence to: Eduard Gorbunov <eduard.gorbunov@phystech.edu>, Peter Richtárik <peter.richtarik@kaust.edu.sa>.

pers contain the state-of-the-art results in the setup when the local loss functions can be arbitrary heterogeneous.

## 1.1. Contributions

We propose several new distributed optimization methods supporting compressed communication, specifically focusing on smooth but nonconvex problems of the form

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

where  $n$  workers/devices/clients/peers are connected in a centralized way with a parameter-server, and client  $i$  has an access to the local loss function  $f_i$  only. We establish strong complexity rates for them and show that they are better than previous state-of-the-art results.

- **MARINA.** The main contribution of our paper is a new distributed method supporting communication compression called MARINA (Alg 1). In this algorithm, workers apply an unbiased compression operator to the *gradient differences* at each iteration with some probability and send them to the server that performs aggregation by averaging. Unlike all known methods operating with unbiased compression operators, this procedure leads to a *biased* gradient estimator. We prove convergence guarantees for MARINA, which are strictly better than previous state-of-the-art methods (see Table 1). For example, MARINA’s rate  $\mathcal{O}(\frac{1+\omega/\sqrt{n}}{\varepsilon^2})$  is  $\mathcal{O}(\sqrt{\omega})$  times better than that of the state-of-the-art method DIANA (Mishchenko et al., 2019), where  $\omega$  is the variance parameter associated with the deployed compressor. For example, in the case of the Rand1 sparsification compressor, we have  $\omega = d - 1$ , and hence we get an improvement by the factor  $\mathcal{O}(\sqrt{d})$ . Since the number  $d$  of features can be truly very large when training modern models, this is a substantial improvement that can even amount to *several orders of magnitude*.

- **Variance Reduction on Nodes.** We generalize MARINA to VR-MARINA, which can handle the situation when the local functions  $f_i$  have either a finite-sum (each  $f_i$  is an average of  $m$  functions) or an expectation form, and when it is more efficient to rely on local stochastic gradients rather than on local gradients. When compared with MARINA, VR-MARINA additionally performs *local variance reduction* on all nodes, progressively removing the variance coming from the stochastic approximation, leading to a better oracle complexity than previous state-of-the-art results (see Table 1). When no compression is used (i.e.,  $\omega = 0$ ), the rate of VR-MARINA is  $\mathcal{O}(\frac{\sqrt{m}}{\sqrt{n\varepsilon^2}})$ , while the rate of the state-of-the-art method VR-DIANA is  $\mathcal{O}(\frac{m^{2/3}}{\varepsilon^2})$ . This is an improvement by the factor  $\mathcal{O}(\sqrt{nm}^{1/6})$ . When much compression is applied, and  $\omega$  is large, our method is faster by the factor  $\mathcal{O}(\frac{m^{2/3}+\omega}{m^{1/2}+\omega^{1/2}})$ . In the special case, when there is just a sin-

gle node ( $n = 1$ ), and no compression is used, VR-MARINA reduces to the PAGE method of Li et al. (2020); this is an optimal first-order algorithm for smooth non-convex finite-sum/online optimization problems.

- **Partial Participation.** We develop a modification of MARINA allowing for *partial participation* of the clients, which is a feature critical in federated learning. The resulting method, PP-MARINA, has superior communication complexity to the existing methods developed for this settings (see Table 1).

- **Convergence Under the Polyak-Łojasiewicz Condition.** We analyze all proposed methods for problems satisfying the Polyak-Łojasiewicz condition (Polyak, 1963; Łojasiewicz, 1963). Again, the obtained results are strictly better than previous ones (see Table 2). Statements and proofs of all these results are in the Appendix.

- **Simple Analysis.** The simplicity and flexibility of our analysis offer several extensions. For example, one can easily generalize our analysis to the case of different quantization operators and different batch sizes used by clients. Moreover, one can combine the ideas of VR-MARINA and PP-MARINA and obtain a single distributed algorithm with compressed communications, variance reduction on nodes, and clients’ sampling. We did not do this to keep the exposition simpler.

## 1.2. Related Work

**Non-Convex Optimization.** Since finding a global minimum of a non-convex function is, in general, an NP-hard problem (Murty & Kabadi, 1987), many researchers in non-convex optimization focus on relaxed goals such as finding an  $\varepsilon$ -stationary point. The theory of stochastic first-order methods for finding  $\varepsilon$ -stationary points is well-developed: it contains lower bounds for expectation minimization without smoothness of stochastic realizations (Arjevani et al., 2019) and for finite-sum/expectation minimization (Fang et al., 2018; Li et al., 2020) as well as optimal methods matching the lower bounds (see (Danilova et al., 2020; Li et al., 2020) for the overview). Recently, distributed variants of such methods were proposed (Sun et al., 2020; Sharma et al., 2019; Khanduri et al., 2020).

**Compressed Communications.** Works on distributed methods supporting communication compression can be roughly split into two large groups: the first group focuses on methods using *unbiased* compression operators (which refer to as quantizations in this paper), such as RandK, and the second one studies methods using *biased* compressors such as TopK. One can find a detailed summary of the most popular compression operators in (Safaryan et al., 2020; Beznosikov et al., 2020).

**Unbiased Compression.** In this line of work, the first con-

Table 1: Summary of the state-of-the-art results for finding an  $\varepsilon$ -stationary point for the problem (1), i.e., such a point  $\hat{x}$  that  $\mathbb{E} [\|\nabla f(\hat{x})\|^2] \leq \varepsilon^2$ . Dependences on the numerical constants, “quality” of the starting point, and smoothness constants are omitted in the complexity bounds. Abbreviations: “PP” = partial participation; “Communication complexity” = the number of communications rounds needed to find an  $\varepsilon$ -stationary point; “Oracle complexity” = the number of (stochastic) first-order oracle calls needed to find an  $\varepsilon$ -stationary point. Notation:  $\omega$  = the quantization parameter (see Def. 1.1);  $n$  = the number of nodes;  $m$  = the size of the local dataset;  $r$  = (expected) number of clients sampled at each iteration;  $b'$  = the batchsize for VR-MARINA at the iterations with compressed communication. To simplify the bounds, we assume that the expected density  $\zeta_Q$  of the quantization operator  $Q$  (see Def. 1.1) satisfies  $\omega + 1 = \Theta(d/\zeta_Q)$  (e.g., this holds for RandK and  $\ell_2$ -quantization, see (Beznosikov et al., 2020)). We notice that (Haddadpour et al., 2020) and (Das et al., 2020) contain also better rates under different assumptions on clients’ similarity.

Setup	Method	Citation	Communication Complexity	Oracle Complexity
(1)	DIANA	(Mishchenko et al., 2019) (Horváth et al., 2019)	$\frac{1+(1+\omega)\sqrt{\omega/n}}{\varepsilon^2}$	$\frac{1+(1+\omega)\sqrt{\omega/n}}{\varepsilon^2}$
	FedCOMGATE <sup>(1)</sup>	(Li & Richtárik, 2020)	$\frac{1+\omega}{\varepsilon^2}$	$\frac{1+\omega}{n\varepsilon^4}$
	FedSTEPH, $r = n$	(Haddadpour et al., 2020)	$\frac{1+\omega/n}{\varepsilon^4}$	$\frac{1+\omega/n}{\varepsilon^4}$
	MARINA (Alg. 1)	Thm. 2.1 & Cor. 2.1 (NEW)	$\frac{1+\omega/\sqrt{n}}{\varepsilon^2}$	$\frac{1+\omega/\sqrt{n}}{\varepsilon^2}$
(1)+(5)	DIANA	(Li & Richtárik, 2020)	$\frac{1+(1+\omega)\sqrt{\omega/n}}{\varepsilon^2} + \frac{1+\omega}{n\varepsilon^4}$	$\frac{1+(1+\omega)\sqrt{\omega/n}}{\varepsilon^2} + \frac{1+\omega}{n\varepsilon^4}$
	VR-DIANA	(Horváth et al., 2019)	$\frac{(m^{2/3}+\omega)\sqrt{1+\omega/n}}{\varepsilon^2}$	$\frac{(m^{2/3}+\omega)\sqrt{1+\omega/n}}{\varepsilon^2}$
	VR-MARINA (Alg. 2), $b' = 1$ <sup>(2)</sup>	Thm. 3.1 & Cor. 3.1 (NEW)	$\frac{1+\max\{\omega, \sqrt{(1+\omega)m}\}/\sqrt{n}}{\varepsilon^2}$	$\frac{1+\max\{\omega, \sqrt{(1+\omega)m}\}/\sqrt{n}}{\varepsilon^2}$
(1)+(6)	DIANA <sup>(3)</sup>	(Mishchenko et al., 2019) (Li & Richtárik, 2020)	$\frac{1+(1+\omega)\sqrt{\omega/n}}{\varepsilon^2} + \frac{1+\omega}{n\varepsilon^4}$	$\frac{1+(1+\omega)\sqrt{\omega/n}}{\varepsilon^2} + \frac{1+\omega}{n\varepsilon^4}$
	FedCOMGATE <sup>(3)</sup>	(Haddadpour et al., 2020)	$\frac{1+\omega}{\varepsilon^2}$	$\frac{1+\omega}{n\varepsilon^4}$
	VR-MARINA (Alg. 2), $b' = 1$	Thm. 3.2 & Cor. 3.2 (NEW)	$\frac{1+\omega/\sqrt{n} + \sqrt{1+\omega}}{\varepsilon^2 n\varepsilon^3}$	$\frac{1+\omega/\sqrt{n} + \sqrt{1+\omega}}{\varepsilon^2 n\varepsilon^3}$
	VR-MARINA (Alg. 2), $b' = \Theta\left(\frac{1}{n\varepsilon^2}\right)$	Thm. 3.2 & Cor. 3.2 (NEW)	$\frac{1+\omega/\sqrt{n}}{\varepsilon^2} + \frac{1+\omega}{n\varepsilon^3}$	$\frac{1+\omega/\sqrt{n}}{\varepsilon^2} + \frac{1+\omega}{n\varepsilon^3}$
PP, (1)	FedSTEPH	(Das et al., 2020)	$\frac{1+\omega/n}{r\varepsilon^4} + \frac{(1+\omega)(n-r)}{r(n-1)\varepsilon^4}$	$\frac{1+\omega/n}{r\varepsilon^4} + \frac{(1+\omega)(n-r)}{r(n-1)\varepsilon^4}$
	PP-MARINA (Alg. 4)	Thm. 4.1 & Cor. 4.1 (NEW)	$\frac{1+(1+\omega)\sqrt{n}/r}{\varepsilon^2}$	$\frac{1+(1+\omega)\sqrt{n}/r}{\varepsilon^2}$

<sup>(1)</sup> The results for FedCOMGATE are derived under assumption that for all vectors  $x_1, \dots, x_n \in \mathbb{R}^d$  the quantization operator  $Q$  satisfies  $\mathbb{E} \left[ \left\| \frac{1}{n} \sum_{i=1}^n Q(x_i) \right\|^2 - \left\| Q\left(\frac{1}{n} \sum_{i=1}^n x_i\right) \right\|^2 \right] \leq G$  for some constant  $G \geq 0$ . In fact, this assumption does not hold for classical quantization operators like RandK and  $\ell_2$ -quantization on  $\mathbb{R}^d$ . The counterexample:  $n = 2$  and  $x_1 = -x_2 = (t, t, \dots, t)^\top$  with arbitrary large  $t > 0$ .

<sup>(2)</sup> One can even further improve the communication complexity by increasing  $b'$ .

<sup>(3)</sup> No assumptions on the smoothness of the stochastic realizations  $f_\varepsilon(x)$  are used.

vergence result in the non-convex case was obtained by Alistarh et al. (2017) for QSGD, under assumptions that the local loss functions are the same for all workers, and the stochastic gradient has uniformly bounded second moment. After that, Mishchenko et al. (2019) proposed DIANA (and its momentum version) and proved its convergence rate for non-convex problems without any assumption on the boundedness of the second moment of the stochastic gradient, but under the assumption that the dissimilarity between local loss functions is bounded. This restriction was later eliminated by Horváth et al. (2019) for the variance reduced version of DIANA called VR-DIANA, and the analysis was extended to a large class of unbiased compressors. Finally, the results for QSGD and DIANA were recently generalized and tightened by Li & Richtárik (2020) in a unifying framework that included many other methods as well.

**Biased Compression.** Biased compression operators are less “optimization-friendly” than unbiased ones. Indeed, one can construct a simple convex quadratic problem for which distributed SGD with Top1 compression diverges exponentially fast (Beznosikov et al., 2020). However, this issue can be resolved using *error compensation* (Seide et al., 2014). The first analysis of error-compensated SGD (EC-

SGD) for non-convex problems was obtained by Karimireddy et al. (2019) for homogeneous problems under the assumption that the second moment of the stochastic gradient is uniformly bounded. The last assumption was recently removed from the analysis of EC-SGD by Stich & Karimireddy (2020); Beznosikov et al. (2020), while the first results without the homogeneity assumption were obtained by Koloskova et al. (2020a) for Choco-SGD, but still under the assumption that the second moment of the stochastic gradient is uniformly bounded. This issue was resolved by Beznosikov et al. (2020). In general, the current understanding of optimization methods with biased compressors is far from complete: even in the strongly convex case, the first linearly converging (Gorbunov et al., 2020) and accelerated (Qian et al., 2020) error-compensated stochastic methods were proposed just recently.

**Other Approaches.** Besides communication compression, there are also different techniques aiming to reduce the overall communication cost of distributed methods. The most popular ones are based on decentralized communications and multiple local steps between communication rounds, where the second technique is very popular in federated learning (Konečný et al., 2016; Kairouz et al., 2019).

Table 2: Summary of the state-of-the-art results for finding an  $\varepsilon$ -solution for the problem (1) satisfying **Polyak-Łojasiewicz condition** (see As. 2.1), i.e., such a point  $\hat{x}$  that  $\mathbf{E}[f(\hat{x}) - f(x^*)] \leq \varepsilon$ . Dependences on the numerical constants and  $\log(1/\varepsilon)$  factors are omitted and all smoothness constant are denoted by  $L$  in the complexity bounds. Abbreviations: “PP” = partial participation; “Communication complexity” = the number of communications rounds needed to find an  $\varepsilon$ -stationary point; “Oracle complexity” = the number of (stochastic) first-order oracle calls needed to find an  $\varepsilon$ -stationary point. Notation:  $\omega$  = the quantization parameter (see Def. 1.1);  $n$  = the number of nodes;  $m$  = the size of the local dataset;  $r$  = (expected) number of clients sampled at each iteration;  $b'$  = the batchsize for VR-MARINA at the iterations with compressed communication. To simplify the bounds, we assume that the expected density  $\zeta_{\mathcal{Q}}$  of the quantization operator  $\mathcal{Q}$  (see Def. 1.1) satisfies  $\omega + 1 = \Theta(d/\zeta_{\mathcal{Q}})$  (e.g., this holds for RandK and  $\ell_2$ -quantization, see (Beznosikov et al., 2020)). We notice that (Haddadpour et al., 2020) and (Das et al., 2020) contain also better rates under different assumptions on clients’ similarity.

Setup	Method	Citation	Communication Complexity	Oracle Complexity
(1)	DIANA	(Li & Richtárik, 2020)	$\frac{L(1+(1+\omega)\sqrt{\omega/n})}{\mu}$	$\frac{L(1+(1+\omega)\sqrt{\omega/n})}{\mu}$
	FedCOMGATE <sup>(1)</sup>	(Haddadpour et al., 2020)	$\frac{L(1+\omega)}{\mu}$	$\frac{L(1+\omega)}{\mu}$
	MARINA (Alg. 1)	Thm. 2.2 & Cor. C.2 (NEW)	$\omega + \frac{L(1+\omega/\sqrt{n})}{\mu}$	$\omega + \frac{L(1+\omega/\sqrt{n})}{\mu}$
(1)+(5)	DIANA	(Li & Richtárik, 2020)	$\frac{L(1+(1+\omega)\sqrt{\omega/n})}{\mu} + \frac{L(1+\omega)}{n\mu} \left( \frac{L}{\mu} + \frac{1}{\varepsilon} \right)$	$\frac{L(1+(1+\omega)\sqrt{\omega/n})}{\mu} + \frac{L(1+\omega)}{n\mu} \left( \frac{L}{\mu} + \frac{1}{\varepsilon} \right)$
	VR-DIANA	(Li & Richtárik, 2020)	$\frac{L(m^{2/3}+\omega)\sqrt{1+\omega/n}}{\mu}$	$\frac{L(m^{2/3}+\omega)\sqrt{1+\omega/n}}{\mu}$
	VR-MARINA (Alg. 2), $b' = 1$ <sup>(2)</sup>	Thm. D.2 & Cor. D.2 (NEW)	$\omega + m + \frac{L(1+\max\{\omega, \sqrt{(1+\omega)m}\}/\sqrt{n})}{\mu}$	$\omega + m + \frac{L(1+\max\{\omega, \sqrt{(1+\omega)m}\}/\sqrt{n})}{\mu}$
(1)+(6)	DIANA <sup>(3)</sup>	(Mishchenko et al., 2019) (Li & Richtárik, 2020)	$\frac{1+(1+\omega)\sqrt{\omega/n}}{\varepsilon^2} + \frac{1+\omega}{n\varepsilon^4}$	$\frac{1+(1+\omega)\sqrt{\omega/n}}{\varepsilon^2} + \frac{1+\omega}{n\varepsilon^4}$
	FedCOMGATE <sup>(3)</sup>	(Haddadpour et al., 2020)	$\frac{L(1+\omega)}{n\mu\varepsilon}$	$\frac{L(1+\omega)}{n\mu\varepsilon}$
	VR-MARINA (Alg. 2), $b' = 1$	Thm. D.4 & Cor. D.4 (NEW)	$\omega + \frac{L(1+\omega/\sqrt{n})}{\mu} + \frac{L\sqrt{1+\omega}}{n\mu\varepsilon}$	$\omega + \frac{L(1+\omega/\sqrt{n})}{\mu} + \frac{L\sqrt{1+\omega}}{n\mu\varepsilon}$
	VR-MARINA (Alg. 2), $b' = \Theta\left(\frac{1}{n\mu\varepsilon}\right)$	Thm. D.4 & Cor. D.4 (NEW)	$\omega + \frac{L(1+\omega/\sqrt{n})}{\mu}$	$\frac{1+\omega}{n\mu\varepsilon} + \frac{L(1+\omega/\sqrt{n})}{n\mu^2\varepsilon} + \frac{L(1+\omega)}{n\mu^2\sqrt{\varepsilon}}$
PP, (1)	FedSTEPH <sup>(4)</sup>	(Das et al., 2020)	$\left(\frac{L}{\mu}\right)^{3/2}$	$\left(\frac{L}{\mu}\right)^{3/2}$
	PP-MARINA (Alg. 4)	Thm. E.2 & Cor. E.2 (NEW)	$\frac{(\omega+1)n}{r} + \frac{L(1+(1+\omega)\sqrt{n}/r)}{\mu}$	$\frac{(\omega+1)n}{r} + \frac{L(1+(1+\omega)\sqrt{n}/r)}{\mu}$

<sup>(1)</sup> The results for FedCOMGATE are derived under assumption that for all vectors  $x_1, \dots, x_n \in \mathbb{R}^d$  the quantization operator  $\mathcal{Q}$  satisfies  $\mathbf{E} \left[ \left\| \frac{1}{n} \sum_{i=1}^n \mathcal{Q}(x_j) \right\|^2 - \left\| \mathcal{Q} \left( \frac{1}{n} \sum_{i=1}^n x_j \right) \right\|^2 \right] \leq G$  for some constant  $G \geq 0$ . In fact, this assumption does not hold for classical quantization operators like RandK and  $\ell_2$ -quantization on  $\mathbb{R}^d$ . The counterexample:  $n = 2$  and  $x_1 = -x_2 = (t, t, \dots, t)^\top$  with arbitrary large  $t > 0$ .

<sup>(2)</sup> One can even further improve the communication complexity by increasing  $b'$ .  
<sup>(3)</sup> No assumptions on the smoothness of the stochastic realizations  $f_\varepsilon(x)$  are used.

<sup>(4)</sup> The rate is derived under assumption that  $r = \Omega((1+\omega)\sqrt{L/\mu} \log(1/\varepsilon))$ .

One can find the state-of-the-art distributed optimization methods using these techniques and their combinations in (Lian et al., 2017; Karimireddy et al., 2020; Li et al., 2019; Koloskova et al., 2020b). Moreover, there exist results based on the combinations of communication compression with either decentralized communication, e.g., Choco-SGD (Koloskova et al., 2020a), or local updates, e.g., Qsparse-Local-SGD (Basu et al., 2019), FedCOMGATE (Haddadpour et al., 2020), FedSTEPH (Das et al., 2020), where in (Basu et al., 2019) the convergence rates were derived under an assumption that the stochastic gradient has uniformly bounded second moment and the results for Choco-SGD, FedCOMGATE, FedSTEPH were described either earlier in the text, or in Table 1.

### 1.3. Preliminaries

We will rely on two key assumptions throughout the text.

**Assumption 1.1** (Uniform lower bound). *There exists  $f_* \in \mathbb{R}$  such that  $f(x) \geq f_*$  for all  $x \in \mathbb{R}^d$ .*

**Assumption 1.2** ( $L$ -smoothness). *We assume that  $f_i$  is  $L_i$ -smooth for all  $i \in [n] = \{1, 2, \dots, n\}$  meaning that the*

*following inequality holds  $\forall x, y \in \mathbb{R}^d, \forall i \in [n]$ :*

$$\|\nabla f_i(x) - \nabla f_i(y)\| \leq L_i \|x - y\|. \quad (2)$$

This assumption implies that  $f$  is  $L_f$ -smooth with  $L_f^2 \leq L^2 = \frac{1}{n} \sum_{i=1}^n L_i^2$ .

Finally, we describe a large class of unbiased compression operators satisfying a certain variance bound, which we will refer to, in this paper, by the name *quantization*.

**Definition 1.1** (Quantization). *We say that a stochastic mapping  $\mathcal{Q} : \mathbb{R}^d \rightarrow \mathbb{R}^d$  is a quantization operator/quantization if there exists  $\omega > 0$  such that for any  $x \in \mathbb{R}^d$ , we have*

$$\mathbf{E}[\mathcal{Q}(x)] = x, \quad \mathbf{E}[\|\mathcal{Q}(x) - x\|^2] \leq \omega \|x\|^2. \quad (3)$$

*For the given quantization operator  $\mathcal{Q}(x)$ , we define the the expected density as  $\zeta_{\mathcal{Q}} = \sup_{x \in \mathbb{R}^d} \mathbf{E}[\|\mathcal{Q}(x)\|_0]$ , where  $\|y\|_0$  is the number of non-zero components of  $y \in \mathbb{R}^d$ .*

Notice that the expected density is well-defined for any quantization operator since  $\|\mathcal{Q}(x)\|_0 \leq d$ .



## 2. MARINA

In this section, we describe the main algorithm of this work: MARINA (see Algorithm 1). At each iteration of MARINA, each worker  $i$  either sends to the server the dense vector  $\nabla f_i(x^{k+1})$  with probability  $p$ , or it sends the quantized gradient difference  $\mathcal{Q}(\nabla f_i(x^{k+1}) - \nabla f_i(x^k))$  with probability  $1-p$ . In the first situation, the server just averages the vectors received from workers and gets  $g^{k+1} = \nabla f(x^{k+1})$ , whereas in the second case, the server averages the quantized differences from all workers and then adds the result to  $g^k$  to get  $g^{k+1}$ . Moreover, if  $\mathcal{Q}$  is identity quantization, i.e.,  $\mathcal{Q}(x) = x$ , then MARINA reduces to Gradient Descent (GD).

---

### Algorithm 1 MARINA

---

- 1: **Input:** starting point  $x^0$ , stepsize  $\gamma$ , probability  $p \in (0, 1]$ , number of iterations  $K$
  - 2: Initialize  $g^0 = \nabla f(x^0)$
  - 3: **for**  $k = 0, 1, \dots, K - 1$  **do**
  - 4:   Sample  $c_k \sim \text{Be}(p)$
  - 5:   Broadcast  $g^k$  to all workers
  - 6:   **for**  $i = 1, \dots, n$  **in parallel do**
  - 7:      $x^{k+1} = x^k - \gamma g^k$
  - 8:     Set  $g_i^{k+1} = \nabla f_i(x^{k+1})$  if  $c_k = 1$ , and  $g_i^{k+1} = g^k + \mathcal{Q}(\nabla f_i(x^{k+1}) - \nabla f_i(x^k))$  otherwise
  - 9:   **end for**
  - 10:  $g^{k+1} = \frac{1}{n} \sum_{i=1}^n g_i^{k+1}$
  - 11: **end for**
  - 12: **Return:**  $\hat{x}^K$  chosen uniformly at random from  $\{x^k\}_{k=0}^{K-1}$
- 

However, for non-trivial quantizations, we have  $\mathbf{E}[g^{k+1} | x^{k+1}] \neq \nabla f(x^{k+1})$  unlike all other distributed methods using exclusively unbiased compressors we know of. That is,  $g^{k+1}$  is a *biased* stochastic estimator of  $\nabla f(x^{k+1})$ . However, MARINA is an example of a rare phenomenon in stochastic optimization when the *bias of the stochastic gradient helps to achieve better complexity*.

