MARINA-P: Superior Performance in Nonsmooth Federated Optimization with Adaptive Stepsizes

Anonymous authorsPaper under double-blind review

ABSTRACT

Non-smooth communication-efficient federated optimization remains largely unexplored theoretically, despite its importance in machine learning applications. We consider a setup focusing on optimizing downlink communication by improving state-of-the-art schemes like EF21-P (Gruntkowska et al., 2023) and MARINA-P (Gruntkowska et al., 2024) in the non-smooth convex setting. Our key contributions include extending the non-smooth convex theory of EF21-P from single-node to distributed settings and generalizing MARINA-P to non-smooth convex optimization. For both algorithms, we prove optimal $\mathcal{O}\left(1/\sqrt{T}\right)$ convergence rates under standard assumptions and establish matching communication complexity bounds with classical subgradient methods. We provide theoretical guarantees under constant, decreasing, and adaptive (Polyak-type) stepsizes. Our experiments demonstrate MARINA-P's superior performance with correlated compressors in both smooth non-convex and non-smooth convex settings. This work presents the first theoretical analysis of distributed non-smooth optimization with server-to-worker compression, including comprehensive analysis for various stepsize schemes.

1 Introduction

In recent years, the machine learning community has witnessed a paradigm shift toward larger models and datasets, spurring major performance gains but also posing new hardware, algorithmic, and software challenges (LeCun et al., 2015; Bottou et al., 2018; Kaplan et al., 2020; Deng et al., 2009).

The Rise of Big Data and Distributed Systems. The sheer volume of data needed for cutting-edge models has driven the adoption of distributed systems (Dean et al., 2012; Khirirat et al., 2018; Lin et al., 2018), since single-machine setups can no longer handle the storage and computational demands. This approach is particularly relevant in supervised learning (Hastie et al., 2009; Shalev-Shwartz & Ben-David, 2014; Vapnik, 2013), often formulated as:

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

where n denotes the number of clients, $x \in \mathbb{R}^d$ is the model's parameter vector, and $f_i(x)$ is the local loss on client i. Throughout, we assume each f_i is convex (possibly non-smooth).

Federated learning (FL) (McMahan et al., 2016; Konečný et al., 2016b;a; McMahan et al., 2017) extends the distributed paradigm to heterogeneous clients with decentralized data, seeking to avoid central data aggregation and preserve privacy. In FL, devices connect to a central server that orchestrates training (Konečný et al., 2016b; Kairouz et al., 2021): each device locally updates parameters using its data, then sends these updates to the server. The server aggregates them, performs global calculations, and broadcasts new parameters back to devices. This process continues until convergence or acceptable performance is reached.

Communication Challenges in Large-scale Model Training. Although distributing data alleviates storage and compute constraints, it introduces substantial communication overhead. Modern gradient-based methods (Bottou, 2012; Kingma & Ba, 2014; Demidovich et al., 2023; Duchi et al., 2011; Robbins & Monro, 1951) require iterative updates for all *d* parameters, making frequent transmission of high-dimensional gradients expensive. Two broad approaches reduce this burden: (i) performing

056

057

058

060

061 062

063

064

065

066

067

068

069

070 071

072

073

074 075 076

077 078

079

080

081

082 083

084

085

087

880

089

090 091

092

093 094

096

097

098

103

104

105

106

107

multiple local gradient steps before communicating, as in LocalSGD (Stich, 2020; Khaled et al., 2020; Woodworth et al., 2020; Yi et al., 2024; Sadiev et al., 2022; Richtárik et al., 2024), and (ii) compressing gradients via lossy transformations (Khirirat et al., 2018; Alistarh et al., 2018b; Mishchenko et al., 2020; 2019; Li et al., 2020; Li & Richtárik, 2021; Richtárik et al., 2021; Fatkhullin et al., 2021; Richtárik et al., 2022; Seide et al., 2014; Alistarh et al., 2017; Panferov et al., 2024). Moreover, studies of 4G LTE and 5G networks (Huang et al., 2012; Narayanan et al., 2021) show that upload/download speeds are often comparable, emphasizing that both server-to-worker and worker-to-server communication must be optimized.