### 2.1. Convergence Results for Generally Non-Convex Problems

We start with the following result.

**Theorem 2.1.** *Let Assumptions 1.1 and 1.2 be satisfied. Then, after*

$$K = \mathcal{O} \left( \frac{\Delta_0 L}{\varepsilon^2} \left( 1 + \sqrt{\frac{(1-p)\omega}{pn}} \right) \right)$$

*iterations with  $\Delta_0 = f(x^0) - f_*$ ,  $L^2 = \frac{1}{n} \sum_{i=1}^n L_i^2$  and the stepsize  $\gamma \leq L^{-1} \left( 1 + \sqrt{(1-p)\omega/(pn)} \right)^{-1}$ , MARINA produces point  $\hat{x}^K$  for which  $\mathbf{E}[\|\nabla f(\hat{x}^K)\|^2] \leq \varepsilon^2$ .*

One can find the full statement of the theorem together with its proof in Section C.1 of the Appendix.

The following corollary provides the bounds on the number of iterations/communication rounds and estimates the total communication cost needed to achieve an  $\varepsilon$ -stationary point in expectation. Moreover, for simplicity, throughout the paper we assume that the communication cost is proportional to the number of non-zero components of transmitted vectors from workers to the server.

**Corollary 2.1.** *Let the assumptions of Theorem 2.1 hold and  $p = \zeta_{\mathcal{Q}}/d$ . If  $\gamma \leq L^{-1} \left( 1 + \sqrt{\omega^{(d-\zeta_{\mathcal{Q}})/(n\zeta_{\mathcal{Q}})} \right)^{-1}$ , then MARINA requires*

$$\mathcal{O} \left( \frac{\Delta_0 L}{\varepsilon^2} \left( 1 + \sqrt{\frac{\omega}{n} \left( \frac{d}{\zeta_{\mathcal{Q}}} - 1 \right)} \right) \right)$$

*iterations/communication rounds in order to achieve  $\mathbf{E}[\|\nabla f(\hat{x}^K)\|^2] \leq \varepsilon^2$ , and the expected total communication cost per worker is  $\mathcal{O}(d + \zeta_{\mathcal{Q}}K)$ .*

Let us clarify the obtained result. First of all, if  $\omega = 0$  (no quantization), then  $\zeta_{\mathcal{Q}} = 0$  and the rate coincides with the rate of Gradient Descent (GD). Since GD is optimal among first-order methods in terms of reducing the norm of the gradient (Carmon et al., 2019), the dependence on  $\varepsilon$  in our bound cannot be improved in general. Next, if  $n$  is large enough, i.e.,  $n \geq \omega(d/\zeta_{\mathcal{Q}} - 1)$ , then<sup>1</sup> the iteration complexity of MARINA (method with compressed communications) and GD (method with dense communications) coincide. This means that in this regime, MARINA is able to reach a provably better communication complexity than GD!

### 2.2. Convergence Results Under Polyak-Łojasiewicz condition

In this section, we provide a complexity bounds for MARINA under the Polyak-Łojasiewicz (PŁ) condition.

**Assumption 2.1** (PŁ condition). *Function  $f$  satisfies Polyak-Łojasiewicz (PŁ) condition with parameter  $\mu$ , i.e.,*

$$\|\nabla f(x)\|^2 \geq 2\mu(f(x) - f(x^*)). \quad (4)$$

*holds for  $x^* = \arg \min_{x \in \mathbb{R}^d} f(x)$  and for all  $x \in \mathbb{R}^d$ .*

Under this and previously introduced assumptions, we derive the following result.

**Theorem 2.2.** *Let Assumptions 1.1, 1.2 and 2.1 be satisfied. Then, after*

$$K = \mathcal{O} \left( \max \left\{ \frac{1}{p}, \frac{L}{\mu} \left( 1 + \sqrt{\frac{(1-p)\omega}{pn}} \right) \right\} \log \frac{\Delta_0}{\varepsilon} \right)$$

<sup>1</sup>For  $\ell_2$ -quantization this requirement is satisfied when  $n \geq d$ .

iterations with  $\Delta_0 = f(x^0) - f(x^*)$ ,  $L^2 = \frac{1}{n} \sum_{i=1}^n L_i^2$  and the stepsize  $\gamma \leq \min \left\{ L^{-1} \left( 1 + \sqrt{2(1-p)\omega/(pn)} \right)^{-1}, p(2\mu)^{-1} \right\}$ , MARINA produces a point  $x^K$  for which  $\mathbf{E}[f(x^K) - f(x^*)] \leq \varepsilon$ .

One can find the full statement of the theorem together with its proof in Section C.2 of the Appendix.

### 3. Variance Reduction

Throughout this section, we assume that the local loss on each node has either a finite-sum form (finite sum case),

$$f_i(x) = \frac{1}{m} \sum_{j=1}^m f_{ij}(x), \quad (5)$$

or an expectation form (online case),

$$f_i(x) = \mathbf{E}_{\xi_i \sim \mathcal{D}_i} [f_{\xi_i}(x)]. \quad (6)$$

#### 3.1. Finite Sum Case

In this section, we generalize MARINA to problems of the form (1)+(5), obtaining VR-MARINA (see Algorithm 2). At

---

#### Algorithm 2 VR-MARINA: finite sum case

---

- 1: **Input:** starting point  $x^0$ , stepsize  $\gamma$ , minibatch size  $b'$ , probability  $p \in (0, 1]$ , number of iterations  $K$
  - 2: Initialize  $g^0 = \nabla f(x^0)$
  - 3: **for**  $k = 0, 1, \dots, K - 1$  **do**
  - 4:   Sample  $c_k \sim \text{Be}(p)$
  - 5:   Broadcast  $g^k$  to all workers
  - 6:   **for**  $i = 1, \dots, n$  **in parallel do**
  - 7:      $x^{k+1} = x^k - \gamma g^k$
  - 8:     Set  $g_i^{k+1} = \nabla f_i(x^{k+1})$  if  $c_k = 1$ , and  $g_i^{k+1} = g^k + \mathcal{Q} \left( \frac{1}{b'} \sum_{j \in I'_{i,k}} (\nabla f_{ij}(x^{k+1}) - \nabla f_{ij}(x^k)) \right)$  otherwise, where  $I'_{i,k}$  is the set of the indices in the minibatch,  $|I'_{i,k}| = b'$
  - 9:   **end for**
  - 10:    $g^{k+1} = \frac{1}{n} \sum_{i=1}^n g_i^{k+1}$
  - 11: **end for**
  - 12: **Return:**  $\hat{x}^K$  chosen uniformly at random from  $\{x^k\}_{k=0}^{K-1}$
- 

each iteration of VR-MARINA, devices are to compute the full gradients  $\nabla f_i(x^{k+1})$  and send them to the server with probability  $p$ . Typically,  $p \leq 1/m$  and  $m$  is large, meaning that workers compute full gradients rarely (once per  $\geq m$  iterations in expectation). At other iterations, workers compute minibatch stochastic gradients evaluated at the current and previous points, compress them using an unbiased compression operator, i.e., quantization/quantization operator, and send the resulting vectors  $g_i^{k+1} - g^k$  to the server. Moreover, if  $\mathcal{Q}$  is the identity quantization, i.e.,  $\mathcal{Q}(x) = x$ , and

$n = 1$ , then MARINA reduces to the optimal method PAGE (Li et al., 2020).

In this part, we will rely on the following average smoothness assumption.

**Assumption 3.1** (Average  $\mathcal{L}$ -smoothness). *For all  $k \geq 0$  and  $i \in [n]$  the minibatch stochastic gradients difference  $\tilde{\Delta}_i^k = \frac{1}{b'} \sum_{j \in I'_{i,k}} (\nabla f_{ij}(x^{k+1}) - \nabla f_{ij}(x^k))$  computed on the  $i$ -th worker satisfies  $\mathbf{E} [\tilde{\Delta}_i^k | x^k, x^{k+1}] = \Delta_i^k$  and*

$$\mathbf{E} \left[ \left\| \tilde{\Delta}_i^k - \Delta_i^k \right\|^2 \mid x^k, x^{k+1} \right] \leq \frac{\mathcal{L}_i^2}{b'} \|x^{k+1} - x^k\|^2 \quad (7)$$

with some  $\mathcal{L}_i \geq 0$ , where  $\Delta_i^k = \nabla f_i(x^{k+1}) - \nabla f_i(x^k)$ .

This assumption is satisfied in many standard minibatch regimes. In particular, if  $I'_{i,k} = \{1, \dots, m\}$ , then  $\mathcal{L}_i = 0$ , and if  $I'_{i,k}$  consists of  $b'$  i.i.d. samples from the uniform distributions on  $\{1, \dots, m\}$  and  $f_{ij}$  are  $L_{ij}$ -smooth, then  $\mathcal{L}_i \leq \max_{j \in [m]} L_{ij}$ .

Under this and the previously introduced assumptions, we derive the following result.

**Theorem 3.1.** *Consider the finite sum case (1)+(5). Let Assumptions 1.1, 1.2 and 3.1 be satisfied. Then, after*

$$K = \mathcal{O} \left( \frac{\Delta_0}{\varepsilon^2} \left( L + \sqrt{\frac{1-p}{pn}} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right) \right) \right)$$

iterations with  $\Delta_0 = f(x^0) - f_*$ ,  $L^2 = \frac{1}{n} \sum_{i=1}^n L_i^2$ ,  $\mathcal{L}^2 = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i^2$  and the stepsize  $\gamma \leq \left( L + \sqrt{(\omega L^2 + (1+\omega)\mathcal{L}^2/b') (1-p)/(pn)} \right)^{-1}$ , VR-MARINA produces such a point  $\hat{x}^K$  that  $\mathbf{E}[\|\nabla f(\hat{x}^K)\|^2] \leq \varepsilon^2$ .

One can find the full statement of the theorem together with its proof in Section D.1.1 of the Appendix.

**Corollary 3.1.** *Let the assumptions of Theorem 3.1 hold and  $p = \min \{\zeta_{\mathcal{Q}}/d, b'/(m+b')\}$ , where  $b' \leq m$ . If  $\gamma \leq \left( L + \sqrt{(\omega L^2 + (1+\omega)\mathcal{L}^2/b') \max\{d/\zeta_{\mathcal{Q}} - 1, m/b'\}/n} \right)^{-1}$  then VR-MARINA requires*

$$\mathcal{O} \left( \frac{\Delta_0}{\varepsilon^2} \left( L \left( 1 + \sqrt{\frac{\omega \max\{d/\zeta_{\mathcal{Q}} - 1, m/b'\}}{n}} \right) + \mathcal{L} \sqrt{\frac{(1+\omega) \max\{d/\zeta_{\mathcal{Q}} - 1, m/b'\}}{nb'}} \right) \right)$$

iterations/communication rounds and  $\mathcal{O}(m + b'K)$  stochastic oracle calls per node in expectation in order to achieve  $\mathbf{E}[\|\nabla f(\hat{x}^K)\|^2] \leq \varepsilon^2$ , and the expected total communication cost per worker is  $\mathcal{O}(d + \zeta_{\mathcal{Q}}K)$ .

First of all, when workers quantize differences of the full gradients, then  $I'_{i,k} = \{1, \dots, m\}$  for all  $i \in [n]$  and  $k \geq 0$ ,

implying  $\mathcal{L} = 0$ . In this case, the complexity bounds for VR-MARINA recover the ones for MARINA. Next, when  $\omega = 0$  (no quantization) and  $n = 1$ , our bounds for iteration and oracle complexities for VR-MARINA recover the bounds for PAGE (Li & Richtárik, 2020), which is optimal for finite-sum smooth non-convex optimization. This observation implies that the dependence on  $\varepsilon$  and  $m$  in the complexity bounds for VR-MARINA cannot be improved in the class of first-order stochastic methods. Next, we notice that up to the differences in smoothness constants, the iteration and oracle complexities for VR-MARINA benefit from the number of workers  $n$ . Finally, as Table 1 shows, the rates for VR-MARINA are strictly better than ones for the previous state-of-the-art method VR-DIANA (Horváth et al., 2019).

We provide the convergence results for VR-MARINA in the finite-sum case under the Polyak-Łojasiewicz condition, together with complete proofs, in Section D.1.2 of the Appendix.

### 3.2. Online Case

In this section, we focus on problems of type (1)+(6). For this type of problems, we consider a slightly modified version of VR-MARINA. That is, we replace line 8 in Algorithm 2 with the following update rule:  $g_i^{k+1} = \frac{1}{b} \sum_{j \in I_{i,k}} \nabla f_{\xi_{ij}^k}(x^{k+1})$  if  $c_k = 1$ , and  $g_i^{k+1} = g_i^k + \mathcal{Q} \left( \frac{1}{b'} \sum_{j \in I'_{i,k}} (\nabla f_{\xi_{ij}^k}(x^{k+1}) - \nabla f_{\xi_{ij}^k}(x^k)) \right)$  otherwise, where  $I_{i,k}, I'_{i,k}$  are the sets of the indices in the minibatches,  $|I_{i,k}| = b, |I'_{i,k}| = b'$ , and  $\xi_{ij}^k$  is independently sampled from  $\mathcal{D}_i$  for  $i \in [n], j \in [m]$  (see Algorithm 3 in the Appendix).

Before we provide our convergence results in this setup, we reformulate Assumption 3.1 for the online case.

**Assumption 3.2** (Average  $\mathcal{L}$ -smoothness). *For all  $k \geq 0$  and  $i \in [n]$  the minibatch stochastic gradients difference  $\tilde{\Delta}_i^k = \frac{1}{b'} \sum_{j \in I'_{i,k}} (\nabla f_{\xi_{ij}^k}(x^{k+1}) - \nabla f_{\xi_{ij}^k}(x^k))$  computed on the  $i$ -th worker satisfies  $\mathbf{E} [\tilde{\Delta}_i^k | x^k, x^{k+1}] = \Delta_i^k$  and*

$$\mathbf{E} \left[ \left\| \tilde{\Delta}_i^k - \Delta_i^k \right\|^2 \mid x^k, x^{k+1} \right] \leq \frac{\mathcal{L}^2}{b'} \|x^{k+1} - x^k\|^2 \quad (8)$$

with some  $\mathcal{L}_i \geq 0$ , where  $\Delta_i^k = \nabla f_i(x^{k+1}) - \nabla f_i(x^k)$ .

Moreover, we assume that the variance of the stochastic gradients on all nodes is uniformly upper bounded.

**Assumption 3.3.** *We assume that for all  $i \in [n]$  there exists such constant  $\sigma_i \in [0, +\infty)$  that for all  $x \in \mathbb{R}^d$*

$$\begin{aligned} \mathbf{E}_{\xi_i \sim \mathcal{D}_i} [\nabla f_{\xi_i}(x)] &= \nabla f_i(x), \quad (9) \\ \mathbf{E}_{\xi_i \sim \mathcal{D}_i} [\|\nabla f_{\xi_i}(x) - \nabla f_i(x)\|^2] &\leq \sigma_i^2. \quad (10) \end{aligned}$$

Under these and previously introduced assumptions, we derive the following result.

**Theorem 3.2.** *Consider the online case (1)+(6). Let Assumptions 1.1, 1.2, 3.2 and 3.3 be satisfied. Then, after*

$$K = \mathcal{O} \left( \frac{\Delta_0}{\varepsilon^2} \left( L + \sqrt{\frac{1-p}{pn} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right)} \right) \right)$$

iterations with  $\Delta_0 = f(x^0) - f_*$ ,  $L^2 = \frac{1}{n} \sum_{i=1}^n L_i^2$ ,  $\mathcal{L}^2 = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i^2$ , the stepsize  $\gamma \leq \left( L + \sqrt{(\omega L^2 + (1+\omega)\mathcal{L}^2/b')^{(1-p)/(pn)}} \right)^{-1}$ , and  $b = \Theta(\sigma^2/(n\varepsilon^2))$ ,  $\sigma^2 = \frac{1}{n} \sum_{i=1}^n \sigma_i^2$ , VR-MARINA produces a point  $\hat{x}^K$  for which  $\mathbf{E}[\|\nabla f(\hat{x}^K)\|^2] \leq \varepsilon^2$ .

One can find the full statement of the theorem, together with its proof, in Section D.2.1 of the Appendix.

**Corollary 3.2.** *Let the assumptions of Theorem 3.2 hold and choose  $p = \min\{\zeta_{\mathcal{Q}}/d, b'/(b+b')\}$ , where  $b' \leq b$ ,  $b = \Theta(\sigma^2/(n\varepsilon^2))$ . If  $\gamma \leq \left( L + \sqrt{(\omega L^2 + (1+\omega)\mathcal{L}^2/b')^{\max\{d/\zeta_{\mathcal{Q}}-1, b/b'\}}/n} \right)^{-1}$ , then VR-MARINA requires*

$$\mathcal{O} \left( \frac{\Delta_0}{\varepsilon^2} \left( L \left( 1 + \sqrt{\frac{\omega}{n} \max\left\{ \frac{d}{\zeta_{\mathcal{Q}}} - 1, \frac{\sigma^2}{nb'\varepsilon^2} \right\}} \right) + \mathcal{L} \sqrt{\frac{(1+\omega)}{nb'}} \max\left\{ \frac{d}{\zeta_{\mathcal{Q}}} - 1, \frac{\sigma^2}{nb'\varepsilon^2} \right\} \right) \right)$$

iterations/communication rounds and  $\mathcal{O}(\zeta_{\mathcal{Q}}K + \sigma^2/(n\varepsilon^2))$  stochastic oracle calls per node in expectation to achieve  $\mathbf{E}[\|\nabla f(\hat{x}^K)\|^2] \leq \varepsilon^2$ , and the expected total communication cost per worker is  $\mathcal{O}(d + \zeta_{\mathcal{Q}}K)$ .

Similarly to the finite-sum case, when  $\omega = 0$  (no quantization) and  $n = 1$ , our bounds for iteration and oracle complexities for VR-MARINA recover the bounds for PAGE (Li & Richtárik, 2020), which is optimal for online smooth non-convex optimization as well. That is, the dependence on  $\varepsilon$  in the complexity bound for VR-MARINA cannot be improved in the class of first-order stochastic methods. As previously, up to the differences in smoothness constants, the iteration and oracle complexities for VR-MARINA benefit from an increase in the number of workers  $n$ .

We provide the convergence results for VR-MARINA in the online case under the Polyak-Łojasiewicz condition, together with complete proofs, in Section D.2.2 of the Appendix.

## 4. Partial Participation

Finally, we propose another modification of MARINA. In particular, we prove an option for *partial participation* of



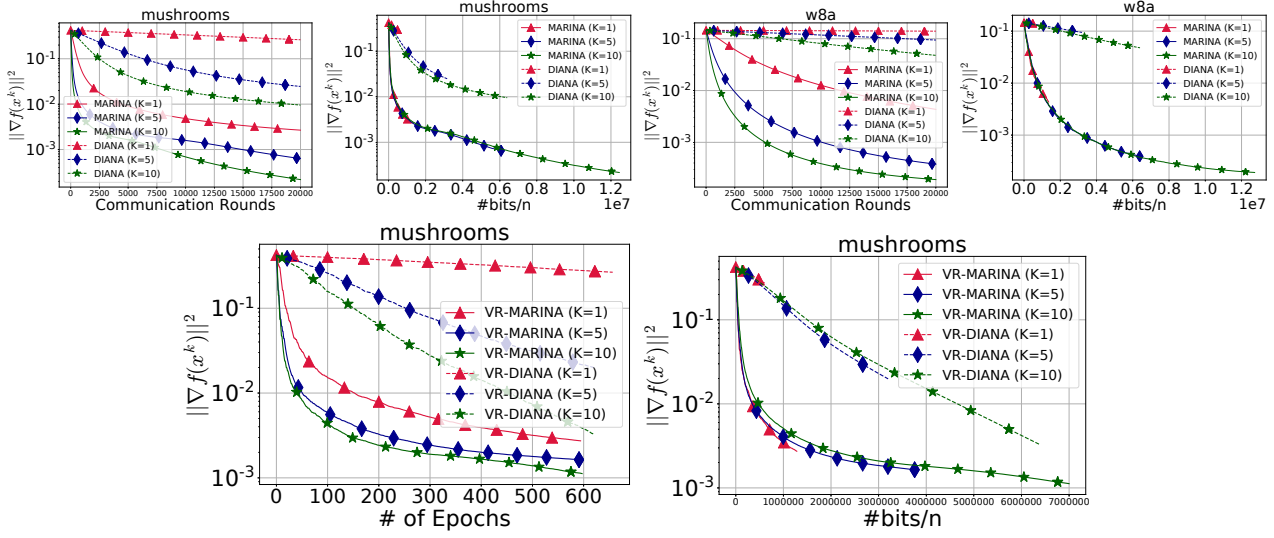


Figure 1: Comparison of MARINA with DIANA, and of VR-MARINA with VR-DIANA, on binary classification problem involving non-convex loss (11) with LibSVM data (Chang & Lin, 2011). Parameter  $n$  is chosen as per Tbl. 3 in the Appendix. Stepsizes for the methods are chosen according to the theory and the batchsizes for VR-MARINA and VR-DIANA are  $\sim m/100$ . In all cases, we used the RandK sparsification operator with  $K \in \{1, 5, 10\}$ .

the clients - a feature important in federated learning. The resulting method is called PP-MARINA (see Algorithm 4 in the Appendix). At each iteration of PP-MARINA, the server receives the quantized gradient differences from  $r$  clients with probability  $1 - p$ , and aggregates full gradients from all clients with probability  $p$ , i.e., PP-MARINA coincides with MARINA up to the following difference:  $g_i^{k+1} = \nabla f_i(x^{k+1})$ ,  $g^{k+1} = \frac{1}{n} \sum_{i=1}^n g_i^{k+1}$  if  $c_k = 1$ , and  $g_i^{k+1} = g^k + \mathcal{Q}(\nabla f_i(x^{k+1}) - \nabla f_i(x^k))$ ,  $g^{k+1} = \frac{1}{r} \sum_{i_k \in I'_k} g_i^{k+1}$  otherwise, where  $I'_k$  is the set of  $r$  i.i.d. samples from the uniform distribution over  $\{1, \dots, n\}$ . That is, if the probability  $p$  is chosen to be small enough, then with high probability the server receives only quantized vectors from a subset of clients at each iteration.

Below, we provide a convergence result for PP-MARINA for smooth non-convex problems.

**Theorem 4.1.** *Let Assumptions 1.1 and 1.2 be satisfied. Then, after*

$$K = \mathcal{O} \left( \frac{\Delta_0 L}{\varepsilon^2} \left( 1 + \sqrt{\frac{(1-p)(1+\omega)}{pr}} \right) \right)$$

iterations with  $\Delta_0 = f(x^0) - f_*$ ,  $L^2 = \frac{1}{n} \sum_{i=1}^n L_i^2$  and the stepsize  $\gamma \leq L^{-1} \left( 1 + \sqrt{\frac{(1-p)(1+\omega)}{pr}} \right)^{-1}$ , PP-MARINA produces a point  $\hat{x}^K$  for which  $\mathbf{E}[\|\nabla f(\hat{x}^K)\|^2] \leq \varepsilon^2$ .

One can find the full statement of the theorem together with its proof in Section E.1 of the appendix.

**Corollary 4.1.** *Let the assumptions of Theorem 4.1 hold and choose  $p = \zeta_{\mathcal{Q}^r}/(dn)$ , where  $r \leq n$ . If  $\gamma \leq$*

$L^{-1} \left( 1 + \sqrt{\frac{(1+\omega)(dn - \zeta_{\mathcal{Q}^r})}{(b/\zeta_{\mathcal{Q}^r})}} \right)^{-1}$ , then PP-MARINA requires

$$\mathcal{O} \left( \frac{\Delta_0 L}{\varepsilon^2} \left( 1 + \sqrt{\frac{1+\omega}{r} \left( \frac{dn}{\zeta_{\mathcal{Q}^r}} - 1 \right)} \right) \right)$$

iterations/communication rounds to achieve  $\mathbf{E}[\|\nabla f(\hat{x}^K)\|^2] \leq \varepsilon^2$ , and the expected total communication cost is  $\mathcal{O}(dn + \zeta_{\mathcal{Q}^r} K)$ .

When  $r = n$ , i.e., all clients participate in communication with the server at each iteration, the rate for PP-MARINA recovers the rate for MARINA under the assumption that  $(1 + \omega)(d/\zeta_{\mathcal{Q}} - 1) = \mathcal{O}(\omega(d/\zeta_{\mathcal{Q}} - 1))$ , which holds for a wide class of quantization operators, e.g., for identical quantization, RandK, and  $\ell_p$ -quantization. In general, the derived complexity is strictly better than previous state-of-the-art one (see Table 1).

We provide the convergence results for PP-MARINA under the Polyak-Łojasiewicz condition, together with complete proofs, in Section E.2 of the Appendix.

## 5. Numerical Experiments

### 5.1. Binary Classification with Non-Convex Loss

We conduct several numerical experiments<sup>2</sup> on binary classification problem involving non-convex loss (Zhao et al., 2010) (used for two-layer neural networks) with LibSVM

<sup>2</sup>Our code is available at <https://github.com/burlachenkok/marina>.

data (Chang & Lin, 2011) to justify the theoretical claims of the paper. That is, we consider the following optimization problem:

$$\min_{x \in \mathbb{R}^d} \left\{ f(x) = \frac{1}{N} \sum_{t=1}^N \ell(a_t^\top x, y_t) \right\}, \quad (11)$$

where  $\{a_t\} \in \mathbb{R}^d$ ,  $y_t \in \{-1, 1\}$  for all  $t = 1, \dots, N$ , and the function  $\ell : \mathbb{R}^d \rightarrow \mathbb{R}$  is defined as

$$\ell(b, c) = \left( 1 - \frac{1}{1 + \exp(-bc)} \right)^2.$$

The distributed environment is simulated in Python 3.8 using MPI4PY and other standard libraries. Additional details about the experimental setup together with extra experiments are deferred to Section A of the Appendix.

In our experiments, we compare MARINA with the full-batch version of DIANA, and then VR-MARINA with VR-DIANA. We exclude FedCOMGATE and FedPATH from this comparison since they have significantly worse oracle complexities (see Table 1). The results are presented in Fig. 1. As our theory predicts, the first row shows the superiority of MARINA to DIANA both in terms of iteration/communication complexity and the total number of transmitted bits to achieve the given accuracy. Next, to study the oracle complexity as well, we consider non-full-batched methods – VR-MARINA and VR-DIANA – since they have better oracle complexity than the full-batched methods in the finite-sum case. Again, the results presented in the second row justify that VR-MARINA outperforms VR-DIANA in terms of oracle complexity and the total number of transmitted bits to achieve the given accuracy.

## 5.2. Image Classification

We also compared the performance of VR-MARINA and VR-DIANA on the training ResNet-18 (He et al., 2016) at CIFAR100 (Krizhevsky et al., 2009) dataset. Formally, the optimization problem is

$$\min_{x \in \mathbb{R}^d} \left\{ f(x) = \frac{1}{N} \sum_{i=1}^N \ell(p(f(a_i, x)), y_i) \right\}, \quad (12)$$

where  $\{(a_i, y_i)\}_{i=1}^N$  encode images and labels from CIFAR100 dataset,  $f(a_i, x)$  is the output of ResNet-18 on image  $a_i$  with weights  $x$ ,  $p$  is softmax function, and  $\ell(\cdot, \cdot)$  is cross-entropy loss. The code is written in Python 3.9 using PyTorch 1.7, and the distributed environment is simulated.

The results are presented in Fig. 2. Again, VR-MARINA converges significantly faster than VR-DIANA both in terms of the oracle complexity and the total number of transmitted bits to achieve the given accuracy. See other details and observations in Section A of the Appendix.

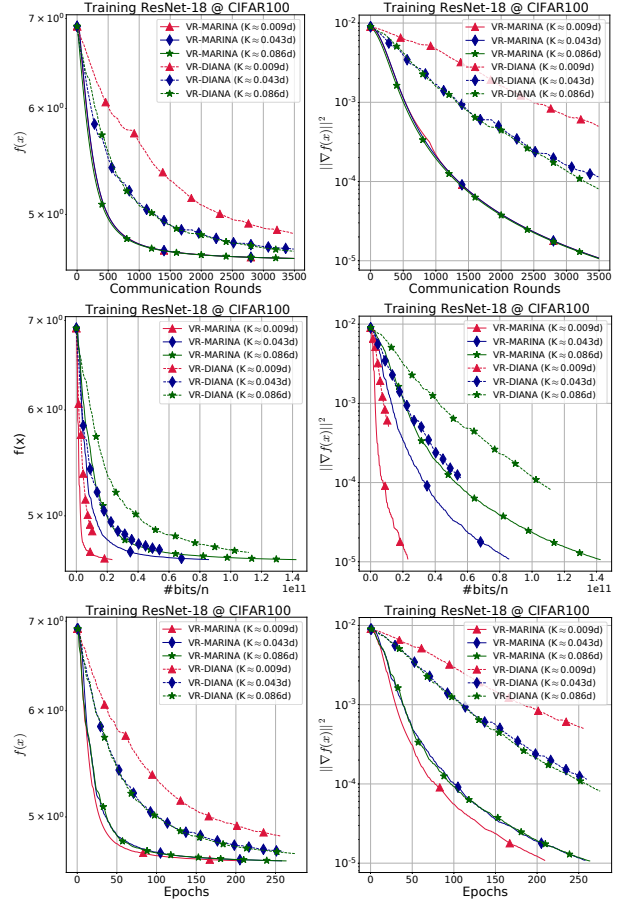


Figure 2: Comparison of VR-MARINA with VR-DIANA on training ResNet-18 at CIFAR100 dataset. Number of workers equals 5. Stepsizes for the methods were tuned and the batchsizes are  $\sim m/50$ . In all cases, we used the RandK sparsification operator, the approximate values of  $K$  are given in the legends ( $d$  is dimension of the problem).

## Acknowledgements

The work of Peter Richtárik, Eduard Gorbunov, Konstantin Burlachenko and Zhize Li was supported by KAUST Baseline Research Fund. The paper was written while E. Gorbunov was a research intern at KAUST. The work of E. Gorbunov in Sections 1, 2, and C was also partially supported by the Ministry of Science and Higher Education of the Russian Federation (Goszadaniye) 075-00337-20-03, project No. 0714-2020-0005, and in Sections 3, 4, D, E – by RFBR, project number 19-31-51001. We thank Konstantin Mishchenko (KAUST) for a suggestion related to the experiments, Elena Bazanova (MIPT) for the suggestions about improving the text, and Slavomír Hanzely (KAUST) and Egor Shulgin (KAUST) for spotting the typos.

## References

- Alistarh, D., Grubic, D., Li, J., Tomioka, R., and Vojnovic, M. QSGD: Communication-efficient SGD via gradient quantization and encoding. In *Advances in Neural Information Processing Systems*, pp. 1709–1720, 2017.
- Arjevani, Y., Carmon, Y., Duchi, J. C., Foster, D. J., Srebro, N., and Woodworth, B. Lower bounds for non-convex stochastic optimization. *arXiv preprint arXiv:1912.02365*, 2019.
- Basu, D., Data, D., Karakus, C., and Diggavi, S. Qsparse-local-SGD: Distributed SGD with quantization, sparsification and local computations. In *Advances in Neural Information Processing Systems*, pp. 14668–14679, 2019.
- Beznosikov, A., Horváth, S., Richtárik, P., and Safaryan, M. On biased compression for distributed learning. *arXiv preprint arXiv:2002.12410*, 2020.
- Bhojanapalli, S., Kyrillidis, A., and Sanghavi, S. Dropping convexity for faster semi-definite optimization. In *Conference on Learning Theory*, pp. 530–582. PMLR, 2016.
- Carmon, Y., Duchi, J. C., Hinder, O., and Sidford, A. Lower bounds for finding stationary points i. *Mathematical Programming*, pp. 1–50, 2019.
- Chang, C.-C. and Lin, C.-J. LIBSVM: A library for support vector machines. *ACM transactions on intelligent systems and technology (TIST)*, 2(3):1–27, 2011.
- Danilova, M., Dvurechensky, P., Gasnikov, A., Gorbunov, E., Guminov, S., Kamzolov, D., and Shibaev, I. Recent theoretical advances in non-convex optimization. *arXiv preprint arXiv:2012.06188*, 2020.
- Das, R., Hashemi, A., Sanghavi, S., and Dhillon, I. S. Improved convergence rates for non-convex federated learning with compression. *arXiv preprint arXiv:2012.04061*, 2020.
- Fang, C., Li, C., Lin, Z., and Zhang, T. Near-optimal non-convex optimization via stochastic path integrated differential estimator. *Advances in Neural Information Processing Systems*, 31:689, 2018.
- Goodfellow, I., Bengio, Y., Courville, A., and Bengio, Y. *Deep learning*, volume 1. MIT press Cambridge, 2016.
- Gorbunov, E., Kovalev, D., Makarenko, D., and Richtárik, P. Linearly converging error compensated sgd. *Advances in Neural Information Processing Systems*, 33, 2020.
- Goyal, P., Dollár, P., Girshick, R., Noordhuis, P., Wesolowski, L., Kyrola, A., Tulloch, A., Jia, Y., and He, K. Accurate, large minibatch SGD: training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017.
- Haddadpour, F., Kamani, M. M., Mokhtari, A., and Mavadavi, M. Federated learning with compression: Unified analysis and sharp guarantees. *arXiv preprint arXiv:2007.01154*, 2020.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Horváth, S., Ho, C.-Y., Ľudovít Horváth, Sahu, A. N., Canini, M., and Richtárik, P. Natural compression for distributed deep learning. *arXiv preprint arXiv:1905.10988*, 2019.
- Horváth, S., Kovalev, D., Mishchenko, K., Stich, S., and Richtárik, P. Stochastic distributed learning with gradient quantization and variance reduction. *arXiv preprint arXiv:1904.05115*, 2019.
- Kairouz, P., McMahan, H. B., Avent, B., Bellet, A., Bennis, M., Bhagoji, A. N., Bonawitz, K., Charles, Z., Cormode, G., Cummings, R., et al. Advances and open problems in federated learning. *arXiv preprint arXiv:1912.04977*, 2019.
- Karimireddy, S. P., Rebjock, Q., Stich, S., and Jaggi, M. Error feedback fixes signSGD and other gradient compression schemes. In *International Conference on Machine Learning*, pp. 3252–3261, 2019.
- Karimireddy, S. P., Kale, S., Mohri, M., Reddi, S., Stich, S., and Suresh, A. T. Scaffold: Stochastic controlled averaging for federated learning. In *International Conference on Machine Learning*, pp. 5132–5143. PMLR, 2020.
- Khanduri, P., Sharma, P., Kafle, S., Bulusu, S., Rajawat, K., and Varshney, P. K. Distributed stochastic non-convex optimization: Momentum-based variance reduction. *arXiv preprint arXiv:2005.00224*, 2020.
- Koloskova, A., Lin, T., Stich, S. U., and Jaggi, M. Decentralized deep learning with arbitrary communication compression. *ICLR*, pp. arXiv:1907.09356, 2020a. URL <https://arxiv.org/abs/1907.09356>.
- Koloskova, A., Loizou, N., Boreiri, S., Jaggi, M., and Stich, S. A unified theory of decentralized sgd with changing topology and local updates. In *International Conference on Machine Learning*, pp. 5381–5393. PMLR, 2020b.
- Konečný, J., McMahan, H. B., Yu, F. X., Richtárik, P., Suresh, A. T., and Bacon, D. Federated learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492*, 2016.
- Konečný, J., McMahan, H. B., Yu, F., Richtárik, P., Suresh, A. T., and Bacon, D. Federated learning: strategies for

- improving communication efficiency. In *NIPS Private Multi-Party Machine Learning Workshop*, 2016.
- Krizhevsky, A., Hinton, G., et al. Learning multiple layers of features from tiny images. 2009.
- Li, X., Yang, W., Wang, S., and Zhang, Z. Communication efficient decentralized training with multiple local updates. *arXiv preprint arXiv:1910.09126*, 5, 2019.
- Li, Z. and Richtárik, P. A unified analysis of stochastic gradient methods for nonconvex federated optimization. *arXiv preprint arXiv:2006.07013*, 2020.
- Li, Z., Bao, H., Zhang, X., and Richtárik, P. Page: A simple and optimal probabilistic gradient estimator for non-convex optimization. *arXiv preprint arXiv:2008.10898*, 2020.
- Lian, X., Zhang, C., Zhang, H., Hsieh, C.-J., Zhang, W., and Liu, J. Can decentralized algorithms outperform centralized algorithms? a case study for decentralized parallel stochastic gradient descent. In *Advances in Neural Information Processing Systems*, pp. 5330–5340, 2017.
- Łojasiewicz, S. A topological property of real analytic subsets. *Coll. du CNRS, Les équations aux dérivées partielles*, 117:87–89, 1963.
- Ma, C., Wang, K., Chi, Y., and Chen, Y. Implicit regularization in nonconvex statistical estimation: Gradient descent converges linearly for phase retrieval and matrix completion. In *International Conference on Machine Learning*, pp. 3345–3354. PMLR, 2018.
- McMahan, H. B., Moore, E., Ramage, D., Hampson, S., and Agüera y Arcas, B. Communication-efficient learning of deep networks from decentralized data. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2017.
- Mishchenko, K., Gorbunov, E., Takáč, M., and Richtárik, P. Distributed learning with compressed gradient differences. *arXiv preprint arXiv:1901.09269*, 2019.
- Murty, K. and Kabadi, S. Some np-complete problems in quadratic and nonlinear programming. *Mathematical programming*, 39(2):117–129, 1987.
- Polyak, B. T. Gradient methods for the minimisation of functionals. *USSR Computational Mathematics and Mathematical Physics*, 3(4):864–878, 1963.
- Qian, X., Richtárik, P., and Zhang, T. Error compensated distributed sgd can be accelerated. *arXiv preprint arXiv:2010.00091*, 2020.
- Safaryan, M., Shulgin, E., and Richtárik, P. Uncertainty principle for communication compression in distributed and federated learning and the search for an optimal compressor. *arXiv preprint arXiv:2002.08958*, 2020.
- Seide, F., Fu, H., Droppo, J., Li, G., and Yu, D. 1-bit stochastic gradient descent and its application to data-parallel distributed training of speech dnns. In *Fifteenth Annual Conference of the International Speech Communication Association*, 2014.
- Sharma, P., Kafle, S., Khanduri, P., Bulusu, S., Rajawat, K., and Varshney, P. K. Parallel restarted spider-communication efficient distributed nonconvex optimization with optimal computation complexity. *arXiv preprint arXiv:1912.06036*, 2019.
- Stich, S. U. and Karimireddy, S. P. The error-feedback framework: Better rates for sgd with delayed gradients and compressed updates. *Journal of Machine Learning Research*, 21:1–36, 2020.
- Sun, H., Lu, S., and Hong, M. Improving the sample and communication complexity for decentralized non-convex optimization: Joint gradient estimation and tracking. In *International Conference on Machine Learning*, pp. 9217–9228. PMLR, 2020.
- Sun, R. Optimization for deep learning: theory and algorithms. *arXiv preprint arXiv:1912.08957*, 2019.
- Suresh, A. T., Yu, F. X., Kumar, S., and McMahan, H. B. Distributed mean estimation with limited communication. In *Proceedings of the 34th International Conference on Machine Learning*, 2017.
- You, Y., Li, J., Reddi, S., Hseu, J., Kumar, S., Bhojanapalli, S., Song, X., Demmel, J., Keutzer, K., and Hsieh, C.-J. Large batch optimization for deep learning: Training bert in 76 minutes. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=Syx4wnEtvH>.
- Zhao, L., Mammadov, M., and Yearwood, J. From convex to nonconvex: a loss function analysis for binary classification. In *2010 IEEE International Conference on Data Mining Workshops*, pp. 1281–1288. IEEE, 2010.



# Appendix

## A. Extra Experiments

### A.1. Binary Classification with Non-Convex Loss

#### A.1.1. SETUP

In Section 5.1, we present the behavior of MARINA, VR-MARINA, DIANA, and VR-DIANA on the binary classification problem involving non-convex loss (Zhao et al., 2010). The datasets were taken from LibSVM (Chang & Lin, 2011) and split into five equal parts among five clients (we excluded  $N - 5 \cdot \lfloor N/5 \rfloor$  last datapoints from each dataset), see the summary in Table 3.

Table 3: Summary of the datasets and splitting of the data among clients (Figure 1).

Dataset	$n$	$N$ (# of datapoints)	$d$ (# of features)
mushrooms	5	8 120	112
w8a	5	49 745	300
phishing	5	11 055	69
a9a	5	32 560	124

The code was written in Python 3.8 using MPI4PY to emulate the distributed environment and then was executed on a machine with 48 cores, each is Intel(R) Xeon(R) Gold 6246 CPU 3.30GHz.

#### A.1.2. EXTRA EXPERIMENTS

In this section, we provide additional numerical results on the comparison of MARINA, VR-MARINA, DIANA, and VR-DIANA on the problem (11). Since one of the main goals of our experiments is to justify the theoretical findings of the paper, in the experiments, we used the stepsizes from the corresponding theoretical results for the methods (for DIANA and VR-DIANA the stepsizes were chosen according to (Horváth et al., 2019; Li & Richtárik, 2020)). Next, to compute the stochastic gradients, we use batchsizes =  $\max\{1, m/100\}$  for VR-MARINA and VR-DIANA.

The results for the full-batched methods are reported in Figure 3, and the comparison of VR-MARINA and VR-DIANA is given in Figure 4. Clearly, in both cases, MARINA and VR-MARINA show faster convergence than the previous state-of-the-art methods, DIANA and VR-DIANA, for distributed non-convex optimization with compression in terms of  $\|\nabla f(x^k)\|^2$  and  $f(x^k)$  decrease w.r.t. the number of communication rounds, oracle calls per node and the total number of transferred bits from workers to the master.

We also tested MARINA and DIANA on mushrooms dataset with a bigger number of workers ( $n = 20$ ). The results are reported in Figure 5. Similarly to the previous numerical tests, MARINA shows its superiority to DIANA with  $n = 20$  as well.

## A.2. Image Classification

#### A.2.1. SETUP

In Section 5.2, we demonstrate the performance of VR-MARINA and VR-DIANA on training ResNet-18 at CIFAR100 dataset. ResNet-18 has  $d = 11\,689\,512$  parameters to train and CIFAR100 contains  $N = 50\,000$  colored images. The dataset is split into 5 parts among 5 workers in such a way that the first four workers get 10 112 samples and the fifth one get 9 552 samples. The code was written in Python 3.9 using PYTORCH 1.7 and then was executed on a machine with NVIDIA GPU Geforce RTX 2080 Ti with 11 GByte onboard global GPU memory.

In all experiments, we use batchsize = 256 on each worker and tune the stepsizes for each method separately. That is, for each method and for each choice of  $K$  for RandK operator we run the method with stepsize  $\gamma \in \{10^{-6}, 0.1, 0.2, 0.5, 1.0, 5.0\}$  to find the interval containing the best stepsize. Next, the obtained interval is split into  $\sim 10$  equal parts and the method is run with corresponding stepsizes. Other parameters of the methods are chosen according to the theory. The summary of used



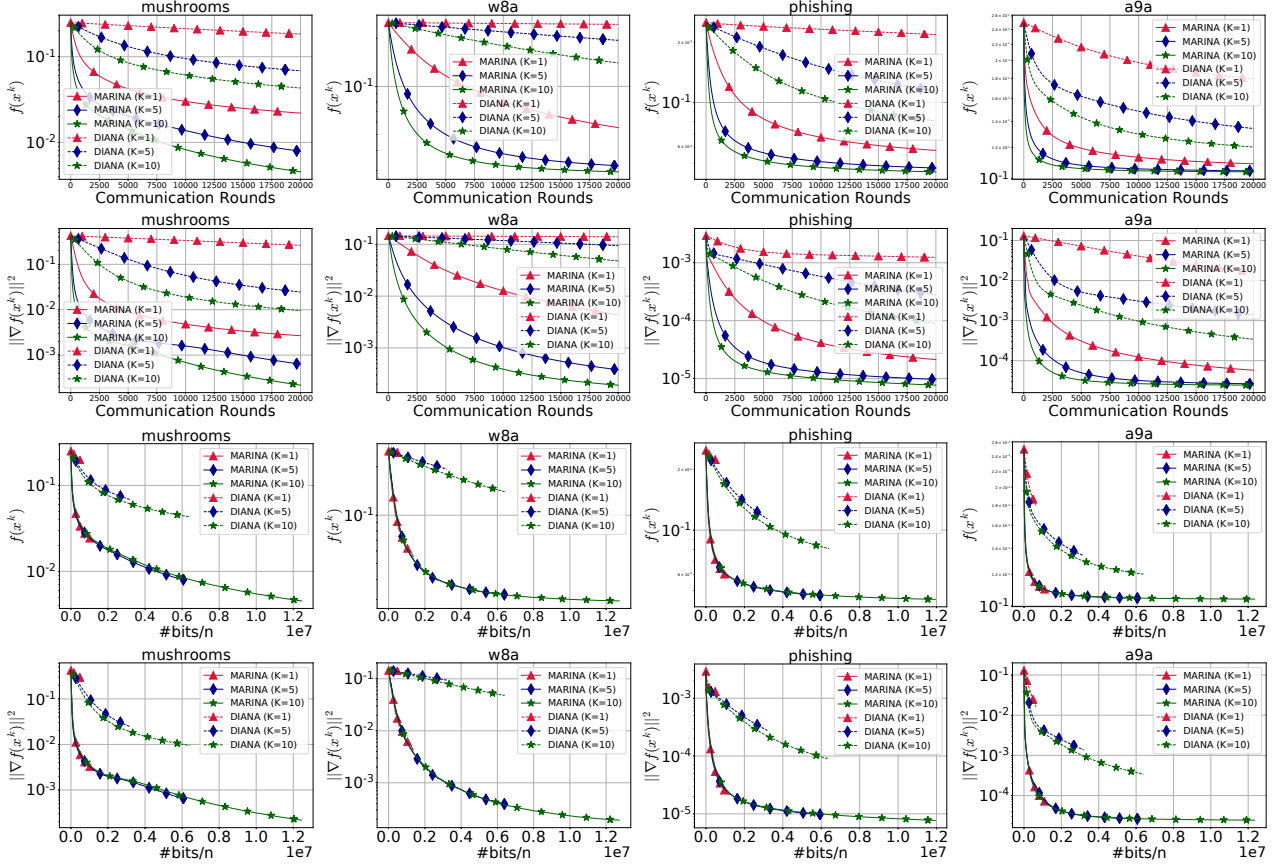


Figure 3: Comparison of MARINA with DIANA on binary classification problem involving non-convex loss (11) with LibSVM data (Chang & Lin, 2011). Parameter  $n$  is chosen as per Tbl. 3 ( $n = 5$ ). Stepsizes for the methods are chosen according to the theory. In all cases, we used the RandK sparsification operator with  $K \in \{1, 5, 10\}$ .

parameters is given in Table 4.

Table 4: Summary of the parameters used in the experiments presented in Fig. 2 and Fig. 6. Stepsizes were tuned, batchsize = 256 on each worker, other parameters were picked according to the theory, except the last line, where  $p$  for VR-MARINA without compression was picked as for VR-MARINA with RandK,  $K = 100\,000$  compression operator.