Prevalence of Non-smooth Objectives in Machine Learning Applications. Despite notable advances in distributed optimization, theoretical work has primarily targeted smooth objectives, leaving non-smooth problems less explored in federated contexts. Non-smoothness arises in various ML scenarios: ReLU activations (Glorot et al., 2011; Nair & Hinton, 2010), L1 regularization for sparsity (Tibshirani, 1996; Zou & Hastie, 2005), hinge loss (Cortes, 1995), total variation (Rudin et al., 1992; Chambolle, 2004), quantile regression (Koenker & Bassett Jr, 1978), max-pooling (Scherer et al., 2010), submodular minimization (Bach, 2013), Huber loss (Huber, 1964), and graph-based learning (Hallac et al., 2015).

Adaptive Stepsizes are Widely Used in Practice. Because constants like L-Lipschitz continuity or smoothness parameters are difficult to determine in deep learning, practitioners rely on adaptive learning rates. Popular methods include AdaGrad (Duchi et al., 2011), RMSProp, Adam (Kingma & Ba, 2014), and AMSGrad (Reddi et al., 2018), all of which adjust per-parameter stepsizes based on observed gradients.

1.1 NOTATION AND ASSUMPTIONS

We denote the set $\{1, 2, \dots, n\}$ by [n]. For vectors, $\|\cdot\|_2$ represents the Euclidean norm, while for matrices, it denotes the spectral norm. The inner product of vectors u and v is denoted by $\langle u, v \rangle$. We use $\mathcal{O}(\cdot)$ to hide absolute constants. We denote $R_0 := ||x^0 - x^*||_2$.

Our analysis relies on the following standard assumptions:

Assumption 1. The function f has at least one minimizer, denoted by x^* .

Assumption 2. The functions f_i are convex for all $i \in [n]$.

In the distributed setting, assuming convexity for individual functions f_i is sufficient, as it implies convexity for f itself.

Assumption 3 (Lipschitz continuity of f_i). Functions f_i are $L_{0,i}$ -Lipschitz continuous for all $i \in [n]$. That is, for all $i \in [n]$, there exists $L_{0,i} > 0$ such that $|f_i(x) - f_i(y)| \le L_{0,i} ||x - y||_2$, $\forall x, y \in \mathbb{R}^d$.

If each f_i is Lipschitz continuous, then by Jensen's inequality, f is L_0 -Lipschitz with $L_0:=$ $\frac{1}{n} \sum_{i=1}^{n} L_{0,i}$ (Nesterov, 2013).

Both convexity and Lipschitz continuity of f are standard assumptions in non-smooth optimization (Vorontsova et al., 2021; Nesterov, 2013; Bubeck, 2015; Beck, 2017; Duchi, 2018; Lan, 2020; Drusvyatskiy, 2020). Moreover, L_0 and $L_{0,i}$ -Lipschitz continuity imply uniformly bounded subgradients (Beck, 2017), a property that will be useful in our proofs:

$$\|\partial f(x)\|_2 \le L_0 \quad \forall x \in \mathbb{R}^d, \tag{2}$$

$$\|\partial f_i(x)\|_2 \le L_{0,i} \quad \forall x \in \mathbb{R}^d \text{ and } \forall i \in [n].$$
 (3)

 $\|\partial f_i(x)\|_2 \leq L_{0,i} \quad \forall x \in \mathbb{R}^d \text{ and } \forall i \in [n]. \tag{3}$ We define $\widetilde{L}_0 := \sqrt{\frac{1}{n} \sum_{i=1}^n L_{0,i}^2}$ and $\overline{L}_0 := \frac{1}{n} \sum_{i=1}^n L_{0,i}$. By the arithmetic-quadratic mean inequality, we have $\overline{L}_0 < \widetilde{L}_0$.