Method	RandK, $K =$	$\gamma$	$p$
VR-MARINA	100 000	0.95	0.008554
VR-MARINA	500 000	0.95	0.024691
VR-MARINA	1 000 000	0.95	0.024691
VR-DIANA	100 000	0.15	0.025316
VR-DIANA	500 000	0.35	0.025316
VR-DIANA	1 000 000	0.35	0.025316
VR-MARINA	11 689 512 ( $K = d$ )	3.5	0.024691
VR-DIANA	11 689 512 ( $K = d$ )	2.5	0.025316
VR-MARINA	11 689 512 ( $K = d$ )	3.5	0.008554

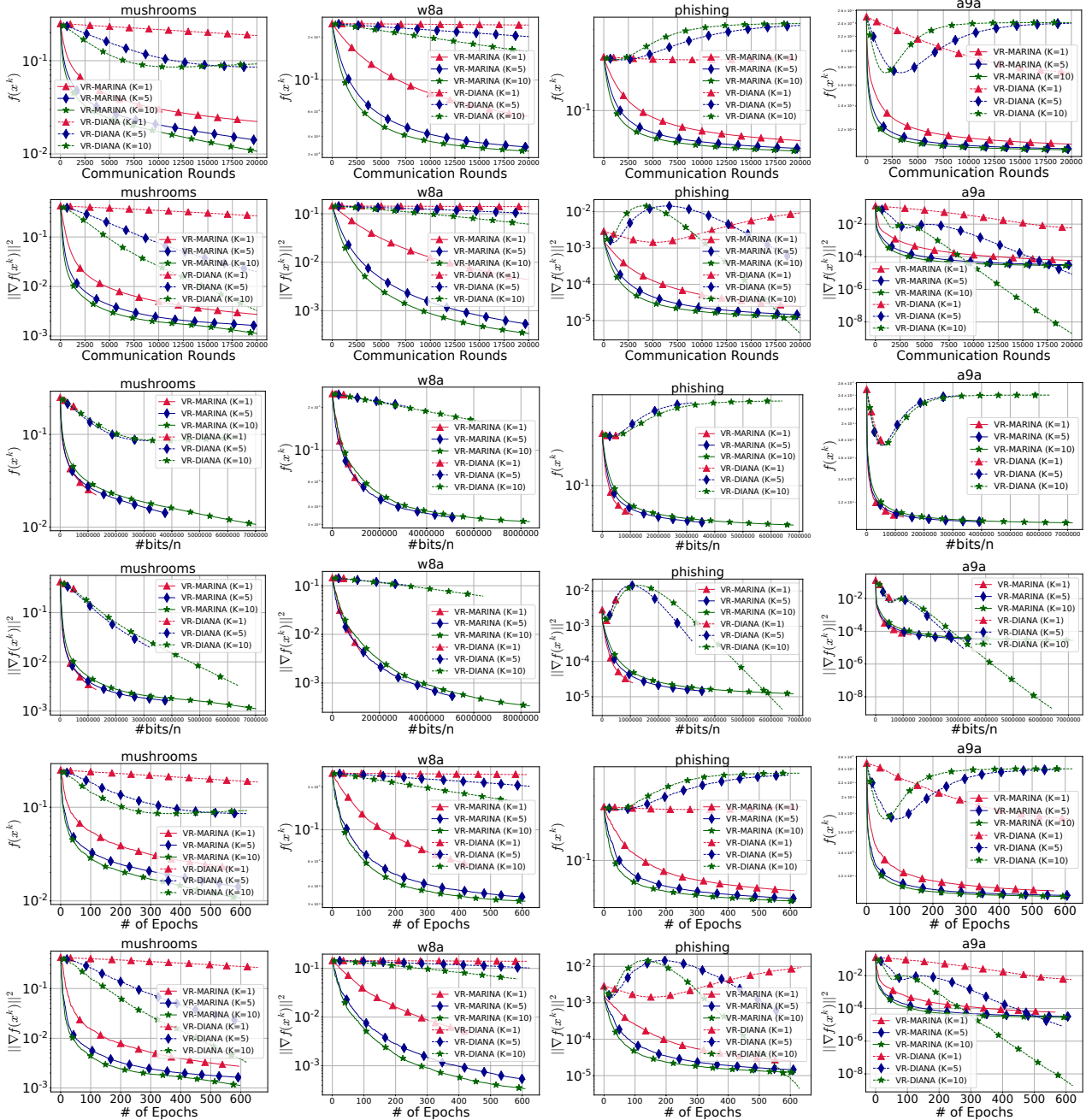


Figure 4: Comparison of VR-MARINA with VR-DIANA on binary classification problem involving non-convex loss (11) with LibSVM data (Chang & Lin, 2011). Parameter  $n$  is chosen as per Tbl. 3 ( $n = 5$ ). Stepsizes for the methods are chosen according to the theory and the batchsizes are  $\sim m/100$ . In all cases, we used the RandK sparsification operator with  $K \in \{1, 5, 10\}$ .

### A.2.2. EXTRA EXPERIMENTS

Results presented in Fig. 2 show the superiority of VR-MARINA to VR-DIANA in training ResNet-18 at CIFAR100. To emphasize the effect of compression we also run VR-MARINA and VR-DIANA without compression, see the results in Fig. 6. First of all, one can notice that the methods do benefit from compression: VR-MARINA and VR-DIANA with compression converge much faster than their non-compressed versions in terms of the total number of transmitted bits to achieve given

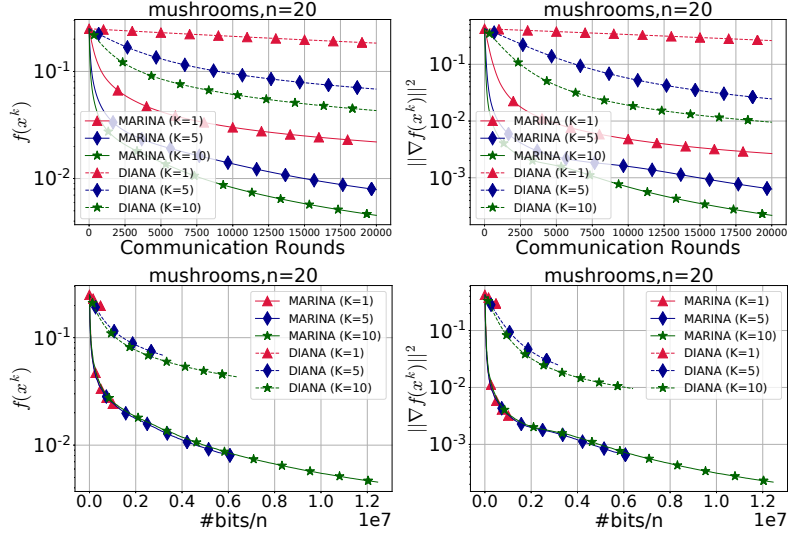


Figure 5: Comparison of MARINA with DIANA on binary classification problem involving non-convex loss (11) with mushrooms dataset and  $n = 20$  workers. Stepsizes for the methods are chosen according to the theory. In all cases, we used the RandK sparsification operator with  $K \in \{1, 5, 10\}$ .

accuracy.

Moreover, as Fig. 2 shows, VR-MARINA with  $K = 100\,000$  converges faster than VR-MARINA with larger  $K$  in terms of the epochs. That is, the method with more aggressive compression requires less oracle calls to achieve the same accuracy. The reason of such an unusual behavior is the choice of  $p$ : when  $K = 100\,000$  the theoretical choice of  $p$  is much smaller than for  $K = 500\,000$  and  $K = 1\,000\,000$ . Therefore, in VR-MARINA with  $K = 100\,000$ , the workers compute the full gradients more rarely than in the case of larger  $K$ . As the result, it turns out, that the total number of oracle calls needed to achieve given accuracy also smaller for  $K = 100\,000$  than for larger  $K$ . Moreover, we see this phenomenon even without applying compression: VR-MARINA without compression and with  $p$  as in the experiment with VR-MARINA with  $K = 100\,000$  converges faster than VR-MARINA without compression and with theoretical choice of  $p$ , which is the same as in the case when  $K = 500\,000, 1\,000\,000$ , see Table 4.

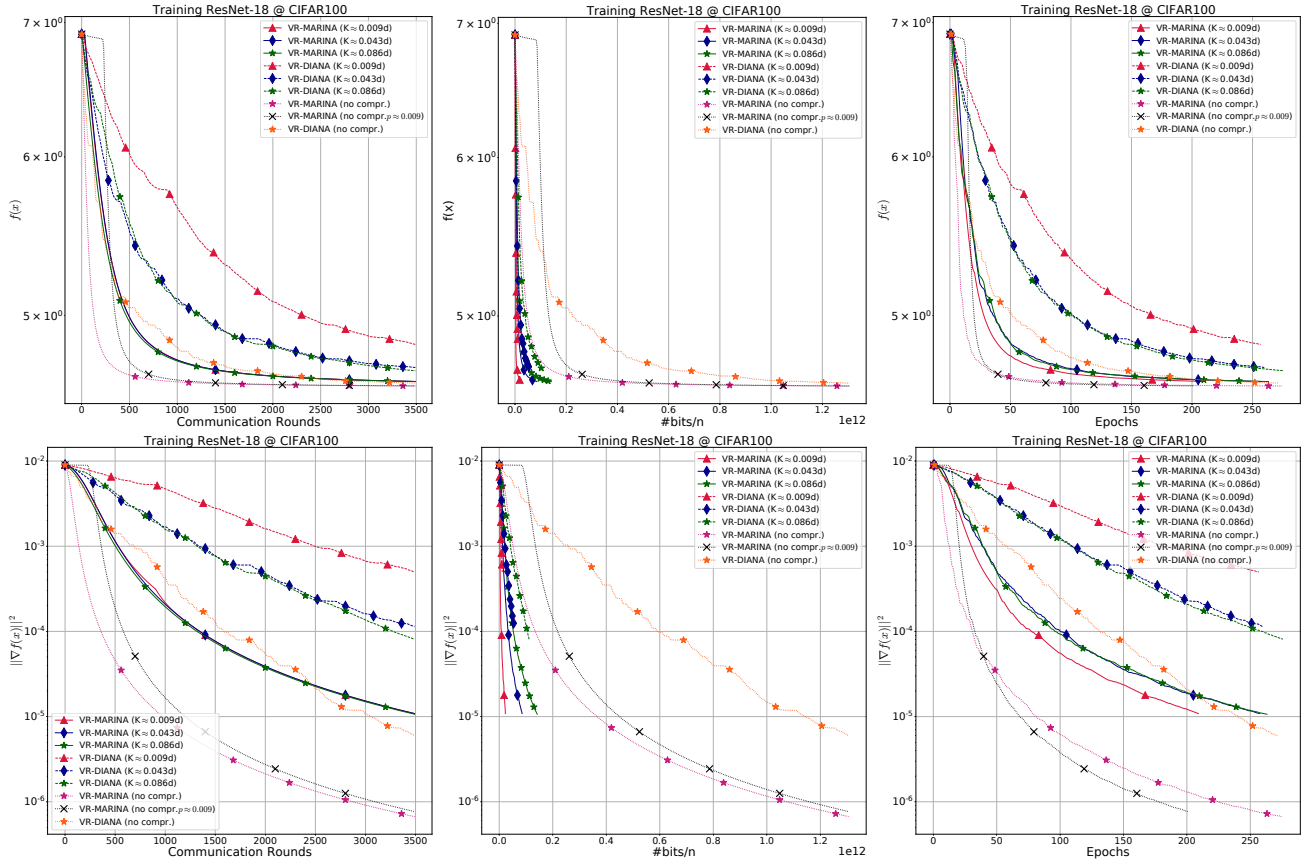


Figure 6: Comparison of VR-MARINA with VR-DIANA on training ResNet-18 at CIFAR100 dataset. Number of workers equals 5. Stepsizes for the methods were tuned and the batchsizes are  $\sim m/50$ . We used the RandK sparsification operator, the approximate values of  $K$  are given in the legends ( $d$  is dimension of the problem). We also show the performance of VR-MARINA and VR-DIANA without compression.

## B. Basic Facts and Auxiliary Results

### B.1. Useful Properties of Expectations

**Variance decomposition.** For a random vector  $\xi \in \mathbb{R}^d$  and any deterministic vector  $x \in \mathbb{R}^d$ , the variance can be decomposed as

$$\mathbf{E} \left[ \|\xi - \mathbf{E}\xi\|^2 \right] = \mathbf{E} [\|\xi - x\|^2] - \|\mathbf{E}\xi - x\|^2 \quad (13)$$

**Tower property of mathematical expectation.** For random variables  $\xi, \eta \in \mathbb{R}^d$ , we have

$$\mathbf{E} [\xi] = \mathbf{E} [\mathbf{E} [\xi \mid \eta]] \quad (14)$$

under an assumption that all expectations in the expression above are well-defined.

### B.2. One Lemma

In this section, we formulate a lemma from (Li et al., 2020), which holds in our settings as well. We omit the proof of this lemmas since it is identical to the one from (Li et al., 2020).

**Lemma B.1** (Lemma 2 from (Li et al., 2020)). *Assume that function  $f$  is  $L$ -smooth and  $x^{k+1} = x^k - \gamma g^k$ . Then*

$$f(x^{k+1}) \leq f(x^k) - \frac{\gamma}{2} \|\nabla f(x^k)\|^2 - \left( \frac{1}{2\gamma} - \frac{L}{2} \right) \|x^{k+1} - x^k\|^2 + \frac{\gamma}{2} \|g^k - \nabla f(x^k)\|^2. \quad (15)$$



## C. Missing Proofs for MARINA

### C.1. Generally Non-Convex Problems

In this section, we provide the full statement of Theorem 2.1 together with the proof of this result.

**Theorem C.1** (Theorem 2.1). *Let Assumptions 1.1 and 1.2 be satisfied and*

$$\gamma \leq \frac{1}{L \left( 1 + \sqrt{\frac{(1-p)\omega}{pn}} \right)}, \quad (16)$$

where  $L^2 = \frac{1}{n} \sum_{i=1}^n L_i^2$ . Then after  $K$  iterations of MARINA we have

$$\mathbf{E} \left[ \|\nabla f(\hat{x}^K)\|^2 \right] \leq \frac{2\Delta_0}{\gamma K}, \quad (17)$$

where  $\hat{x}^K$  is chosen uniformly at random from  $x^0, \dots, x^{K-1}$  and  $\Delta_0 = f(x^0) - f_*$ . That is, after

$$K = \mathcal{O} \left( \frac{\Delta_0 L}{\varepsilon^2} \left( 1 + \sqrt{\frac{(1-p)\omega}{pn}} \right) \right) \quad (18)$$

iterations MARINA produces such a point  $\hat{x}^K$  that  $\mathbf{E}[\|\nabla f(\hat{x}^K)\|^2] \leq \varepsilon^2$ . Moreover, under an assumption that the communication cost is proportional to the number of non-zero components of transmitted vectors from workers to the server, we have that the expected total communication cost per worker equals

$$d + K(pd + (1-p)\zeta_{\mathcal{Q}}) = \mathcal{O} \left( d + \frac{\Delta_0 L}{\varepsilon^2} \left( 1 + \sqrt{\frac{(1-p)\omega}{pn}} \right) (pd + (1-p)\zeta_{\mathcal{Q}}) \right), \quad (19)$$

where  $\zeta_{\mathcal{Q}}$  is the expected density of the quantization (see Def. 1.1).

*Proof of Theorem 2.1.* The scheme of the proof is similar to the proof of Theorem 1 from (Li et al., 2020). From Lemma B.1, we have

$$\mathbf{E}[f(x^{k+1})] \leq \mathbf{E}[f(x^k)] - \frac{\gamma}{2} \mathbf{E}[\|\nabla f(x^k)\|^2] - \left( \frac{1}{2\gamma} - \frac{L}{2} \right) \mathbf{E}[\|x^{k+1} - x^k\|^2] + \frac{\gamma}{2} \mathbf{E}[\|g^k - \nabla f(x^k)\|^2]. \quad (20)$$

Next, we need to derive an upper bound for  $\mathbf{E}[\|g^{k+1} - \nabla f(x^{k+1})\|^2]$ . By definition of  $g^{k+1}$ , we have

$$g^{k+1} = \begin{cases} \nabla f(x^{k+1}) & \text{with probability } p, \\ g^k + \frac{1}{n} \sum_{i=1}^n \mathcal{Q}(\nabla f_i(x^{k+1}) - \nabla f_i(x^k)) & \text{with probability } 1-p. \end{cases}$$

Using this, variance decomposition (13) and tower property (14), we derive:

$$\begin{aligned} \mathbf{E}[\|g^{k+1} - \nabla f(x^{k+1})\|^2] &\stackrel{(14)}{=} (1-p) \mathbf{E} \left[ \left\| g^k + \frac{1}{n} \sum_{i=1}^n \mathcal{Q}(\nabla f_i(x^{k+1}) - \nabla f_i(x^k)) - \nabla f(x^{k+1}) \right\|^2 \right] \\ &\stackrel{(14),(13)}{=} (1-p) \mathbf{E} \left[ \left\| \frac{1}{n} \sum_{i=1}^n \mathcal{Q}(\nabla f_i(x^{k+1}) - \nabla f_i(x^k)) - \nabla f(x^{k+1}) + \nabla f(x^k) \right\|^2 \right] \\ &\quad + (1-p) \mathbf{E}[\|g^k - \nabla f(x^k)\|^2]. \end{aligned}$$

Since  $\mathcal{Q}(\nabla f_1(x^{k+1}) - \nabla f_1(x^k)), \dots, \mathcal{Q}(\nabla f_n(x^{k+1}) - \nabla f_n(x^k))$  are independent random vectors for fixed  $x^k$  and  $x^{k+1}$  we have

$$\begin{aligned}
 \mathbf{E} [\|g^{k+1} - \nabla f(x^{k+1})\|^2] &= (1-p)\mathbf{E} \left[ \left\| \frac{1}{n} \sum_{i=1}^n (\mathcal{Q}(\nabla f_i(x^{k+1}) - \nabla f_i(x^k)) - \nabla f_i(x^{k+1}) + \nabla f_i(x^k)) \right\|^2 \right] \\
 &\quad + (1-p)\mathbf{E} [\|g^k - \nabla f(x^k)\|^2] \\
 &= \frac{1-p}{n^2} \sum_{i=1}^n \mathbf{E} [\|\mathcal{Q}(\nabla f_i(x^{k+1}) - \nabla f_i(x^k)) - \nabla f_i(x^{k+1}) + \nabla f_i(x^k)\|^2] \\
 &\quad + (1-p)\mathbf{E} [\|g^k - \nabla f(x^k)\|^2] \\
 &\stackrel{(3)}{\leq} \frac{(1-p)\omega}{n^2} \sum_{i=1}^n \mathbf{E} [\|\nabla f_i(x^{k+1}) - \nabla f_i(x^k)\|^2] + (1-p)\mathbf{E} [\|g^k - \nabla f(x^k)\|^2].
 \end{aligned}$$

Using  $L$ -smoothness (2) of  $f_i$  together with the tower property (14), we obtain

$$\begin{aligned}
 \mathbf{E} [\|g^{k+1} - \nabla f(x^{k+1})\|^2] &\leq \frac{(1-p)\omega}{n^2} \sum_{i=1}^n L_i^2 \mathbf{E} [\|x^{k+1} - x^k\|^2] + (1-p)\mathbf{E} [\|g^k - \nabla f(x^k)\|^2] \\
 &= \frac{(1-p)\omega L^2}{n} \mathbf{E} [\|x^{k+1} - x^k\|^2] + (1-p)\mathbf{E} [\|g^k - \nabla f(x^k)\|^2]. \tag{21}
 \end{aligned}$$

Next, we introduce a new notation:  $\Phi_k = f(x^k) - f_* + \frac{\gamma}{2p}\|g^k - \nabla f(x^k)\|^2$ . Using this and inequalities (20) and (21), we establish the following inequality:

$$\begin{aligned}
 \mathbf{E} [\Phi_{k+1}] &\leq \mathbf{E} \left[ f(x^k) - f_* - \frac{\gamma}{2}\|\nabla f(x^k)\|^2 - \left( \frac{1}{2\gamma} - \frac{L}{2} \right) \|x^{k+1} - x^k\|^2 + \frac{\gamma}{2}\|g^k - \nabla f(x^k)\|^2 \right] \\
 &\quad + \frac{\gamma}{2p} \mathbf{E} \left[ \frac{(1-p)\omega L^2}{n} \|x^{k+1} - x^k\|^2 + (1-p)\|g^k - \nabla f(x^k)\|^2 \right] \\
 &= \mathbf{E} [\Phi_k] - \frac{\gamma}{2} \mathbf{E} [\|\nabla f(x^k)\|^2] + \left( \frac{\gamma(1-p)\omega L^2}{2pn} - \frac{1}{2\gamma} + \frac{L}{2} \right) \mathbf{E} [\|x^{k+1} - x^k\|^2] \\
 &\stackrel{(16)}{\leq} \mathbf{E} [\Phi_k] - \frac{\gamma}{2} \mathbf{E} [\|\nabla f(x^k)\|^2], \tag{22}
 \end{aligned}$$

where in the last inequality, we use  $\frac{\gamma(1-p)\omega L^2}{2pn} - \frac{1}{2\gamma} + \frac{L}{2} \leq 0$  following from (16). Summing up inequalities (22) for  $k = 0, 1, \dots, K-1$  and rearranging the terms, we derive

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbf{E} [\|\nabla f(x^k)\|^2] \leq \frac{2}{\gamma K} \sum_{k=0}^{K-1} (\mathbf{E} [\Phi_k] - \mathbf{E} [\Phi_{k+1}]) = \frac{2(\mathbf{E} [\Phi_0] - \mathbf{E} [\Phi_K])}{\gamma K} = \frac{2\Delta_0}{\gamma K},$$

since  $g^0 = \nabla f(x^0)$  and  $\Phi_{k+1} \geq 0$ . Finally, using the tower property (14) and the definition of  $\hat{x}^K$ , we obtain (17) that implies (18) and (19).  $\square$

**Corollary C.1** (Corollary 2.1). *Let the assumptions of Theorem 2.1 hold and  $p = \frac{\zeta_{\mathcal{Q}}}{d}$ , where  $\zeta_{\mathcal{Q}}$  is the expected density of the quantization (see Def. 1.1). If*

$$\gamma \leq \frac{1}{L \left( 1 + \sqrt{\frac{\omega}{n} \left( \frac{d}{\zeta_{\mathcal{Q}}} - 1 \right)} \right)},$$

then MARINA requires

$$K = \mathcal{O} \left( \frac{\Delta_0 L}{\varepsilon^2} \left( 1 + \sqrt{\frac{\omega}{n} \left( \frac{d}{\zeta_{\mathcal{Q}}} - 1 \right)} \right) \right)$$

iterations/communication rounds to achieve  $\mathbf{E}[\|\nabla f(\hat{x}^K)\|^2] \leq \varepsilon^2$ , and the expected total communication cost per worker is

$$\mathcal{O} \left( d + \frac{\Delta_0 L}{\varepsilon^2} \left( \zeta_{\mathcal{Q}} + \sqrt{\frac{\omega \zeta_{\mathcal{Q}}}{n} (d - \zeta_{\mathcal{Q}})} \right) \right)$$

under an assumption that the communication cost is proportional to the number of non-zero components of transmitted vectors from workers to the server.

*Proof of Corollary 2.1.* The choice of  $p = \frac{\zeta_{\mathcal{Q}}}{d}$  implies

$$\begin{aligned} \frac{1-p}{p} &= \frac{d}{\zeta_{\mathcal{Q}}} - 1, \\ pd + (1-p)\zeta_{\mathcal{Q}} &\leq \zeta_{\mathcal{Q}} + \left(1 - \frac{\zeta_{\mathcal{Q}}}{d}\right) \cdot \zeta_{\mathcal{Q}} \leq 2\zeta_{\mathcal{Q}}. \end{aligned}$$

Plugging these relations in (16), (18), and (19), we get that if

$$\gamma \leq \frac{1}{L \left( 1 + \sqrt{\frac{\omega}{n} \left( \frac{d}{\zeta_{\mathcal{Q}}} - 1 \right)} \right)},$$

then MARINA requires

$$\begin{aligned} K &= \mathcal{O} \left( \frac{\Delta_0 L}{\varepsilon^2} \left( 1 + \sqrt{\frac{(1-p)\omega}{pn}} \right) \right) \\ &= \mathcal{O} \left( \frac{\Delta_0 L}{\varepsilon^2} \left( 1 + \sqrt{\frac{\omega}{n} \left( \frac{d}{\zeta_{\mathcal{Q}}} - 1 \right)} \right) \right) \end{aligned}$$

iterations/communication rounds in order to achieve  $\mathbf{E}[\|\nabla f(\hat{x}^K)\|^2] \leq \varepsilon^2$ , and the expected total communication cost per worker is

$$\begin{aligned} d + K(pd + (1-p)\zeta_{\mathcal{Q}}) &= \mathcal{O} \left( d + \frac{\Delta_0 L}{\varepsilon^2} \left( 1 + \sqrt{\frac{(1-p)\omega}{pn}} \right) (pd + (1-p)\zeta_{\mathcal{Q}}) \right) \\ &= \mathcal{O} \left( d + \frac{\Delta_0 L}{\varepsilon^2} \left( \zeta_{\mathcal{Q}} + \sqrt{\frac{\omega \zeta_{\mathcal{Q}}}{n} (d - \zeta_{\mathcal{Q}})} \right) \right) \end{aligned}$$

under an assumption that the communication cost is proportional to the number of non-zero components of transmitted vectors from workers to the server.  $\square$

### C.2. Convergence Results Under Polyak-Łojasiewicz condition

In this section, we provide the full statement of Theorem 2.2 together with the proof of this result.

**Theorem C.2** (Theorem 2.2). *Let Assumptions 1.1, 1.2 and 2.1 be satisfied and*

$$\gamma \leq \min \left\{ \frac{1}{L \left( 1 + \sqrt{\frac{2(1-p)\omega}{pn}} \right)}, \frac{p}{2\mu} \right\}, \quad (23)$$

where  $L^2 = \frac{1}{n} \sum_{i=1}^n L_i^2$ . Then after  $K$  iterations of MARINA we have

$$\mathbf{E} [f(x^K) - f(x^*)] \leq (1 - \gamma\mu)^K \Delta_0, \quad (24)$$

where  $\Delta_0 = f(x^0) - f(x^*)$ . That is, after

$$K = \mathcal{O} \left( \max \left\{ \frac{1}{p}, \frac{L}{\mu} \left( 1 + \sqrt{\frac{(1-p)\omega}{pn}} \right) \right\} \log \frac{\Delta_0}{\varepsilon} \right) \quad (25)$$

iterations MARINA produces such a point  $x^K$  that  $\mathbf{E}[f(x^K) - f(x^*)] \leq \varepsilon$ . Moreover, under an assumption that the communication cost is proportional to the number of non-zero components of transmitted vectors from workers to the server, we have that the expected total communication cost per worker equals

$$d + K(pd + (1-p)\zeta_{\mathcal{Q}}) = \mathcal{O} \left( d + \max \left\{ \frac{1}{p}, \frac{L}{\mu} \left( 1 + \sqrt{\frac{(1-p)\omega}{pn}} \right) \right\} (pd + (1-p)\zeta_{\mathcal{Q}}) \log \frac{\Delta_0}{\varepsilon} \right), \quad (26)$$

where  $\zeta_{\mathcal{Q}}$  is the expected density of the quantization (see Def. 1.1).

*Proof of Theorem 2.2.* The proof is very similar to the proof of Theorem 2.1. From Lemma B.1 and PL condition, we have

$$\begin{aligned} \mathbf{E}[f(x^{k+1}) - f(x^*)] &\leq \mathbf{E}[f(x^k) - f(x^*)] - \frac{\gamma}{2} \mathbf{E}[\|\nabla f(x^k)\|^2] - \left( \frac{1}{2\gamma} - \frac{L}{2} \right) \mathbf{E}[\|x^{k+1} - x^k\|^2] \\ &\quad + \frac{\gamma}{2} \mathbf{E}[\|g^k - \nabla f(x^k)\|^2] \\ &\stackrel{(4)}{\leq} (1 - \gamma\mu) \mathbf{E}[f(x^k) - f(x^*)] - \left( \frac{1}{2\gamma} - \frac{L}{2} \right) \mathbf{E}[\|x^{k+1} - x^k\|^2] + \frac{\gamma}{2} \mathbf{E}[\|g^k - \nabla f(x^k)\|^2]. \end{aligned}$$

Using the same arguments as in the proof of (21), we obtain

$$\mathbf{E}[\|g^{k+1} - \nabla f(x^{k+1})\|^2] \leq \frac{(1-p)\omega L^2}{n} \mathbf{E}[\|x^{k+1} - x^k\|^2] + (1-p) \mathbf{E}[\|g^k - \nabla f(x^k)\|^2].$$

Putting all together, we derive that the sequence  $\Phi_k = f(x^k) - f(x^*) + \frac{\gamma}{p} \|g^k - \nabla f(x^k)\|^2$  satisfies

$$\begin{aligned} \mathbf{E}[\Phi_{k+1}] &\leq \mathbf{E} \left[ (1 - \gamma\mu)(f(x^k) - f(x^*)) - \left( \frac{1}{2\gamma} - \frac{L}{2} \right) \|x^{k+1} - x^k\|^2 + \frac{\gamma}{2} \|g^k - \nabla f(x^k)\|^2 \right] \\ &\quad + \frac{\gamma}{p} \mathbf{E} \left[ \frac{(1-p)\omega L^2}{n} \|x^{k+1} - x^k\|^2 + (1-p) \|g^k - \nabla f(x^k)\|^2 \right] \\ &= \mathbf{E} \left[ (1 - \gamma\mu)(f(x^k) - f(x^*)) + \left( \frac{\gamma}{2} + \frac{\gamma}{p}(1-p) \right) \|g^k - \nabla f(x^k)\|^2 \right] \\ &\quad + \left( \frac{\gamma(1-p)\omega L^2}{pn} - \frac{1}{2\gamma} + \frac{L}{2} \right) \mathbf{E}[\|x^{k+1} - x^k\|^2] \\ &\stackrel{(23)}{\leq} (1 - \gamma\mu) \mathbf{E}[\Phi_k], \end{aligned}$$

where in the last inequality, we use  $\frac{\gamma(1-p)\omega L^2}{pn} - \frac{1}{2\gamma} + \frac{L}{2} \leq 0$  and  $\frac{\gamma}{2} + \frac{\gamma}{p}(1-p) \leq (1 - \gamma\mu)\frac{\gamma}{p}$  following from (23). Unrolling the recurrence and using  $g^0 = \nabla f(x^0)$ , we obtain

$$\mathbf{E}[f(x^K) - f(x^*)] \leq \mathbf{E}[\Phi_K] \leq (1 - \gamma\mu)^K \Phi_0 = (1 - \gamma\mu)^K (f(x^0) - f(x^*))$$

that implies (25) and (26).  $\square$

**Corollary C.2.** Let the assumptions of Theorem 2.2 hold and  $p = \frac{\zeta_{\mathcal{Q}}}{d}$ , where  $\zeta_{\mathcal{Q}}$  is the expected density of the quantization (see Def. 1.1). If

$$\gamma \leq \min \left\{ \frac{1}{L \left( 1 + \sqrt{\frac{2\omega}{n} \left( \frac{d}{\zeta_{\mathcal{Q}}} - 1 \right)} \right)}, \frac{p}{2\mu} \right\},$$

then MARINA requires

$$K = \mathcal{O} \left( \max \left\{ \frac{d}{\zeta_{\mathcal{Q}}}, \frac{L}{\mu} \left( 1 + \sqrt{\frac{\omega}{n} \left( \frac{d}{\zeta_{\mathcal{Q}}} - 1 \right)} \right) \right\} \log \frac{\Delta_0}{\varepsilon} \right)$$

iterations/communication rounds to achieve  $\mathbf{E}[f(x^K) - f(x^*)] \leq \varepsilon$ , and the expected total communication cost per worker is

$$\mathcal{O} \left( d + \max \left\{ d, \frac{L}{\mu} \left( \zeta_{\mathcal{Q}} + \sqrt{\frac{\omega \zeta_{\mathcal{Q}}}{n} (d - \zeta_{\mathcal{Q}})} \right) \right\} \log \frac{\Delta_0}{\varepsilon} \right)$$

under an assumption that the communication cost is proportional to the number of non-zero components of transmitted vectors from workers to the server.

*Proof.* The choice of  $p = \frac{\zeta_{\mathcal{Q}}}{d}$  implies

$$\begin{aligned} \frac{1-p}{p} &= \frac{d}{\zeta_{\mathcal{Q}}} - 1, \\ pd + (1-p)\zeta_{\mathcal{Q}} &\leq \zeta_{\mathcal{Q}} + \left(1 - \frac{\zeta_{\mathcal{Q}}}{d}\right) \cdot \zeta_{\mathcal{Q}} \leq 2\zeta_{\mathcal{Q}}. \end{aligned}$$

Plugging these relations in (23), (25), and (26), we get that if

$$\gamma \leq \min \left\{ \frac{1}{L \left( 1 + \sqrt{\frac{2\omega}{n} \left( \frac{d}{\zeta_{\mathcal{Q}}} - 1 \right)} \right)}, \frac{p}{2\mu} \right\},$$

then MARINA requires

$$\begin{aligned} K &= \mathcal{O} \left( \max \left\{ \frac{1}{p}, \frac{L}{\mu} \left( 1 + \sqrt{\frac{(1-p)\omega}{pn}} \right) \right\} \log \frac{\Delta_0}{\varepsilon} \right) \\ &= \mathcal{O} \left( \max \left\{ \frac{d}{\zeta_{\mathcal{Q}}}, \frac{L}{\mu} \left( 1 + \sqrt{\frac{\omega}{n} \left( \frac{d}{\zeta_{\mathcal{Q}}} - 1 \right)} \right) \right\} \log \frac{\Delta_0}{\varepsilon} \right) \end{aligned}$$

iterations/communication rounds in order to achieve  $\mathbf{E}[f(x^K) - f(x^*)] \leq \varepsilon$ , and the expected total communication cost per worker is

$$\begin{aligned} d + K(pd + (1-p)\zeta_{\mathcal{Q}}) &= \mathcal{O} \left( d + \max \left\{ \frac{1}{p}, \frac{L}{\mu} \left( 1 + \sqrt{\frac{(1-p)\omega}{pn}} \right) \right\} (pd + (1-p)\zeta_{\mathcal{Q}}) \log \frac{\Delta_0}{\varepsilon} \right) \\ &= \mathcal{O} \left( d + \max \left\{ d, \frac{L}{\mu} \left( \zeta_{\mathcal{Q}} + \sqrt{\frac{\omega \zeta_{\mathcal{Q}}}{n} (d - \zeta_{\mathcal{Q}})} \right) \right\} \log \frac{\Delta_0}{\varepsilon} \right) \end{aligned}$$

under an assumption that the communication cost is proportional to the number of non-zero components of transmitted vectors from workers to the server.  $\square$



## D. Missing Proofs for VR-MARINA

### D.1. Finite Sum Case

#### D.1.1. GENERALLY NON-CONVEX PROBLEMS

In this section, we provide the full statement of Theorem 3.1 together with the proof of this result.

**Theorem D.1** (Theorem 3.1). *Consider the finite sum case (1)+(5). Let Assumptions 1.1, 1.2 and 3.1 be satisfied and*

$$\gamma \leq \frac{1}{L + \sqrt{\frac{1-p}{pn} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right)}}, \quad (27)$$

where  $L^2 = \frac{1}{n} \sum_{i=1}^n L_i^2$  and  $\mathcal{L}^2 = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i^2$ . Then after  $K$  iterations of VR-MARINA we have

$$\mathbf{E} \left[ \|\nabla f(\hat{x}^K)\|^2 \right] \leq \frac{2\Delta_0}{\gamma K}, \quad (28)$$

where  $\hat{x}^K$  is chosen uniformly at random from  $x^0, \dots, x^{K-1}$  and  $\Delta_0 = f(x^0) - f_*$ . That is, after

$$K = \mathcal{O} \left( \frac{\Delta_0}{\varepsilon^2} \left( L + \sqrt{\frac{1-p}{pn} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right)} \right) \right) \quad (29)$$

iterations VR-MARINA produces such a point  $\hat{x}^K$  that  $\mathbf{E}[\|\nabla f(\hat{x}^K)\|^2] \leq \varepsilon^2$ , and the expected total number of stochastic oracle calls per node equals

$$m + K(pm + 2(1-p)b') = \mathcal{O} \left( m + \frac{\Delta_0}{\varepsilon^2} \left( L + \sqrt{\frac{1-p}{pn} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right)} \right) (pm + (1-p)b') \right). \quad (30)$$

Moreover, under an assumption that the communication cost is proportional to the number of non-zero components of transmitted vectors from workers to the server, we have that the expected total communication cost per worker equals

$$d + K(pd + (1-p)\zeta_{\mathcal{Q}}) = \mathcal{O} \left( d + \frac{\Delta_0}{\varepsilon^2} \left( L + \sqrt{\frac{1-p}{pn} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right)} \right) (pd + (1-p)\zeta_{\mathcal{Q}}) \right), \quad (31)$$

where  $\zeta_{\mathcal{Q}}$  is the expected density of the quantization (see Def. 1.1).

*Proof of Theorem 3.1.* The proof of this theorem is a generalization of the proof of Theorem 2.1. From Lemma B.1, we have

$$\mathbf{E}[f(x^{k+1})] \leq \mathbf{E}[f(x^k)] - \frac{\gamma}{2} \mathbf{E}[\|\nabla f(x^k)\|^2] - \left( \frac{1}{2\gamma} - \frac{L}{2} \right) \mathbf{E}[\|x^{k+1} - x^k\|^2] + \frac{\gamma}{2} \mathbf{E}[\|g^k - \nabla f(x^k)\|^2]. \quad (32)$$

Next, we need to derive an upper bound for  $\mathbf{E}[\|g^{k+1} - \nabla f(x^{k+1})\|^2]$ . Since  $g^{k+1} = \frac{1}{n} \sum_{i=1}^n g_i^{k+1}$ , we get the following representation of  $g^{k+1}$ :

$$g^{k+1} = \begin{cases} \nabla f(x^{k+1}) & \text{with probability } p, \\ g^k + \frac{1}{n} \sum_{i=1}^n \mathcal{Q} \left( \frac{1}{b'} \sum_{j \in I'_{i,k}} (\nabla f_{ij}(x^{k+1}) - \nabla f_{ij}(x^k)) \right) & \text{with probability } 1-p. \end{cases}$$

Using this, variance decomposition (13) and tower property (14), we derive:

$$\begin{aligned}
 \mathbf{E} [\|g^{k+1} - \nabla f(x^{k+1})\|^2] &\stackrel{(14)}{=} (1-p)\mathbf{E} \left[ \left\| g^k + \frac{1}{n} \sum_{i=1}^n \mathcal{Q} \left( \frac{1}{b'} \sum_{j \in I'_{i,k}} (\nabla f_{ij}(x^{k+1}) - \nabla f_{ij}(x^k)) \right) - \nabla f(x^{k+1}) \right\|^2 \right] \\
 &\stackrel{(14),(13)}{=} (1-p)\mathbf{E} \left[ \left\| \frac{1}{n} \sum_{i=1}^n \mathcal{Q} \left( \frac{1}{b'} \sum_{j \in I'_{i,k}} (\nabla f_{ij}(x^{k+1}) - \nabla f_{ij}(x^k)) \right) - \nabla f(x^{k+1}) + \nabla f(x^k) \right\|^2 \right] \\
 &\quad + (1-p)\mathbf{E} [\|g^k - \nabla f(x^k)\|^2].
 \end{aligned}$$

Next, we use the notation:  $\tilde{\Delta}_i^k = \frac{1}{b'} \sum_{j \in I'_{i,k}} (\nabla f_{ij}(x^{k+1}) - \nabla f_{ij}(x^k))$  and  $\Delta_i^k = \nabla f_i(x^{k+1}) - \nabla f_i(x^k)$ . These vectors satisfy  $\mathbf{E} [\tilde{\Delta}_i^k | x^k, x^{k+1}] = \Delta_i^k$  for all  $i \in [n]$ . Moreover,  $\mathcal{Q}(\tilde{\Delta}_1^k), \dots, \mathcal{Q}(\tilde{\Delta}_n^k)$  are independent random vectors for fixed  $x^k$  and  $x^{k+1}$ . These observations imply

$$\begin{aligned}
 \mathbf{E} [\|g^{k+1} - \nabla f(x^{k+1})\|^2] &= (1-p)\mathbf{E} \left[ \left\| \frac{1}{n} \sum_{i=1}^n (\mathcal{Q}(\tilde{\Delta}_i^k) - \Delta_i^k) \right\|^2 \right] + (1-p)\mathbf{E} [\|g^k - \nabla f(x^k)\|^2] \\
 &= \frac{1-p}{n^2} \sum_{i=1}^n \mathbf{E} \left[ \left\| \mathcal{Q}(\tilde{\Delta}_i^k) - \tilde{\Delta}_i^k + \tilde{\Delta}_i^k - \Delta_i^k \right\|^2 \right] + (1-p)\mathbf{E} [\|g^k - \nabla f(x^k)\|^2] \\
 &\stackrel{(14),(13)}{=} \frac{1-p}{n^2} \sum_{i=1}^n \left( \mathbf{E} \left[ \left\| \mathcal{Q}(\tilde{\Delta}_i^k) - \tilde{\Delta}_i^k \right\|^2 \right] + \mathbf{E} \left[ \left\| \tilde{\Delta}_i^k - \Delta_i^k \right\|^2 \right] \right) \\
 &\quad + (1-p)\mathbf{E} [\|g^k - \nabla f(x^k)\|^2] \\
 &\stackrel{(14),(3)}{=} \frac{1-p}{n^2} \sum_{i=1}^n \left( \omega \mathbf{E} \left[ \left\| \tilde{\Delta}_i^k \right\|^2 \right] + \mathbf{E} \left[ \left\| \tilde{\Delta}_i^k - \Delta_i^k \right\|^2 \right] \right) + (1-p)\mathbf{E} [\|g^k - \nabla f(x^k)\|^2] \\
 &\stackrel{(14),(13)}{=} \frac{1-p}{n^2} \sum_{i=1}^n \left( \omega \mathbf{E} \left[ \left\| \Delta_i^k \right\|^2 \right] + (1+\omega) \mathbf{E} \left[ \left\| \tilde{\Delta}_i^k - \Delta_i^k \right\|^2 \right] \right) \\
 &\quad + (1-p)\mathbf{E} [\|g^k - \nabla f(x^k)\|^2].
 \end{aligned}$$

Using  $L$ -smoothness (2) and average  $\mathcal{L}$ -smoothness (7) of  $f_i$  together with the tower property (14), we get

$$\begin{aligned}
 \mathbf{E} [\|g^{k+1} - \nabla f(x^{k+1})\|^2] &\leq \frac{1-p}{n^2} \sum_{i=1}^n \left( \omega L_i^2 + \frac{(1+\omega)\mathcal{L}_i^2}{b'} \right) \mathbf{E} [\|x^{k+1} - x^k\|^2] \\
 &\quad + (1-p)\mathbf{E} [\|g^k - \nabla f(x^k)\|^2] \\
 &= \frac{1-p}{n} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right) \mathbf{E} [\|x^{k+1} - x^k\|^2] \\
 &\quad + (1-p)\mathbf{E} [\|g^k - \nabla f(x^k)\|^2]. \tag{33}
 \end{aligned}$$

Next, we introduce new notation:  $\Phi_k = f(x^k) - f_* + \frac{\gamma}{2p} \|g^k - \nabla f(x^k)\|^2$ . Using this and inequalities (32) and (33), we

establish the following inequality:

$$\begin{aligned}
 \mathbf{E}[\Phi_{k+1}] &\leq \mathbf{E}\left[f(x^k) - f_* - \frac{\gamma}{2}\|\nabla f(x^k)\|^2 - \left(\frac{1}{2\gamma} - \frac{L}{2}\right)\|x^{k+1} - x^k\|^2 + \frac{\gamma}{2}\|g^k - \nabla f(x^k)\|^2\right] \\
 &\quad + \frac{\gamma}{2p}\mathbf{E}\left[\frac{1-p}{n}\left(\omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'}\right)\|x^{k+1} - x^k\|^2 + (1-p)\|g^k - \nabla f(x^k)\|^2\right] \\
 &= \mathbf{E}[\Phi_k] - \frac{\gamma}{2}\mathbf{E}[\|\nabla f(x^k)\|^2] + \left(\frac{\gamma(1-p)}{2pn}\left(\omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'}\right) - \frac{1}{2\gamma} + \frac{L}{2}\right)\mathbf{E}[\|x^{k+1} - x^k\|^2] \\
 &\stackrel{(27)}{\leq} \mathbf{E}[\Phi_k] - \frac{\gamma}{2}\mathbf{E}[\|\nabla f(x^k)\|^2], \tag{34}
 \end{aligned}$$

where in the last inequality, we use  $\frac{\gamma(1-p)}{2pn}\left(\omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'}\right) - \frac{1}{2\gamma} + \frac{L}{2} \leq 0$  following from (27). Summing up inequalities (34) for  $k = 0, 1, \dots, K-1$  and rearranging the terms, we derive

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbf{E}[\|\nabla f(x^k)\|^2] \leq \frac{2}{\gamma K} \sum_{k=0}^{K-1} (\mathbf{E}[\Phi_k] - \mathbf{E}[\Phi_{k+1}]) = \frac{2(\mathbf{E}[\Phi_0] - \mathbf{E}[\Phi_K])}{\gamma K} = \frac{2\Delta_0}{\gamma K},$$

since  $g^0 = \nabla f(x^0)$  and  $\Phi_{k+1} \geq 0$ . Finally, using the tower property (14) and the definition of  $\hat{x}^K$ , we obtain (28) that implies (29), (30), and (31).  $\square$

**Remark D.1** (About batchsizes dissimilarity). *We notice that our analysis can be easily extended to handle the version of VR-MARINA with different batchsizes  $b'_1, \dots, b'_n$  on different workers, i.e., when  $|I'_{i,k}| = b'_i$  and  $\tilde{\Delta}_i^k = \frac{1}{b'_i} \sum_{j \in I'_{i,k}} (\nabla f_{ij}(x^{k+1}) - \nabla f_{ij}(x^k))$ . In this case, the statement of Theorem 3.1 remains the same with the small modification: instead of  $\frac{\mathcal{L}^2}{b'}$  the complexity bounds will have  $\frac{1}{n} \sum_{i=1}^n \frac{\mathcal{L}_i^2}{b'_i}$ .*

**Corollary D.1** (Corollary 3.1). *Let the assumptions of Theorem 3.1 hold and  $p = \min\left\{\frac{\zeta_{\mathcal{Q}}}{d}, \frac{b'}{m+b'}\right\}$ , where  $b' \leq m$  and  $\zeta_{\mathcal{Q}}$  is the expected density of the quantization (see Def. 1.1). If*

$$\gamma \leq \frac{1}{L + \sqrt{\frac{\max\{d/\zeta_{\mathcal{Q}} - 1, m/b'\}}{n} \left(\omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'}\right)}},$$

then VR-MARINA requires

$$\mathcal{O}\left(\frac{\Delta_0}{\varepsilon^2} \left(L \left(1 + \sqrt{\frac{\omega \max\{d/\zeta_{\mathcal{Q}} - 1, m/b'\}}{n}}\right) + \mathcal{L} \sqrt{\frac{(1+\omega) \max\{d/\zeta_{\mathcal{Q}} - 1, m/b'\}}{nb'}}\right)\right)$$

iterations/communication rounds,

$$\mathcal{O}\left(m + \frac{\Delta_0}{\varepsilon^2} \left(L \left(b' + \sqrt{\frac{\omega \max\{(d/\zeta_{\mathcal{Q}} - 1)(b')^2, mb'\}}{n}}\right) + \mathcal{L} \sqrt{\frac{(1+\omega) \max\{(d/\zeta_{\mathcal{Q}} - 1)b', m\}}{n}}\right)\right)$$

stochastic oracle calls per node in expectation in order to achieve  $\mathbf{E}[\|\nabla f(\hat{x}^K)\|^2] \leq \varepsilon^2$ , and the expected total communication cost per worker is

$$\mathcal{O}\left(d + \frac{\Delta_0 \zeta_{\mathcal{Q}}}{\varepsilon^2} \left(L \left(1 + \sqrt{\frac{\omega \max\{d/\zeta_{\mathcal{Q}} - 1, m/b'\}}{n}}\right) + \mathcal{L} \sqrt{\frac{(1+\omega) \max\{d/\zeta_{\mathcal{Q}} - 1, m/b'\}}{nb'}}\right)\right)$$

under an assumption that the communication cost is proportional to the number of non-zero components of transmitted vectors from workers to the server.

*Proof of Corollary 3.1.* The choice of  $p = \min \left\{ \frac{\zeta_{\mathcal{Q}}}{d}, \frac{b'}{m+b'} \right\}$  implies

$$\begin{aligned} \frac{1-p}{p} &= \max \left\{ \frac{d}{\zeta_{\mathcal{Q}}} - 1, \frac{m}{b'} \right\}, \\ pm + (1-p)b' &\leq \frac{2mb'}{m+b'} \leq 2b', \\ pd + (1-p)\zeta_{\mathcal{Q}} &\leq \frac{\zeta_{\mathcal{Q}}}{d} \cdot d + \left(1 - \frac{\zeta_{\mathcal{Q}}}{d}\right) \cdot \zeta_{\mathcal{Q}} \leq 2\zeta_{\mathcal{Q}}. \end{aligned}$$

Plugging these relations in (27), (29), (30) and (31) and using  $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$ , we get that if

$$\gamma \leq \frac{1}{L + \sqrt{\frac{\max\{d/\zeta_{\mathcal{Q}}-1, m/b'\}}{n} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right)}},$$

then VR-MARINA requires

$$\begin{aligned} K &= \mathcal{O} \left( \frac{\Delta_0}{\varepsilon^2} \left( L + \sqrt{\frac{1-p}{pn} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right)} \right) \right) \\ &= \mathcal{O} \left( \frac{\Delta_0}{\varepsilon^2} \left( L + \sqrt{L^2 \frac{\omega \max\{d/\zeta_{\mathcal{Q}}-1, m/b'\}}{n} + \mathcal{L}^2 \frac{(1+\omega) \max\{d/\zeta_{\mathcal{Q}}-1, m/b'\}}{nb'}} \right) \right) \\ &= \mathcal{O} \left( \frac{\Delta_0}{\varepsilon^2} \left( L \left( 1 + \sqrt{\frac{\omega \max\{d/\zeta_{\mathcal{Q}}-1, m/b'\}}{n}} \right) + \mathcal{L} \sqrt{\frac{(1+\omega) \max\{d/\zeta_{\mathcal{Q}}-1, m/b'\}}{nb'}} \right) \right) \end{aligned}$$

iterations/communication rounds and

$$\begin{aligned} m + K(pm + 2(1-p)b') &= \mathcal{O} \left( m + \frac{\Delta_0}{\varepsilon^2} \left( L + \sqrt{\frac{1-p}{pn} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right)} \right) (pm + (1-p)b') \right) \\ &= \mathcal{O} \left( m + \frac{\Delta_0}{\varepsilon^2} \left( L \left( 1 + \sqrt{\frac{\omega \max\{d/\zeta_{\mathcal{Q}}-1, m/b'\}}{n}} \right) \right. \right. \\ &\quad \left. \left. + \mathcal{L} \sqrt{\frac{(1+\omega) \max\{d/\zeta_{\mathcal{Q}}-1, m/b'\}}{nb'}} \right) b' \right) \\ &= \mathcal{O} \left( m + \frac{\Delta_0}{\varepsilon^2} \left( L \left( b' + \sqrt{\frac{\omega \max\{(d/\zeta_{\mathcal{Q}}-1)(b')^2, mb'\}}{n}} \right) \right. \right. \\ &\quad \left. \left. + \mathcal{L} \sqrt{\frac{(1+\omega) \max\{(d/\zeta_{\mathcal{Q}}-1)b', m\}}{n}} \right) \right) \end{aligned}$$

stochastic oracle calls per node in expectation in order to achieve  $\mathbf{E}[\|\nabla f(\hat{x}^K)\|^2] \leq \varepsilon^2$ , and the expected total communication cost per worker is

$$\begin{aligned} d + K(pd + (1-p)\zeta_{\mathcal{Q}}) &= \mathcal{O} \left( d + \frac{\Delta_0}{\varepsilon^2} \left( L + \sqrt{\frac{1-p}{pn} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right)} \right) (pd + (1-p)\zeta_{\mathcal{Q}}) \right) \\ &= \mathcal{O} \left( d + \frac{\Delta_0 \zeta_{\mathcal{Q}}}{\varepsilon^2} \left( L \left( 1 + \sqrt{\frac{\omega \max\{d/\zeta_{\mathcal{Q}}-1, m/b'\}}{n}} \right) \right. \right. \\ &\quad \left. \left. + \mathcal{L} \sqrt{\frac{(1+\omega) \max\{d/\zeta_{\mathcal{Q}}-1, m/b'\}}{nb'}} \right) \right) \end{aligned}$$

under an assumption that the communication cost is proportional to the number of non-zero components of transmitted vectors from workers to the server.  $\square$

## D.1.2. CONVERGENCE RESULTS UNDER POLYAK-ŁOJASIEWICZ CONDITION

In this section, we provide an analysis of VR-MARINA under the Polyak-Łojasiewicz condition in the finite sum case.

**Theorem D.2.** *Consider the finite sum case (1)+(5). Let Assumptions 1.1, 1.2, 3.1 and 2.1 be satisfied and*

$$\gamma \leq \min \left\{ \frac{1}{L + \sqrt{\frac{2(1-p)}{pn} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right)}}, \frac{p}{2\mu} \right\}, \quad (35)$$

where  $L^2 = \frac{1}{n} \sum_{i=1}^n L_i^2$  and  $\mathcal{L}^2 = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i^2$ . Then after  $K$  iterations of VR-MARINA, we have

$$\mathbf{E} [f(x^K) - f(x^*)] \leq (1 - \gamma\mu)^K \Delta_0, \quad (36)$$

where  $\Delta_0 = f(x^0) - f(x^*)$ . That is, after

$$K = \mathcal{O} \left( \max \left\{ \frac{1}{p}, \frac{L + \sqrt{\frac{1-p}{pn} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right)}}{\mu} \right\} \log \frac{\Delta_0}{\varepsilon} \right) \quad (37)$$

iterations VR-MARINA produces such a point  $x^K$  that  $\mathbf{E} [f(x^K) - f(x^*)] \leq \varepsilon$ , and the expected total number of stochastic oracle calls per node equals

$$m + K(pm + 2(1-p)b') = \mathcal{O} \left( m + \max \left\{ \frac{1}{p}, \frac{L + \sqrt{\frac{1-p}{pn} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right)}}{\mu} \right\} (pm + (1-p)b') \log \frac{\Delta_0}{\varepsilon} \right). \quad (38)$$

Moreover, under an assumption that the communication cost is proportional to the number of non-zero components of transmitted vectors from workers to the server we have that the expected total communication cost per worker equals

$$d + K(pd + (1-p)\zeta_{\mathcal{Q}}) = \mathcal{O} \left( d + \max \left\{ \frac{1}{p}, \frac{L + \sqrt{\frac{1-p}{pn} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right)}}{\mu} \right\} (pd + (1-p)\zeta_{\mathcal{Q}}) \log \frac{\Delta_0}{\varepsilon} \right), \quad (39)$$

where  $\zeta_{\mathcal{Q}}$  is the expected density of the quantization (see Def. 1.1).

*Proof.* The proof is very similar to the proof of Theorem 3.1. From Lemma B.1 and PL condition, we have

$$\begin{aligned} \mathbf{E}[f(x^{k+1}) - f(x^*)] &\leq \mathbf{E}[f(x^k) - f(x^*)] - \frac{\gamma}{2} \mathbf{E} [\|\nabla f(x^k)\|^2] - \left( \frac{1}{2\gamma} - \frac{L}{2} \right) \mathbf{E} [\|x^{k+1} - x^k\|^2] \\ &\quad + \frac{\gamma}{2} \mathbf{E} [\|g^k - \nabla f(x^k)\|^2] \\ &\stackrel{(4)}{\leq} (1 - \gamma\mu) \mathbf{E} [f(x^k) - f(x^*)] - \left( \frac{1}{2\gamma} - \frac{L}{2} \right) \mathbf{E} [\|x^{k+1} - x^k\|^2] + \frac{\gamma}{2} \mathbf{E} [\|g^k - \nabla f(x^k)\|^2]. \end{aligned}$$

Using the same arguments as in the proof of (33), we obtain

$$\mathbf{E} [\|g^{k+1} - \nabla f(x^{k+1})\|^2] \leq \frac{1-p}{n} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right) \mathbf{E} [\|x^{k+1} - x^k\|^2] + (1-p) \mathbf{E} [\|g^k - \nabla f(x^k)\|^2].$$

Putting all together we derive that the sequence  $\Phi_k = f(x^k) - f(x^*) + \frac{\gamma}{p} \|g^k - \nabla f(x^k)\|^2$  satisfies

$$\begin{aligned}
 \mathbf{E}[\Phi_{k+1}] &\leq \mathbf{E} \left[ (1 - \gamma\mu)(f(x^k) - f(x^*)) - \left( \frac{1}{2\gamma} - \frac{L}{2} \right) \|x^{k+1} - x^k\|^2 + \frac{\gamma}{2} \|g^k - \nabla f(x^k)\|^2 \right] \\
 &\quad + \frac{\gamma}{p} \mathbf{E} \left[ \frac{1-p}{n} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right) \|x^{k+1} - x^k\|^2 + (1-p) \|g^k - \nabla f(x^k)\|^2 \right] \\
 &= \mathbf{E} \left[ (1 - \gamma\mu)(f(x^k) - f(x^*)) + \left( \frac{\gamma}{2} + \frac{\gamma(1-p)}{p} \right) \|g^k - \nabla f(x^k)\|^2 \right] \\
 &\quad + \left( \frac{\gamma(1-p)}{pn} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right) - \frac{1}{2\gamma} + \frac{L}{2} \right) \mathbf{E} [\|x^{k+1} - x^k\|^2] \\
 &\stackrel{(35)}{\leq} (1 - \gamma\mu) \mathbf{E}[\Phi_k],
 \end{aligned}$$

where in the last inequality we use  $\frac{\gamma(1-p)}{pn} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right) - \frac{1}{2\gamma} + \frac{L}{2} \leq 0$  and  $\frac{\gamma}{2} + \frac{\gamma(1-p)}{p} \leq (1 - \gamma\mu) \frac{\gamma}{p}$  following from (35). Unrolling the recurrence and using  $g^0 = \nabla f(x^0)$ , we obtain

$$\mathbf{E} [f(x^{k+1}) - f(x^*)] \leq \mathbf{E}[\Phi_{k+1}] \leq (1 - \gamma\mu)^{k+1} \Phi_0 = (1 - \gamma\mu)^{k+1} (f(x^0) - f(x^*))$$

that implies (37), (38), and (39).  $\square$

**Corollary D.2.** *Let the assumptions of Theorem D.2 hold and  $p = \min \left\{ \frac{\zeta_{\mathcal{Q}}}{d}, \frac{b'}{m+b'} \right\}$ , where  $b' \leq m$  and  $\zeta_{\mathcal{Q}}$  is the expected density of the quantization (see Def. 1.1). If*

$$\gamma \leq \min \left\{ \frac{1}{L + \sqrt{\frac{2 \max\{d/\zeta_{\mathcal{Q}} - 1, m/b'\}}{n} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right)}}, \frac{p}{2\mu} \right\},$$

then VR-MARINA requires

$$\mathcal{O} \left( \max \left\{ \frac{1}{p}, \frac{L}{\mu} \left( 1 + \sqrt{\frac{\omega \max\{d/\zeta_{\mathcal{Q}} - 1, m/b'\}}{n}} \right) + \frac{\mathcal{L}}{\mu} \sqrt{\frac{(1+\omega) \max\{d/\zeta_{\mathcal{Q}} - 1, m/b'\}}{nb'}} \right\} \log \frac{\Delta_0}{\varepsilon} \right)$$

iterations/communication rounds,

$$\mathcal{O} \left( m + \max \left\{ \frac{b'}{p}, \frac{L}{\mu} \left( b' + \sqrt{\frac{\omega \max\{(d/\zeta_{\mathcal{Q}} - 1)(b')^2, mb'\}}{n}} \right) + \frac{\mathcal{L}}{\mu} \sqrt{\frac{(1+\omega) \max\{(d/\zeta_{\mathcal{Q}} - 1)b', m\}}{n}} \right\} \log \frac{\Delta_0}{\varepsilon} \right)$$

stochastic oracle calls per node in expectation to achieve  $\mathbf{E}[f(x^K) - f(x^*)] \leq \varepsilon$ , and the expected total communication cost per worker is

$$\mathcal{O} \left( d + \zeta_{\mathcal{Q}} \max \left\{ \frac{1}{p}, \frac{L}{\mu} \left( 1 + \sqrt{\frac{\omega \max\{d/\zeta_{\mathcal{Q}} - 1, m/b'\}}{n}} \right) + \frac{\mathcal{L}}{\mu} \sqrt{\frac{(1+\omega) \max\{d/\zeta_{\mathcal{Q}} - 1, m/b'\}}{nb'}} \right\} \log \frac{\Delta_0}{\varepsilon} \right)$$

under an assumption that the communication cost is proportional to the number of non-zero components of transmitted vectors from workers to the server.

*Proof.* The choice of  $p = \min \left\{ \frac{\zeta_{\mathcal{Q}}}{d}, \frac{b'}{m+b'} \right\}$  implies

$$\begin{aligned}
 \frac{1-p}{p} &= \max \left\{ \frac{d}{\zeta_{\mathcal{Q}}} - 1, \frac{m}{b'} \right\}, \\
 pm + (1-p)b' &\leq \frac{2mb'}{m+b'} \leq 2b', \\
 pd + (1-p)\zeta_{\mathcal{Q}} &\leq \frac{\zeta_{\mathcal{Q}}}{d} \cdot d + \left( 1 - \frac{\zeta_{\mathcal{Q}}}{d} \right) \cdot \zeta_{\mathcal{Q}} \leq 2\zeta_{\mathcal{Q}}.
 \end{aligned}$$



Plugging these relations in (35), (37), (38) and (39) and using  $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$ , we get that if

$$\gamma \leq \min \left\{ \frac{1}{L + \sqrt{\frac{2 \max\{d/\zeta_{\mathcal{Q}} - 1, m/b'\}}{n} (\omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'})}}, \frac{p}{2\mu} \right\},$$

then VR-MARINA requires

$$\begin{aligned} K &= \mathcal{O} \left( \max \left\{ \frac{1}{p}, \frac{L + \sqrt{\frac{1-p}{pn} (\omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'})}}{\mu} \right\} \log \frac{\Delta_0}{\varepsilon} \right) \\ &= \mathcal{O} \left( \max \left\{ \frac{1}{p}, \frac{L + \sqrt{L^2 \frac{\omega \max\{d/\zeta_{\mathcal{Q}} - 1, m/b'\}}{n} + \mathcal{L}^2 \frac{(1+\omega) \max\{d/\zeta_{\mathcal{Q}} - 1, m/b'\}}{nb'}}}{\mu} \right\} \log \frac{\Delta_0}{\varepsilon} \right) \\ &= \mathcal{O} \left( \max \left\{ \frac{1}{p}, \frac{L}{\mu} \left( 1 + \sqrt{\frac{\omega \max\{d/\zeta_{\mathcal{Q}} - 1, m/b'\}}{n}} \right) + \frac{\mathcal{L}}{\mu} \sqrt{\frac{(1+\omega) \max\{d/\zeta_{\mathcal{Q}} - 1, m/b'\}}{nb'}} \right\} \log \frac{\Delta_0}{\varepsilon} \right) \end{aligned}$$

iterations/communication rounds and

$$\begin{aligned} m + K(pm + 2(1-p)b') &= \mathcal{O} \left( m + \max \left\{ \frac{1}{p}, \frac{L + \sqrt{\frac{1-p}{pn} (\omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'})}}{\mu} \right\} (pm + (1-p)b') \log \frac{\Delta_0}{\varepsilon} \right) \\ &= \mathcal{O} \left( m + \max \left\{ \frac{1}{p}, \frac{L}{\mu} \left( 1 + \sqrt{\frac{\omega \max\{d/\zeta_{\mathcal{Q}} - 1, m/b'\}}{n}} \right) \right. \right. \\ &\quad \left. \left. + \frac{\mathcal{L}}{\mu} \sqrt{\frac{(1+\omega) \max\{d/\zeta_{\mathcal{Q}} - 1, m/b'\}}{nb'}} \right\} b' \log \frac{\Delta_0}{\varepsilon} \right) \\ &= \mathcal{O} \left( m + \max \left\{ \frac{b'}{p}, \frac{L}{\mu} \left( b' + \sqrt{\frac{\omega \max\{(d/\zeta_{\mathcal{Q}} - 1)(b')^2, mb'\}}{n}} \right) \right. \right. \\ &\quad \left. \left. + \frac{\mathcal{L}}{\mu} \sqrt{\frac{(1+\omega) \max\{(d/\zeta_{\mathcal{Q}} - 1)b', m\}}{n}} \right\} \log \frac{\Delta_0}{\varepsilon} \right) \end{aligned}$$

stochastic oracle calls per node in expectation in order to achieve  $\mathbf{E}[f(x^K) - f(x^*)] \leq \varepsilon$ , and the expected total communication cost per worker is

$$\begin{aligned} d + K(pd + (1-p)\zeta_{\mathcal{Q}}) &= \mathcal{O} \left( d + \max \left\{ \frac{1}{p}, \frac{L + \sqrt{\frac{1-p}{pn} (\omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'})}}{\mu} \right\} (pd + (1-p)\zeta_{\mathcal{Q}}) \log \frac{\Delta_0}{\varepsilon} \right) \\ &= \mathcal{O} \left( d + \zeta_{\mathcal{Q}} \max \left\{ \frac{1}{p}, \frac{L}{\mu} \left( 1 + \sqrt{\frac{\omega \max\{d/\zeta_{\mathcal{Q}} - 1, m/b'\}}{n}} \right) \right. \right. \\ &\quad \left. \left. + \frac{\mathcal{L}}{\mu} \sqrt{\frac{(1+\omega) \max\{d/\zeta_{\mathcal{Q}} - 1, m/b'\}}{nb'}} \right\} \log \frac{\Delta_0}{\varepsilon} \right) \end{aligned}$$

under an assumption that the communication cost is proportional to the number of non-zero components of transmitted vectors from workers to the server.  $\square$

## D.2. Online Case

**Algorithm 3** VR-MARINA: online case

- 1: **Input:** starting point  $x^0$ , stepsize  $\gamma$ , minibatch sizes  $b, b' < b$ , probability  $p \in (0, 1]$ , number of iterations  $K$
- 2: Initialize  $g^0 = \frac{1}{nb} \sum_{i=1}^n \sum_{j \in I_{i,0}} \nabla f_{\xi_{ij}^0}(x^{k+1})$ , where  $I_{i,0}$  is the set of the indices in the minibatch,  $|I_{i,0}| = b$ , and  $\xi_{ij}^0$  is independently sampled from  $\mathcal{D}_i$  for  $i \in [n], j \in [m]$
- 3: **for**  $k = 0, 1, \dots, K - 1$  **do**
- 4:   Sample  $c_k \sim \text{Be}(p)$
- 5:   Broadcast  $g^k$  to all workers
- 6:   **for**  $i = 1, \dots, n$  **in parallel do**
- 7:      $x^{k+1} = x^k - \gamma g^k$
- 8:     Set  $g_i^{k+1} = \begin{cases} \frac{1}{b} \sum_{j \in I_{i,k}} \nabla f_{\xi_{ij}^k}(x^{k+1}), & \text{if } c_k = 1, \\ g^k + \mathcal{Q}\left(\frac{1}{b'} \sum_{j \in I'_{i,k}} (\nabla f_{\xi_{ij}^k}(x^{k+1}) - \nabla f_{\xi_{ij}^k}(x^k))\right), & \text{if } c_k = 0, \end{cases}$  where  $I_{i,k}, I'_{i,k}$  are the sets of the indices in the minibatches,  $|I_{i,k}| = b, |I'_{i,k}| = b'$ , and  $\xi_{ij}^k$  is independently sampled from  $\mathcal{D}_i$  for  $i \in [n], j \in [m]$
- 9:   **end for**
- 10:    $g^{k+1} = \frac{1}{n} \sum_{i=1}^n g_i^{k+1}$
- 11: **end for**
- 12: **Return:**  $\hat{x}^K$  chosen uniformly at random from  $\{x^k\}_{k=0}^{K-1}$

## D.2.1. GENERALLY NON-CONVEX PROBLEMS

In this section, we provide the full statement of Theorem 3.2 together with the proof of this result.

**Theorem D.3** (Theorem 3.2). *Consider the finite sum case (1)+(6). Let Assumptions 1.1, 1.2 and 3.2 be satisfied and*

$$\gamma \leq \frac{1}{L + \sqrt{\frac{1-p}{pn} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right)}}, \quad (40)$$

where  $L^2 = \frac{1}{n} \sum_{i=1}^n L_i^2$  and  $\mathcal{L}^2 = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i^2$ . Then after  $K$  iterations of VR-MARINA, we have

$$\mathbf{E} \left[ \|\nabla f(\hat{x}^K)\|^2 \right] \leq \frac{2\Delta_0}{\gamma K} + \frac{\sigma^2}{nb}, \quad (41)$$

where  $\hat{x}^K$  is chosen uniformly at random from  $x^0, \dots, x^{K-1}$  and  $\Delta_0 = f(x^0) - f_*$ . That is, after

$$K = \mathcal{O} \left( \frac{\Delta_0}{\varepsilon^2} \left( L + \sqrt{\frac{1-p}{pn} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right)} \right) \right) \quad (42)$$

iterations with  $b = \Theta\left(\frac{\sigma^2}{n\varepsilon^2}\right)$  VR-MARINA produces such a point  $\hat{x}^K$  that  $\mathbf{E}[\|\nabla f(\hat{x}^K)\|^2] \leq \varepsilon^2$ , and the expected total number of stochastic oracle calls per node equals

$$b + K(pb + 2(1-p)b') = \mathcal{O} \left( \frac{\sigma^2}{n\varepsilon^2} + \frac{\Delta_0}{\varepsilon^2} \left( L + \sqrt{\frac{1-p}{pn} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right)} \right) \left( p \frac{\sigma^2}{n\varepsilon^2} + (1-p)b' \right) \right). \quad (43)$$

Moreover, under an assumption that the communication cost is proportional to the number of non-zero components of transmitted vectors from workers to the server we have that the expected total communication cost per worker equals

$$d + K(pd + (1-p)\zeta_{\mathcal{Q}}) = \mathcal{O} \left( d + \frac{\Delta_0}{\varepsilon^2} \left( L + \sqrt{\frac{1-p}{pn} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right)} \right) (pd + (1-p)\zeta_{\mathcal{Q}}) \right), \quad (44)$$

where  $\zeta_{\mathcal{Q}}$  is the expected density of the quantization (see Def. 1.1).

*Proof of Theorem 3.2.* The proof follows the same steps as the proof of Theorem 3.1. From Lemma B.1, we have

$$\mathbf{E}[f(x^{k+1})] \leq \mathbf{E}[f(x^k)] - \frac{\gamma}{2} \mathbf{E}[\|\nabla f(x^k)\|^2] - \left( \frac{1}{2\gamma} - \frac{L}{2} \right) \mathbf{E}[\|x^{k+1} - x^k\|^2] + \frac{\gamma}{2} \mathbf{E}[\|g^k - \nabla f(x^k)\|^2]. \quad (45)$$

Next, we need to derive an upper bound for  $\mathbf{E} [\|g^{k+1} - \nabla f(x^{k+1})\|^2]$ . Since  $g^{k+1} = \frac{1}{n} \sum_{i=1}^n g_i^{k+1}$ , we get the following representation of  $g^{k+1}$ :

$$g^{k+1} = \begin{cases} \frac{1}{nb} \sum_{i=1}^n \sum_{j \in I_{i,k}} \nabla f_{\xi_{ij}^k}(x^{k+1}) & \text{with probability } p, \\ g^k + \frac{1}{n} \sum_{i=1}^n \mathcal{Q} \left( \frac{1}{b'} \sum_{j \in I'_{i,k}} (\nabla f_{\xi_{ij}^k}(x^{k+1}) - \nabla f_{\xi_{ij}^k}(x^k)) \right) & \text{with probability } 1 - p. \end{cases}$$

Using this, variance decomposition (13), tower property (14), and independence of  $\xi_{ij}^k$  for  $i \in [n]$ ,  $j \in I_{i,k}$ , we derive:

$$\begin{aligned} \mathbf{E} [\|g^{k+1} - \nabla f(x^{k+1})\|^2] &\stackrel{(14)}{=} (1-p) \mathbf{E} \left[ \left\| g^k + \frac{1}{n} \sum_{i=1}^n \mathcal{Q} \left( \frac{1}{b'} \sum_{j \in I'_{i,k}} (\nabla f_{ij}(x^{k+1}) - \nabla f_{ij}(x^k)) \right) - \nabla f(x^{k+1}) \right\|^2 \right] \\ &\quad + \frac{p}{n^2 b^2} \mathbf{E} \left[ \left\| \sum_{i=1}^n \sum_{j \in I_{i,k}} (\nabla f_{\xi_{ij}^k}(x^{k+1}) - \nabla f(x^{k+1})) \right\|^2 \right] \\ &\stackrel{(14),(13)}{=} (1-p) \mathbf{E} \left[ \left\| \frac{1}{n} \sum_{i=1}^n \mathcal{Q} \left( \frac{1}{b'} \sum_{j \in I'_{i,k}} (\nabla f_{ij}(x^{k+1}) - \nabla f_{ij}(x^k)) \right) - \nabla f(x^{k+1}) + \nabla f(x^k) \right\|^2 \right] \\ &\quad + (1-p) \mathbf{E} [\|g^k - \nabla f(x^k)\|^2] + \frac{p}{n^2 b^2} \sum_{i=1}^n \sum_{j \in I_{i,k}} \mathbf{E} [\|\nabla f_{\xi_{ij}^k}(x^{k+1}) - \nabla f(x^{k+1})\|^2] \\ &\stackrel{(14),(10)}{=} (1-p) \mathbf{E} \left[ \left\| \frac{1}{n} \sum_{i=1}^n \mathcal{Q} \left( \frac{1}{b'} \sum_{j \in I'_{i,k}} (\nabla f_{ij}(x^{k+1}) - \nabla f_{ij}(x^k)) \right) - \nabla f(x^{k+1}) + \nabla f(x^k) \right\|^2 \right] \\ &\quad + (1-p) \mathbf{E} [\|g^k - \nabla f(x^k)\|^2] + \frac{p\sigma^2}{nb}, \end{aligned}$$

where  $\sigma^2 = \frac{1}{n} \sum_{i=1}^n \sigma_i^2$ . Applying the same arguments as in the proof of inequality (33), we obtain

$$\begin{aligned} \mathbf{E} [\|g^{k+1} - \nabla f(x^{k+1})\|^2] &\leq \frac{1-p}{n} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right) \mathbf{E} [\|x^{k+1} - x^k\|^2] \\ &\quad + (1-p) \mathbf{E} [\|g^k - \nabla f(x^k)\|^2] + \frac{p\sigma^2}{nb}. \end{aligned} \quad (46)$$

Next, we introduce new notation:  $\Phi_k = f(x^k) - f_* + \frac{\gamma}{2p} \|g^k - \nabla f(x^k)\|^2$ . Using this and inequalities (45) and (46), we establish the following inequality:

$$\begin{aligned} \mathbf{E} [\Phi_{k+1}] &\leq \mathbf{E} \left[ f(x^k) - f_* - \frac{\gamma}{2} \|\nabla f(x^k)\|^2 - \left( \frac{1}{2\gamma} - \frac{L}{2} \right) \|x^{k+1} - x^k\|^2 + \frac{\gamma}{2} \|g^k - \nabla f(x^k)\|^2 \right] \\ &\quad + \frac{\gamma}{2p} \mathbf{E} \left[ \frac{1-p}{n} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right) \|x^{k+1} - x^k\|^2 + (1-p) \|g^k - \nabla f(x^k)\|^2 + \frac{p\sigma^2}{nb} \right] \\ &= \mathbf{E} [\Phi_k] - \frac{\gamma}{2} \mathbf{E} [\|\nabla f(x^k)\|^2] + \left( \frac{\gamma(1-p)}{2pn} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right) - \frac{1}{2\gamma} + \frac{L}{2} \right) \mathbf{E} [\|x^{k+1} - x^k\|^2] + \frac{\gamma\sigma^2}{2nb} \\ &\stackrel{(40)}{\leq} \mathbf{E} [\Phi_k] - \frac{\gamma}{2} \mathbf{E} [\|\nabla f(x^k)\|^2] + \frac{\gamma\sigma^2}{2nb}, \end{aligned} \quad (47)$$

where in the last inequality, we use  $\frac{\gamma(1-p)}{2pn} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right) - \frac{1}{2\gamma} + \frac{L}{2} \leq 0$  following from (40). Summing up inequalities (47) for  $k = 0, 1, \dots, K-1$  and rearranging the terms, we derive

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbf{E} [\|\nabla f(x^k)\|^2] \leq \frac{2}{\gamma K} \sum_{k=0}^{K-1} (\mathbf{E} [\Phi_k] - \mathbf{E} [\Phi_{k+1}]) + \frac{\sigma^2}{nb} = \frac{2(\mathbf{E} [\Phi_0] - \mathbf{E} [\Phi_K])}{\gamma K} + \frac{\sigma^2}{nb} = \frac{2\Delta_0}{\gamma K} + \frac{\sigma^2}{nb},$$

since  $g^0 = \nabla f(x^0)$  and  $\Phi_{k+1} \geq 0$ . Finally, using the tower property (14) and the definition of  $\hat{x}^K$ , we obtain (41) that implies (42), (43), and (44).  $\square$

**Remark D.2** (About batchsizes dissimilarity). *Similarly to the finite sum case, our analysis can be easily extended to handle the version of VR-MARINA with different batchsizes  $b_1, \dots, b_n$  and  $b'_1, \dots, b'_n$  on different workers, i.e., when  $|I_{i,k}| = b_i$ ,  $|I'_{i,k}| = b'_i$  for  $i \in [n]$ . In this case, the statement of Theorem 3.2 remains the same with the small modification: instead of  $\frac{\mathcal{L}^2}{b'}$  the complexity bounds will have  $\frac{1}{n} \sum_{i=1}^n \frac{\mathcal{L}_i^2}{b'_i}$ , and instead of the requirement  $b = \Theta\left(\frac{\sigma^2}{n\varepsilon}\right)$  it will have  $\frac{1}{n^2} \sum_{i=1}^n \frac{\sigma_i^2}{b_i} = \Theta(\varepsilon^2)$ .*

**Corollary D.3** (Corollary 3.2). *Let the assumptions of Theorem 3.2 hold and  $p = \min\left\{\frac{\zeta_{\mathcal{Q}}}{d}, \frac{b'}{b+b'}\right\}$ , where  $b' \leq b$ ,  $b = \Theta\left(\frac{\sigma^2}{(n\varepsilon^2)}\right)$  and  $\zeta_{\mathcal{Q}}$  is the expected density of the quantization (see Def. 1.1). If*

$$\gamma \leq \frac{1}{L + \sqrt{\frac{\max\{d/\zeta_{\mathcal{Q}}-1, b/b'\}}{n} \left(\omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'}\right)}},$$

then VR-MARINA requires

$$\mathcal{O}\left(\frac{\Delta_0}{\varepsilon^2} \left(L \left(1 + \sqrt{\frac{\omega}{n} \max\left\{\frac{d}{\zeta_{\mathcal{Q}}} - 1, \frac{\sigma^2}{nb'\varepsilon^2}\right\}}\right) + \mathcal{L} \sqrt{\frac{(1+\omega)}{nb'} \max\left\{\frac{d}{\zeta_{\mathcal{Q}}} - 1, \frac{\sigma^2}{nb'\varepsilon^2}\right\}}\right)\right)$$

iterations/communication rounds and

$$\mathcal{O}\left(\frac{\sigma^2}{n\varepsilon^2} + \frac{\Delta_0 L b'}{\varepsilon^2} + \frac{\Delta_0 L}{\varepsilon^2} \sqrt{\frac{\omega b'}{n} \max\left\{\left(\frac{d}{\zeta_{\mathcal{Q}}} - 1\right) b', \frac{\sigma^2}{n\varepsilon^2}\right\}} + \frac{\Delta_0 \mathcal{L}}{\varepsilon^2} \sqrt{\frac{1+\omega}{n} \max\left\{\left(\frac{d}{\zeta_{\mathcal{Q}}} - 1\right) b', \frac{\sigma^2}{n\varepsilon^2}\right\}}\right)$$

stochastic oracle calls per node in expectation to achieve  $\mathbf{E}[\|\nabla f(\hat{x}^K)\|^2] \leq \varepsilon^2$ , and the expected total communication cost per worker is

$$\mathcal{O}\left(d + \frac{\Delta_0 \zeta_{\mathcal{Q}}}{\varepsilon^2} \left(L \left(1 + \sqrt{\frac{\omega}{n} \max\left\{\frac{d}{\zeta_{\mathcal{Q}}} - 1, \frac{\sigma^2}{nb'\varepsilon^2}\right\}}\right) + \mathcal{L} \sqrt{\frac{1+\omega}{nb'} \max\left\{\frac{d}{\zeta_{\mathcal{Q}}} - 1, \frac{\sigma^2}{nb'\varepsilon^2}\right\}}\right)\right)$$

under an assumption that the communication cost is proportional to the number of non-zero components of transmitted vectors from workers to the server.

*Proof of Corollary 3.1.* The choice of  $p = \min\left\{\frac{\zeta_{\mathcal{Q}}}{d}, \frac{b'}{b+b'}\right\}$  implies

$$\begin{aligned} \frac{1-p}{p} &= \max\left\{\frac{d}{\zeta_{\mathcal{Q}}} - 1, \frac{b}{b'}\right\}, \\ pm + (1-p)b' &\leq \frac{2mb'}{m+b'} \leq 2b', \\ pd + (1-p)\zeta_{\mathcal{Q}} &\leq \frac{\zeta_{\mathcal{Q}}}{d} \cdot d + \left(1 - \frac{\zeta_{\mathcal{Q}}}{d}\right) \cdot \zeta_{\mathcal{Q}} \leq 2\zeta_{\mathcal{Q}}. \end{aligned}$$

Plugging these relations in (40), (42), (43) and (44) and using  $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$ , we get that if

$$\gamma \leq \frac{1}{L + \sqrt{\frac{\max\{d/\zeta_{\mathcal{Q}}-1, b/b'\}}{n} \left(\omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'}\right)}},$$

then VR-MARINA requires

$$\begin{aligned}
 K &= \mathcal{O} \left( \frac{\Delta_0}{\varepsilon^2} \left( L + \sqrt{\frac{1-p}{pn} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right)} \right) \right) \\
 &= \mathcal{O} \left( \frac{\Delta_0}{\varepsilon^2} \left( L + \sqrt{L^2 \frac{\omega \max\{d/\zeta_{\mathcal{Q}} - 1, b/b'\}}{n} + \mathcal{L}^2 \frac{(1+\omega) \max\{d/\zeta_{\mathcal{Q}} - 1, b/b'\}}{nb'}} \right) \right) \\
 &= \mathcal{O} \left( \frac{\Delta_0}{\varepsilon^2} \left( L \left( 1 + \sqrt{\frac{\omega}{n} \max\left\{ \frac{d}{\zeta_{\mathcal{Q}}} - 1, \frac{\sigma^2}{nb'\varepsilon^2} \right\}} \right) + \mathcal{L} \sqrt{\frac{(1+\omega)}{nb'} \max\left\{ \frac{d}{\zeta_{\mathcal{Q}}} - 1, \frac{\sigma^2}{nb'\varepsilon^2} \right\}} \right) \right)
 \end{aligned}$$

iterations/communication rounds and

$$\begin{aligned}
 b + K(pb + 2(1-p)b') &= \mathcal{O} \left( b + \frac{\Delta_0}{\varepsilon^2} \left( L + \sqrt{\frac{1-p}{pn} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right)} \right) (pb + (1-p)b') \right) \\
 &= \mathcal{O} \left( b + \frac{\Delta_0}{\varepsilon^2} \left( L \left( 1 + \sqrt{\frac{\omega \max\{d/\zeta_{\mathcal{Q}} - 1, b/b'\}}{n}} \right) \right. \right. \\
 &\quad \left. \left. + \mathcal{L} \sqrt{\frac{(1+\omega) \max\{d/\zeta_{\mathcal{Q}} - 1, b/b'\}}{nb'}} \right) b' \right) \\
 &= \mathcal{O} \left( \frac{\sigma^2}{n\varepsilon^2} + \frac{\Delta_0}{\varepsilon^2} \left( L \left( b' + \sqrt{\frac{\omega b'}{n} \max\left\{ \left( \frac{d}{\zeta_{\mathcal{Q}}} - 1 \right) b', \frac{\sigma^2}{n\varepsilon^2} \right\}} \right) \right. \right. \\
 &\quad \left. \left. + \mathcal{L} \sqrt{\frac{1+\omega}{n} \max\left\{ \left( \frac{d}{\zeta_{\mathcal{Q}}} - 1 \right) b', \frac{\sigma^2}{n\varepsilon^2} \right\}} \right) \right)
 \end{aligned}$$

stochastic oracle calls per node in expectation to achieve  $\mathbf{E}[\|\nabla f(\hat{x}^K)\|^2] \leq \varepsilon^2$ , and the expected total communication cost per worker is

$$\begin{aligned}
 d + K(pd + (1-p)\zeta_{\mathcal{Q}}) &= \mathcal{O} \left( d + \frac{\Delta_0}{\varepsilon^2} \left( L + \sqrt{\frac{1-p}{pn} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right)} \right) (pd + (1-p)\zeta_{\mathcal{Q}}) \right) \\
 &= \mathcal{O} \left( d + \frac{\Delta_0 \zeta_{\mathcal{Q}}}{\varepsilon^2} \left( L \left( 1 + \sqrt{\frac{\omega}{n} \max\left\{ \frac{d}{\zeta_{\mathcal{Q}}} - 1, \frac{\sigma^2}{nb'\varepsilon^2} \right\}} \right) \right. \right. \\
 &\quad \left. \left. + \mathcal{L} \sqrt{\frac{1+\omega}{nb'} \max\left\{ \frac{d}{\zeta_{\mathcal{Q}}} - 1, \frac{\sigma^2}{nb'\varepsilon^2} \right\}} \right) \right)
 \end{aligned}$$

under an assumption that the communication cost is proportional to the number of non-zero components of transmitted vectors from workers to the server.  $\square$

#### D.2.2. CONVERGENCE RESULTS UNDER POLYAK-ŁOJASIEWICZ CONDITION

In this section, we provide an analysis of VR-MARINA under Polyak-Łojasiewicz condition in the online case.

**Theorem D.4.** *Consider the finite sum case (1)+(6). Let Assumptions 1.1, 1.2, 3.2, 2.1 and 3.3 be satisfied and*

$$\gamma \leq \min \left\{ \frac{1}{L + \sqrt{\frac{2(1-p)}{pn} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right)}}, \frac{p}{2\mu} \right\}, \quad (48)$$

where  $L^2 = \frac{1}{n} \sum_{i=1}^n L_i^2$  and  $\mathcal{L}^2 = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i^2$ . Then after  $K$  iterations of VR-MARINA, we have

$$\mathbf{E} [f(x^K) - f(x^*)] \leq (1 - \gamma\mu)^K \Delta_0 + \frac{\sigma^2}{nb\mu}, \quad (49)$$

where  $\Delta_0 = f(x^0) - f(x^*)$ . That is, after

$$K = \mathcal{O} \left( \max \left\{ \frac{1}{p}, \frac{L + \sqrt{\frac{1-p}{pn} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right)}}{\mu} \right\} \log \frac{\Delta_0}{\varepsilon} \right) \quad (50)$$

iterations with  $b = \Theta \left( \frac{\sigma^2}{n\mu\varepsilon} \right)$  VR-MARINA produces such a point  $x^K$  that  $\mathbf{E} [f(x^K) - f(x^*)] \leq \varepsilon$ , and the expected total number of stochastic oracle calls per node equals

$$b + K(pb + 2(1-p)b') = \mathcal{O} \left( m + \max \left\{ \frac{1}{p}, \frac{L + \sqrt{\frac{1-p}{pn} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right)}}{\mu} \right\} (pb + (1-p)b') \log \frac{\Delta_0}{\varepsilon} \right). \quad (51)$$

Moreover, under an assumption that the communication cost is proportional to the number of non-zero components of transmitted vectors from workers to the server, we have that the expected total communication cost per worker equals

$$d + K(pd + (1-p)\zeta_{\mathcal{Q}}) = \mathcal{O} \left( d + \max \left\{ \frac{1}{p}, \frac{L + \sqrt{\frac{1-p}{pn} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right)}}{\mu} \right\} (pd + (1-p)\zeta_{\mathcal{Q}}) \log \frac{\Delta_0}{\varepsilon} \right), \quad (52)$$

where  $\zeta_{\mathcal{Q}}$  is the expected density of the quantization (see Def. 1.1).

*Proof.* The proof is very similar to the proof of Theorem 3.2. From Lemma B.1 and PL condition, we have

$$\begin{aligned} \mathbf{E}[f(x^{k+1}) - f(x^*)] &\leq \mathbf{E}[f(x^k) - f(x^*)] - \frac{\gamma}{2} \mathbf{E} [\|\nabla f(x^k)\|^2] - \left( \frac{1}{2\gamma} - \frac{L}{2} \right) \mathbf{E} [\|x^{k+1} - x^k\|^2] \\ &\quad + \frac{\gamma}{2} \mathbf{E} [\|g^k - \nabla f(x^k)\|^2] \\ &\stackrel{(4)}{\leq} (1 - \gamma\mu) \mathbf{E} [f(x^k) - f(x^*)] - \left( \frac{1}{2\gamma} - \frac{L}{2} \right) \mathbf{E} [\|x^{k+1} - x^k\|^2] + \frac{\gamma}{2} \mathbf{E} [\|g^k - \nabla f(x^k)\|^2]. \end{aligned}$$

Using the same arguments as in the proof of (46), we obtain

$$\begin{aligned} \mathbf{E} [\|g^{k+1} - \nabla f(x^{k+1})\|^2] &\leq \frac{1-p}{n} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right) \mathbf{E} [\|x^{k+1} - x^k\|^2] \\ &\quad + (1-p) \mathbf{E} [\|g^k - \nabla f(x^k)\|^2] + \frac{p\sigma^2}{nb}. \end{aligned} \quad (53)$$

Putting all together, we derive that the sequence  $\Phi_k = f(x^k) - f(x^*) + \frac{\gamma}{p} \|g^k - \nabla f(x^k)\|^2$  satisfies

$$\begin{aligned} \mathbf{E} [\Phi_{k+1}] &\leq \mathbf{E} \left[ (1 - \gamma\mu)(f(x^k) - f(x^*)) - \left( \frac{1}{2\gamma} - \frac{L}{2} \right) \|x^{k+1} - x^k\|^2 + \frac{\gamma}{2} \|g^k - \nabla f(x^k)\|^2 \right] \\ &\quad + \frac{\gamma}{p} \mathbf{E} \left[ \frac{1-p}{n} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right) \|x^{k+1} - x^k\|^2 + (1-p) \|g^k - \nabla f(x^k)\|^2 + \frac{p\sigma^2}{nb} \right] \\ &= \mathbf{E} \left[ (1 - \gamma\mu)(f(x^k) - f(x^*)) + \left( \frac{\gamma}{2} + \frac{\gamma}{p}(1-p) \right) \|g^k - \nabla f(x^k)\|^2 \right] + \frac{\gamma\sigma^2}{nb} \\ &\quad + \left( \frac{\gamma(1-p)}{pn} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right) - \frac{1}{2\gamma} + \frac{L}{2} \right) \mathbf{E} [\|x^{k+1} - x^k\|^2] \\ &\stackrel{(35)}{\leq} (1 - \gamma\mu) \mathbf{E} [\Phi_k] + \frac{\gamma\sigma^2}{nb}, \end{aligned}$$



where in the last inequality we use  $\frac{\gamma(1-p)}{pn} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right) - \frac{1}{2\gamma} + \frac{L}{2} \leq 0$  and  $\frac{\gamma}{2} + \frac{\gamma}{p}(1-p) \leq (1-\gamma\mu)\frac{\gamma}{p}$  following from (48). Unrolling the recurrence and using  $g^0 = \nabla f(x^0)$ , we obtain

$$\begin{aligned} \mathbf{E} [f(x^K) - f(x^*)] &\leq \mathbf{E}[\Phi_K] \leq (1-\gamma\mu)^K \Phi_0 + \frac{\gamma\sigma^2}{nb} \sum_{k=0}^{K-1} (1-\gamma\mu)^k \\ &\leq (1-\gamma\mu)^K (f(x^0) - f(x^*)) + \frac{\gamma\sigma^2}{nb} \sum_{k=0}^{\infty} (1-\gamma\mu)^k \\ &\leq (1-\gamma\mu)^K (f(x^0) - f(x^*)) + \frac{\sigma^2}{nb\mu}. \end{aligned}$$

Together with  $b = \Theta\left(\frac{\sigma^2}{n\mu\varepsilon}\right)$  it implies (50), (51), and (52).  $\square$

**Corollary D.4.** *Let the assumptions of Theorem D.4 hold and  $p = \min\left\{\frac{\zeta_{\mathcal{Q}}}{d}, \frac{b'}{b+b'}\right\}$ , where  $b' \leq b$  and  $\zeta_{\mathcal{Q}}$  is the expected density of the quantization (see Def. 1.1). If*

$$\gamma \leq \min \left\{ \frac{1}{L + \sqrt{\frac{2 \max\{d/\zeta_{\mathcal{Q}} - 1, b'/b'\}}{n} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right)}}, \frac{p}{2\mu} \right\}$$

and

$$b = \Theta\left(\frac{\sigma^2}{n\mu\varepsilon}\right), \quad \sigma^2 = \frac{1}{n} \sum_{i=1}^n \sigma_i^2,$$

then VR-MARINA requires

$$\mathcal{O} \left( \max \left\{ \frac{1}{p}, \frac{L}{\mu} \left( 1 + \sqrt{\frac{\omega}{n} \max \left\{ \frac{d}{\zeta_{\mathcal{Q}}} - 1, \frac{\sigma^2}{nb'\mu} \right\}} \right) + \frac{\mathcal{L}}{\mu} \sqrt{\frac{1+\omega}{nb'} \max \left\{ \frac{d}{\zeta_{\mathcal{Q}}} - 1, \frac{\sigma^2}{nb'\mu} \right\}} \right\} \log \frac{\Delta_0}{\varepsilon} \right)$$

iterations/communication rounds,

$$\begin{aligned} \mathcal{O} \left( \frac{\sigma^2}{n\mu\varepsilon} + \max \left\{ \frac{b'}{p}, \frac{L}{\mu} \left( b' + \sqrt{\frac{\omega b'}{n} \max \left\{ \left( \frac{d}{\zeta_{\mathcal{Q}}} - 1 \right) b', \frac{\sigma^2}{n\mu\varepsilon} \right\}} \right) \right. \right. \\ \left. \left. + \frac{\mathcal{L}}{\mu} \sqrt{\frac{1+\omega}{n} \max \left\{ \left( \frac{d}{\zeta_{\mathcal{Q}}} - 1 \right) b', \frac{\sigma^2}{n\mu\varepsilon} \right\}} \right\} \log \frac{\Delta_0}{\varepsilon} \right) \end{aligned}$$

stochastic oracle calls per node in expectation to achieve  $\mathbf{E}[f(x^K) - f(x^*)] \leq \varepsilon$ , and the expected total communication cost per worker is

$$\mathcal{O} \left( d + \zeta_{\mathcal{Q}} \max \left\{ \frac{1}{p}, \frac{L}{\mu} \left( 1 + \sqrt{\frac{\omega}{n} \max \left\{ \frac{d}{\zeta_{\mathcal{Q}}} - 1, \frac{\sigma^2}{nb'\mu} \right\}} \right) + \frac{\mathcal{L}}{\mu} \sqrt{\frac{1+\omega}{nb'} \max \left\{ \frac{d}{\zeta_{\mathcal{Q}}} - 1, \frac{\sigma^2}{nb'\mu} \right\}} \right\} \log \frac{\Delta_0}{\varepsilon} \right)$$

under an assumption that the communication cost is proportional to the number of non-zero components of transmitted vectors from workers to the server.

*Proof.* The choice of  $p = \min\left\{\frac{\zeta_{\mathcal{Q}}}{d}, \frac{b'}{b+b'}\right\}$  implies

$$\begin{aligned} \frac{1-p}{p} &= \max \left\{ \frac{d}{\zeta_{\mathcal{Q}}} - 1, \frac{b}{b'} \right\}, \\ pm + (1-p)b' &\leq \frac{2bb'}{b+b'} \leq 2b', \\ pd + (1-p)\zeta_{\mathcal{Q}} &\leq \frac{\zeta_{\mathcal{Q}}}{d} \cdot d + \left(1 - \frac{\zeta_{\mathcal{Q}}}{d}\right) \cdot \zeta_{\mathcal{Q}} \leq 2\zeta_{\mathcal{Q}}. \end{aligned}$$

Plugging these relations in (48), (50), (51) and (52) and using  $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$ , we get that if

$$\gamma \leq \min \left\{ \frac{1}{L + \sqrt{\frac{2 \max\{d/\zeta_Q - 1, b/b'\}}{n}} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right)}, \frac{p}{2\mu} \right\},$$

then VR-MARINA requires

$$\begin{aligned} K &= \mathcal{O} \left( \max \left\{ \frac{1}{p}, \frac{L + \sqrt{\frac{1-p}{pn}} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right)}{\mu} \right\} \log \frac{\Delta_0}{\varepsilon} \right) \\ &= \mathcal{O} \left( \max \left\{ \frac{1}{p}, \frac{L + \sqrt{L^2 \frac{\omega \max\{d/\zeta_Q - 1, b/b'\}}{n}} + \mathcal{L}^2 \frac{(1+\omega) \max\{d/\zeta_Q - 1, b/b'\}}{nb'}}{\mu} \right\} \log \frac{\Delta_0}{\varepsilon} \right) \\ &= \mathcal{O} \left( \max \left\{ \frac{1}{p}, \frac{L}{\mu} \left( 1 + \sqrt{\frac{\omega}{n}} \max \left\{ \frac{d}{\zeta_Q} - 1, \frac{\sigma^2}{nb'\mu} \right\} \right) + \frac{\mathcal{L}}{\mu} \sqrt{\frac{1+\omega}{nb'}} \max \left\{ \frac{d}{\zeta_Q} - 1, \frac{\sigma^2}{nb'\mu} \right\} \right\} \log \frac{\Delta_0}{\varepsilon} \right) \end{aligned}$$

iterations/communication rounds and

$$\begin{aligned} b + K(pb + 2(1-p)b') &= \mathcal{O} \left( b + \max \left\{ \frac{1}{p}, \frac{L + \sqrt{\frac{1-p}{pn}} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right)}{\mu} \right\} (pb + (1-p)b') \log \frac{\Delta_0}{\varepsilon} \right) \\ &= \mathcal{O} \left( b + \max \left\{ \frac{1}{p}, \frac{L}{\mu} \left( 1 + \sqrt{\frac{\omega \max\{d/\zeta_Q - 1, b/b'\}}{n}} \right) \right. \right. \\ &\quad \left. \left. + \frac{\mathcal{L}}{\mu} \sqrt{\frac{(1+\omega) \max\{d/\zeta_Q - 1, b/b'\}}{nb'}} \right\} b' \log \frac{\Delta_0}{\varepsilon} \right) \\ &= \mathcal{O} \left( \frac{\sigma^2}{n\mu\varepsilon} + \max \left\{ \frac{b'}{p}, \frac{L}{\mu} \left( b' + \sqrt{\frac{\omega b'}{n}} \max \left\{ \left( \frac{d}{\zeta_Q} - 1 \right) b', \frac{\sigma^2}{n\mu\varepsilon} \right\} \right) \right. \right. \\ &\quad \left. \left. + \frac{\mathcal{L}}{\mu} \sqrt{\frac{1+\omega}{n}} \max \left\{ \left( \frac{d}{\zeta_Q} - 1 \right) b', \frac{\sigma^2}{n\mu\varepsilon} \right\} \right\} \log \frac{\Delta_0}{\varepsilon} \right) \end{aligned}$$

stochastic oracle calls per node in expectation to achieve  $\mathbf{E}[f(x^K) - f(x^*)] \leq \varepsilon$ , and the expected total communication cost per worker is

$$\begin{aligned} d + K(pd + (1-p)\zeta_Q) &= \mathcal{O} \left( d + \max \left\{ \frac{1}{p}, \frac{L + \sqrt{\frac{1-p}{pn}} \left( \omega L^2 + \frac{(1+\omega)\mathcal{L}^2}{b'} \right)}{\mu} \right\} (pd + (1-p)\zeta_Q) \log \frac{\Delta_0}{\varepsilon} \right) \\ &= \mathcal{O} \left( d + \zeta_Q \max \left\{ \frac{1}{p}, \frac{L}{\mu} \left( 1 + \sqrt{\frac{\omega}{n}} \max \left\{ \frac{d}{\zeta_Q} - 1, \frac{\sigma^2}{nb'\mu} \right\} \right) \right. \right. \\ &\quad \left. \left. + \frac{\mathcal{L}}{\mu} \sqrt{\frac{1+\omega}{nb'}} \max \left\{ \frac{d}{\zeta_Q} - 1, \frac{\sigma^2}{nb'\mu} \right\} \right\} \log \frac{\Delta_0}{\varepsilon} \right) \end{aligned}$$

under an assumption that the communication cost is proportional to the number of non-zero components of transmitted vectors from workers to the server.  $\square$

## E. Missing Proofs for PP-MARINA

---

### Algorithm 4 PP-MARINA

---

- 1: **Input:** starting point  $x^0$ , stepsize  $\gamma$ , probability  $p \in (0, 1]$ , number of iterations  $K$ , clients-batchsize  $r \leq n$
  - 2: Initialize  $g^0 = \nabla f(x^0)$
  - 3: **for**  $k = 0, 1, \dots, K - 1$  **do**
  - 4:   Sample  $c_k \sim \text{Be}(p)$
  - 5:   Choose  $I'_k = \{1, \dots, n\}$  if  $c_k = 1$ , and choose  $I'_k$  as the set of  $r$  i.i.d. samples from the uniform distribution over  $\{1, \dots, n\}$  otherwise
  - 6:   Broadcast  $g^k$  to all workers
  - 7:   **for**  $i = 1, \dots, n$  in parallel **do**
  - 8:      $x^{k+1} = x^k - \gamma g^k$
  - 9:     Set  $g_i^{k+1} = \begin{cases} \nabla f_i(x^{k+1}) & \text{if } c_k = 1, \\ g^k + \mathcal{Q}(\nabla f_i(x^{k+1}) - \nabla f_i(x^k)) & \text{if } c_k = 0. \end{cases}$
  - 10:   **end for**
  - 11:   Set  $g^{k+1} = \begin{cases} \nabla f(x^{k+1}) & \text{if } c_k = 1, \\ g^k + \frac{1}{r} \sum_{i_k \in I'_k} \mathcal{Q}(\nabla f_{i_k}(x^{k+1}) - \nabla f_{i_k}(x^k)) & \text{if } c_k = 0. \end{cases}$
  - 12: **end for**
  - 13: **Return:**  $\hat{x}^K$  chosen uniformly at random from  $\{x^k\}_{k=0}^{K-1}$
- 

### E.1. Generally Non-Convex Problems

In this section, we provide the full statement of Theorem 4.1 together with the proof of this result.

**Theorem E.1** (Theorem 4.1). *Let Assumptions 1.1 and 1.2 be satisfied and*

$$\gamma \leq \frac{1}{L \left( 1 + \sqrt{\frac{(1-p)(1+\omega)}{pr}} \right)}, \quad (54)$$

where  $L^2 = \frac{1}{n} \sum_{i=1}^n L_i^2$ . Then after  $K$  iterations of PP-MARINA, we have

$$\mathbf{E} \left[ \|\nabla f(\hat{x}^K)\|^2 \right] \leq \frac{2\Delta_0}{\gamma K}, \quad (55)$$

where  $\hat{x}^K$  is chosen uniformly at random from  $x^0, \dots, x^{K-1}$  and  $\Delta_0 = f(x^0) - f_*$ . That is, after

$$K = \mathcal{O} \left( \frac{\Delta_0 L}{\varepsilon^2} \left( 1 + \sqrt{\frac{(1-p)(1+\omega)}{pr}} \right) \right) \quad (56)$$

iterations PP-MARINA produces such a point  $\hat{x}^K$  that  $\mathbf{E}[\|\nabla f(\hat{x}^K)\|^2] \leq \varepsilon^2$ . Moreover, under an assumption that the communication cost is proportional to the number of non-zero components of transmitted vectors from workers to the server, we have that the expected total communication cost (for all workers) equals

$$dn + K(pdn + (1-p)\zeta_{\mathcal{Q}}r) = \mathcal{O} \left( dn + \frac{\Delta_0 L}{\varepsilon^2} \left( 1 + \sqrt{\frac{(1-p)(1+\omega)}{pr}} \right) (pdn + (1-p)\zeta_{\mathcal{Q}}r) \right), \quad (57)$$

where  $\zeta_{\mathcal{Q}}$  is the expected density of the quantization (see Def. 1.1).

*Proof of Theorem 4.1.* The proof is very similar to the proof of Theorem 3.1. From Lemma B.1, we have

$$\mathbf{E}[f(x^{k+1})] \leq \mathbf{E}[f(x^k)] - \frac{\gamma}{2} \mathbf{E}[\|\nabla f(x^k)\|^2] - \left( \frac{1}{2\gamma} - \frac{L}{2} \right) \mathbf{E}[\|x^{k+1} - x^k\|^2] + \frac{\gamma}{2} \mathbf{E}[\|g^k - \nabla f(x^k)\|^2]. \quad (58)$$

Next, we need to derive an upper bound for  $\mathbf{E} [\|g^{k+1} - \nabla f(x^{k+1})\|^2]$ . By definition of  $g^{k+1}$ , we have

$$g^{k+1} = \begin{cases} \nabla f(x^{k+1}) & \text{with probability } p, \\ g^k + \frac{1}{r} \sum_{i_k \in I'_k} \mathcal{Q}(\nabla f_{i_k}(x^{k+1}) - \nabla f_{i_k}(x^k)) & \text{with probability } 1-p. \end{cases}$$

Using this, variance decomposition (13) and tower property (14), we derive:

$$\begin{aligned} \mathbf{E} [\|g^{k+1} - \nabla f(x^{k+1})\|^2] &\stackrel{(14)}{=} (1-p) \mathbf{E} \left[ \left\| g^k + \frac{1}{r} \sum_{i_k \in I'_k} \mathcal{Q}(\nabla f_{i_k}(x^{k+1}) - \nabla f_{i_k}(x^k)) - \nabla f(x^{k+1}) \right\|^2 \right] \\ &\stackrel{(14),(13)}{=} (1-p) \mathbf{E} \left[ \left\| \frac{1}{r} \sum_{i_k \in I'_k} \mathcal{Q}(\nabla f_{i_k}(x^{k+1}) - \nabla f_{i_k}(x^k)) - \nabla f(x^{k+1}) + \nabla f(x^k) \right\|^2 \right] \\ &\quad + (1-p) \mathbf{E} [\|g^k - \nabla f(x^k)\|^2]. \end{aligned}$$

Next, we use the notation:  $\Delta_{i_k}^k = \nabla f_{i_k}(x^{k+1}) - \nabla f_{i_k}(x^k)$  for  $i_k \in [n]$  and  $\Delta^k = \nabla f(x^{k+1}) - \nabla f(x^k)$ . These vectors satisfy  $\mathbf{E} [\Delta_{i_k}^k | x^k, x^{k+1}] = \Delta^k$  for all  $i_k \in I'_k$ . Moreover,  $\mathcal{Q}(\Delta_{i_k}^k)$  for  $i_k \in I'_k$  are independent random vectors for fixed  $x^k$  and  $x^{k+1}$ . These observations imply

$$\begin{aligned} \mathbf{E} [\|g^{k+1} - \nabla f(x^{k+1})\|^2] &= (1-p) \mathbf{E} \left[ \left\| \frac{1}{r} \sum_{i_k \in I'_k} (\mathcal{Q}(\Delta_{i_k}^k) - \Delta^k) \right\|^2 \right] + (1-p) \mathbf{E} [\|g^k - \nabla f(x^k)\|^2] \\ &= \frac{1-p}{r} \mathbf{E} [\|\mathcal{Q}(\Delta_{i_k}^k) - \Delta_{i_k}^k + \Delta_{i_k}^k - \Delta^k\|^2] + (1-p) \mathbf{E} [\|g^k - \nabla f(x^k)\|^2] \\ &\stackrel{(14),(13)}{=} \frac{1-p}{r} \left( \mathbf{E} [\|\mathcal{Q}(\Delta_{i_k}^k) - \Delta_{i_k}^k\|^2] + \mathbf{E} [\|\Delta_{i_k}^k - \Delta^k\|^2] \right) \\ &\quad + (1-p) \mathbf{E} [\|g^k - \nabla f(x^k)\|^2] \\ &\stackrel{(14),(3)}{=} \frac{1-p}{r} \left( \omega \mathbf{E} [\|\Delta_{i_k}^k\|^2] + \mathbf{E} [\|\Delta_{i_k}^k - \Delta^k\|^2] \right) + (1-p) \mathbf{E} [\|g^k - \nabla f(x^k)\|^2] \\ &\stackrel{(14),(13)}{=} \frac{(1-p)(1+\omega)}{r} \mathbf{E} [\|\Delta_{i_k}^k\|^2] + (1-p) \mathbf{E} [\|g^k - \nabla f(x^k)\|^2]. \end{aligned}$$

Using  $L$ -smoothness (2) of  $f_i$  together with the tower property (14), we get

$$\begin{aligned} \mathbf{E} [\|g^{k+1} - \nabla f(x^{k+1})\|^2] &\leq \frac{(1-p)(1+\omega)}{nr} \sum_{i=1}^n L_i^2 \mathbf{E} [\|x^{k+1} - x^k\|^2] + (1-p) \mathbf{E} [\|g^k - \nabla f(x^k)\|^2] \\ &= \frac{(1-p)(1+\omega)L^2}{r} \mathbf{E} [\|x^{k+1} - x^k\|^2] + (1-p) \mathbf{E} [\|g^k - \nabla f(x^k)\|^2]. \end{aligned} \quad (59)$$

Next, we introduce new notation:  $\Phi_k = f(x^k) - f_* + \frac{\gamma}{2p} \|g^k - \nabla f(x^k)\|^2$ . Using this and inequalities (58) and (59), we establish the following inequality:

$$\begin{aligned} \mathbf{E} [\Phi_{k+1}] &\leq \mathbf{E} \left[ f(x^k) - f_* - \frac{\gamma}{2} \|\nabla f(x^k)\|^2 - \left( \frac{1}{2\gamma} - \frac{L}{2} \right) \|x^{k+1} - x^k\|^2 + \frac{\gamma}{2} \|g^k - \nabla f(x^k)\|^2 \right] \\ &\quad + \frac{\gamma}{2p} \mathbf{E} \left[ \frac{(1-p)(1+\omega)L^2}{r} \|x^{k+1} - x^k\|^2 + (1-p) \|g^k - \nabla f(x^k)\|^2 \right] \\ &= \mathbf{E} [\Phi_k] - \frac{\gamma}{2} \mathbf{E} [\|\nabla f(x^k)\|^2] + \left( \frac{\gamma(1-p)(1+\omega)L^2}{2pn} - \frac{1}{2\gamma} + \frac{L}{2} \right) \mathbf{E} [\|x^{k+1} - x^k\|^2] \\ &\stackrel{(54)}{\leq} \mathbf{E} [\Phi_k] - \frac{\gamma}{2} \mathbf{E} [\|\nabla f(x^k)\|^2], \end{aligned} \quad (60)$$

where in the last inequality we use  $\frac{\gamma(1-p)(1+\omega)L^2}{2pn} - \frac{1}{2\gamma} + \frac{L}{2} \leq 0$  following from (54). Summing up inequalities (34) for  $k = 0, 1, \dots, K-1$  and rearranging the terms, we derive

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbf{E} [\|\nabla f(x^k)\|^2] \leq \frac{2}{\gamma K} \sum_{k=0}^{K-1} (\mathbf{E}[\Phi_k] - \mathbf{E}[\Phi_{k+1}]) = \frac{2(\mathbf{E}[\Phi_0] - \mathbf{E}[\Phi_K])}{\gamma K} = \frac{2\Delta_0}{\gamma K},$$

since  $g^0 = \nabla f(x^0)$  and  $\Phi_{k+1} \geq 0$ . Finally, using the tower property (14) and the definition of  $\hat{x}^K$ , we obtain (55) that implies (56) and (57).  $\square$

**Corollary E.1** (Corollary 4.1). *Let the assumptions of Theorem 4.1 hold and  $p = \frac{\zeta_{\mathcal{Q}} r}{dn}$ , where  $r \leq n$  and  $\zeta_{\mathcal{Q}}$  is the expected density of the quantization (see Def. 1.1). If*

$$\gamma \leq \frac{1}{L \left( 1 + \sqrt{\frac{1+\omega}{r} \left( \frac{dn}{\zeta_{\mathcal{Q}} r} - 1 \right)} \right)},$$

then PP-MARINA requires

$$K = \mathcal{O} \left( \frac{\Delta_0 L}{\varepsilon^2} \left( 1 + \sqrt{\frac{1+\omega}{r} \left( \frac{dn}{\zeta_{\mathcal{Q}} r} - 1 \right)} \right) \right)$$

iterations/communication rounds to achieve  $\mathbf{E}[\|\nabla f(\hat{x}^K)\|^2] \leq \varepsilon^2$ , and the expected total communication cost is

$$\mathcal{O} \left( dn + \frac{\Delta_0 L}{\varepsilon^2} \left( \zeta_{\mathcal{Q}} r + \sqrt{(1+\omega)\zeta_{\mathcal{Q}}(dn - \zeta_{\mathcal{Q}} r)} \right) \right)$$

under an assumption that the communication cost is proportional to the number of non-zero components of transmitted vectors from workers to the server.

*Proof of Corollary 4.1.* The choice of  $p = \frac{\zeta_{\mathcal{Q}} r}{dn}$  implies

$$\begin{aligned} \frac{1-p}{p} &= \frac{dn}{\zeta_{\mathcal{Q}} r} - 1, \\ pdn + (1-p)\zeta_{\mathcal{Q}} r &\leq \zeta_{\mathcal{Q}} r + \left( 1 - \frac{\zeta_{\mathcal{Q}} r}{dn} \right) \cdot \zeta_{\mathcal{Q}} r \leq 2\zeta_{\mathcal{Q}} r. \end{aligned}$$

Plugging these relations in (54), (56), and (57), we get that if

$$\gamma \leq \frac{1}{L \left( 1 + \sqrt{\frac{1+\omega}{r} \left( \frac{dn}{\zeta_{\mathcal{Q}} r} - 1 \right)} \right)},$$

then PP-MARINA requires

$$\begin{aligned} K &= \mathcal{O} \left( \frac{\Delta_0 L}{\varepsilon^2} \left( 1 + \sqrt{\frac{(1-p)(1+\omega)}{pr}} \right) \right) \\ &= \mathcal{O} \left( \frac{\Delta_0 L}{\varepsilon^2} \left( 1 + \sqrt{\frac{1+\omega}{r} \left( \frac{dn}{\zeta_{\mathcal{Q}} r} - 1 \right)} \right) \right) \end{aligned}$$

iterations/communication rounds in order to achieve  $\mathbf{E}[\|\nabla f(\hat{x}^K)\|^2] \leq \varepsilon^2$ , and the expected total communication cost is

$$\begin{aligned} dn + K(pd n + (1-p)\zeta_{\mathcal{Q}} r) &= \mathcal{O} \left( dn + \frac{\Delta_0 L}{\varepsilon^2} \left( 1 + \sqrt{\frac{(1-p)(1+\omega)}{pr}} \right) (pd n + (1-p)\zeta_{\mathcal{Q}} r) \right) \\ &= \mathcal{O} \left( dn + \frac{\Delta_0 L}{\varepsilon^2} \left( \zeta_{\mathcal{Q}} r + \sqrt{(1+\omega)\zeta_{\mathcal{Q}}(dn - \zeta_{\mathcal{Q}} r)} \right) \right) \end{aligned}$$

under an assumption that the communication cost is proportional to the number of non-zero components of transmitted vectors from workers to the server.  $\square$

## E.2. Convergence Results Under Polyak-Łojasiewicz condition

In this section, we provide an analysis of PP-MARINA under Polyak-Łojasiewicz condition.

**Theorem E.2.** *Let Assumptions 1.1, 1.2 and 2.1 be satisfied and*

$$\gamma \leq \min \left\{ \frac{1}{L \left( 1 + \sqrt{\frac{2(1-p)(1+\omega)}{pr}} \right)}, \frac{p}{2\mu} \right\}, \quad (61)$$

where  $L^2 = \frac{1}{n} \sum_{i=1}^n L_i^2$ . Then after  $K$  iterations of PP-MARINA, we have

$$\mathbf{E} [f(x^K) - f(x^*)] \leq (1 - \gamma\mu)^K \Delta_0, \quad (62)$$

where  $\Delta_0 = f(x^0) - f(x^*)$ . That is, after

$$K = \mathcal{O} \left( \max \left\{ \frac{1}{p}, \frac{L}{\mu} \left( 1 + \sqrt{\frac{(1-p)(1+\omega)}{pr}} \right) \right\} \log \frac{\Delta_0}{\varepsilon} \right) \quad (63)$$

iterations PP-MARINA produces such a point  $x^K$  that  $\mathbf{E}[f(x^K) - f(x^*)] \leq \varepsilon$ . Moreover, under an assumption that the communication cost is proportional to the number of non-zero components of transmitted vectors from workers to the server, we have that the expected total communication cost (for all workers) equals

$$dn + K(pdn + (1-p)\zeta_{\mathcal{Q}}r) = \mathcal{O} \left( dn + \max \left\{ \frac{1}{p}, \frac{L}{\mu} \left( 1 + \sqrt{\frac{(1-p)(1+\omega)}{pr}} \right) \right\} (pdn + (1-p)\zeta_{\mathcal{Q}}r) \log \frac{\Delta_0}{\varepsilon} \right), \quad (64)$$

where  $\zeta_{\mathcal{Q}}$  is the expected density of the quantization (see Def. 1.1).

*Proof.* The proof is very similar to the proof of Theorem 4.1. From Lemma B.1 and PŁ condition we have

$$\begin{aligned} \mathbf{E}[f(x^{k+1}) - f(x^*)] &\leq \mathbf{E}[f(x^k) - f(x^*)] - \frac{\gamma}{2} \mathbf{E} [\|\nabla f(x^k)\|^2] - \left( \frac{1}{2\gamma} - \frac{L}{2} \right) \mathbf{E} [\|x^{k+1} - x^k\|^2] \\ &\quad + \frac{\gamma}{2} \mathbf{E} [\|g^k - \nabla f(x^k)\|^2] \\ &\stackrel{(4)}{\leq} (1 - \gamma\mu) \mathbf{E} [f(x^k) - f(x^*)] - \left( \frac{1}{2\gamma} - \frac{L}{2} \right) \mathbf{E} [\|x^{k+1} - x^k\|^2] + \frac{\gamma}{2} \mathbf{E} [\|g^k - \nabla f(x^k)\|^2]. \end{aligned}$$

Using the same arguments as in the proof of (59), we obtain

$$\mathbf{E} [\|g^{k+1} - \nabla f(x^{k+1})\|^2] \leq \frac{(1-p)(1+\omega)L^2}{r} \mathbf{E} [\|x^{k+1} - x^k\|^2] + (1-p) \mathbf{E} [\|g^k - \nabla f(x^k)\|^2].$$

Putting all together, we derive that the sequence  $\Phi_k = f(x^k) - f(x^*) + \frac{\gamma}{p} \|g^k - \nabla f(x^k)\|^2$  satisfies

$$\begin{aligned} \mathbf{E} [\Phi_{k+1}] &\leq \mathbf{E} \left[ (1 - \gamma\mu)(f(x^k) - f(x^*)) - \left( \frac{1}{2\gamma} - \frac{L}{2} \right) \|x^{k+1} - x^k\|^2 + \frac{\gamma}{2} \|g^k - \nabla f(x^k)\|^2 \right] \\ &\quad + \frac{\gamma}{p} \mathbf{E} \left[ \frac{(1-p)(1+\omega)L^2}{r} \|x^{k+1} - x^k\|^2 + (1-p) \|g^k - \nabla f(x^k)\|^2 \right] \\ &= \mathbf{E} \left[ (1 - \gamma\mu)(f(x^k) - f(x^*)) + \left( \frac{\gamma}{2} + \frac{\gamma}{p}(1-p) \right) \|g^k - \nabla f(x^k)\|^2 \right] \\ &\quad + \left( \frac{\gamma(1-p)(1+\omega)L^2}{pr} - \frac{1}{2\gamma} + \frac{L}{2} \right) \mathbf{E} [\|x^{k+1} - x^k\|^2] \\ &\stackrel{(61)}{\leq} (1 - \gamma\mu) \mathbf{E} [\Phi_k], \end{aligned}$$



where in the last inequality we use  $\frac{\gamma(1-p)(1+\omega)L^2}{pr} - \frac{1}{2\gamma} + \frac{L}{2} \leq 0$  and  $\frac{\gamma}{2} + \frac{\gamma}{p}(1-p) \leq (1-\gamma\mu)\frac{\gamma}{p}$  following from (61). Unrolling the recurrence and using  $g^0 = \nabla f(x^0)$ , we obtain

$$\mathbf{E} [f(x^K) - f(x^*)] \leq \mathbf{E}[\Phi_K] \leq (1-\gamma\mu)^K \Phi_0 = (1-\gamma\mu)^K (f(x^0) - f(x^*))$$

that implies (63) and (64).  $\square$

**Corollary E.2.** *Let the assumptions of Theorem E.2 hold and  $p = \frac{\zeta_{\mathcal{Q}} r}{dn}$ , where  $r \leq n$  and  $\zeta_{\mathcal{Q}}$  is the expected density of the quantization (see Def. 1.1). If*

$$\gamma \leq \min \left\{ \frac{1}{L \left( 1 + \sqrt{\frac{2(1+\omega)}{r} \left( \frac{dn}{\zeta_{\mathcal{Q}} r} - 1 \right)} \right)}, \frac{p}{2\mu} \right\},$$

then PP-MARINA requires

$$K = \mathcal{O} \left( \max \left\{ \frac{dn}{\zeta_{\mathcal{Q}} r} \frac{L}{\mu} \left( 1 + \sqrt{\frac{1+\omega}{r} \left( \frac{dn}{\zeta_{\mathcal{Q}} r} - 1 \right)} \right) \right\} \log \frac{\Delta_0}{\varepsilon} \right)$$

iterations/communication rounds to achieve  $\mathbf{E}[f(x^K) - f(x^*)] \leq \varepsilon$ , and the expected total communication cost is

$$\mathcal{O} \left( dn + \max \left\{ dn, \frac{L}{\mu} \left( \zeta_{\mathcal{Q}} r + \sqrt{(1+\omega)\zeta_{\mathcal{Q}} (dn - \zeta_{\mathcal{Q}} r)} \right) \right\} \log \frac{\Delta_0}{\varepsilon} \right)$$

under an assumption that the communication cost is proportional to the number of non-zero components of transmitted vectors from workers to the server.

*Proof.* The choice of  $p = \frac{\zeta_{\mathcal{Q}} r}{dn}$  implies

$$\begin{aligned} \frac{1-p}{p} &= \frac{dn}{\zeta_{\mathcal{Q}} r} - 1, \\ pdn + (1-p)\zeta_{\mathcal{Q}} r &\leq \zeta_{\mathcal{Q}} r + \left( 1 - \frac{\zeta_{\mathcal{Q}} r}{dn} \right) \cdot \zeta_{\mathcal{Q}} r \leq 2\zeta_{\mathcal{Q}} r. \end{aligned}$$

Plugging these relations in (61), (63), and (64), we get that if

$$\gamma \leq \min \left\{ \frac{1}{L \left( 1 + \sqrt{\frac{2(1+\omega)}{r} \left( \frac{dn}{\zeta_{\mathcal{Q}} r} - 1 \right)} \right)}, \frac{p}{2\mu} \right\},$$

then PP-MARINA requires

$$\begin{aligned} K &= \mathcal{O} \left( \max \left\{ \frac{1}{p}, \frac{L}{\mu} \left( 1 + \sqrt{\frac{(1-p)(1+\omega)}{pr}} \right) \right\} \log \frac{\Delta_0}{\varepsilon} \right) \\ &= \mathcal{O} \left( \max \left\{ \frac{dn}{\zeta_{\mathcal{Q}} r} \frac{L}{\mu} \left( 1 + \sqrt{\frac{1+\omega}{r} \left( \frac{dn}{\zeta_{\mathcal{Q}} r} - 1 \right)} \right) \right\} \log \frac{\Delta_0}{\varepsilon} \right) \end{aligned}$$

iterations/communication rounds to achieve  $\mathbf{E}[f(x^K) - f(x^*)] \leq \varepsilon$ , and the expected total communication cost is

$$\begin{aligned} dn + K(pd n + (1-p)\zeta_{\mathcal{Q}} r) &= \mathcal{O} \left( dn + \max \left\{ \frac{1}{p}, \frac{L}{\mu} \left( 1 + \sqrt{\frac{(1-p)(1+\omega)}{pr}} \right) \right\} (pd n + (1-p)\zeta_{\mathcal{Q}} r) \log \frac{\Delta_0}{\varepsilon} \right) \\ &= \mathcal{O} \left( dn + \max \left\{ dn, \frac{L}{\mu} \left( \zeta_{\mathcal{Q}} r + \sqrt{(1+\omega)\zeta_{\mathcal{Q}} (dn - \zeta_{\mathcal{Q}} r)} \right) \right\} \log \frac{\Delta_0}{\varepsilon} \right) \end{aligned}$$

under an assumption that the communication cost is proportional to the number of non-zero components of transmitted vectors from workers to the server.  $\square$