Following classical optimization literature (Nemirovski et al., 2009; Beck, 2017; Duchi, 2018; Lan, 2020; Drusvyatskiy, 2020), for non-smooth convex objectives, we aim to find an ε -suboptimal solution: a random vector $\hat{x} \in \mathbb{R}^d$ satisfying $\mathbb{E}[f(\hat{x}) - f(x^*)] \leq \varepsilon$, where $\mathbb{E}[\cdot]$ denotes the expectation over algorithmic randomness.

To assess the efficiency of distributed subgradient-based algorithms, we primarily use two metrics:

- 1. Communication complexity (alternatively, communication cost): The expected total number of floats per worker required to communicate to reach an ε -suboptimal solution. In this paper, we focus on server-to-worker communication compression.
- 2. Iteration complexity: The number of communication rounds needed to achieve an ε-suboptimal solution.

1.2 Related work

108

110

111

112

113 114

115 116

117

118

119

120

121

122

123

124

125

126

127

128 129

130

131

132 133

134 135

136

137

138

139

140

141

142

143 144

145

146

147

148

149

150

151 152

153 154

155 156

157

158 159

160

161

Subgradient Methods in Non-smooth Convex Optimization. Subgradient methods, pioneered in the 1960s (Shor et al., 1985; Polyak, 1987), remain central to non-smooth convex optimization. Classic theory establishes $\mathcal{O}(1/\sqrt{T})$ rates for general convex objectives (Nesterov, 2013; Vorontsova et al., 2021; Bubeck, 2015; Beck, 2017; Duchi, 2018; Lan, 2020; Drusvyatskiy, 2020) and $\mathcal{O}(1/T)$ for strongly convex problems (Beck, 2017; Drusvyatskiy, 2020). For unknown T, decreasing stepsizes of order $\mathcal{O}(1/\sqrt{t})$ or $\mathcal{O}(1/t)$ add a logarithmic factor, yielding $\mathcal{O}(\log T/\sqrt{T})$ (Nesterov, 2013) and $\mathcal{O}(\log T/T)$ (Hazan et al., 2007; Hazan & Kale, 2014). Nevertheless, recent works (Zhu et al., 2024; Lacoste-Julien et al., 2012; Rakhlin et al., 2011) have removed these factors, attaining optimal rates in convex and strongly convex settings. In the stochastic regime, mirror-descent methods also achieve $\mathcal{O}(1/\sqrt{T})$ (Nemirovski et al., 2009). Beyond averaged-iterate convergence, tighter last-iterate analyses (Jain et al., 2019; Zamani & Glineur, 2023) provide stronger guarantees. Subgradient methods remain crucial for large-scale machine learning tasks, including support vector machines and structured prediction (Shalev-Shwartz et al., 2007; Ratliff et al., 2007).

Communication Compression. Before discussing more advanced optimization methods, let us consider the simplest baseline: the standard subgradient method (SM) 1, which iteratively performs updates 2

$$x^{t+1} = x^t - \frac{\gamma_t}{n} \sum_{i=1}^n g_i^t, \tag{4}$$

 $x^{t+1} = x^t - \frac{\gamma_t}{n} \sum_{i=1}^n g_i^t,$ where $g_i^t = \partial f_i(x^t)$ is a subgradient of f_i at x^t . In the distributed setting, the method can be implemented as follows: each worker calculates g_i^t and sends it to the server, where the subgradients are aggregated. The server takes the step and broadcasts x^{t+1} back to the workers. With stepsize $\gamma_t := \frac{R_0}{L_0\sqrt{T}}$, where $R_0 := \|x^0 - x^*\|_2$ and T is the total number of iterations, SM finds an ε -approximate solution after $\mathcal{O}\left(L_0^2 R_0^2/\varepsilon^2\right)$ steps (Nesterov, 2013; Drusvyatskiy, 2020). Since at each step the workers and the server send $\Theta(d)$ coordinates/floats, the worker-to-server and serverto-worker communication costs are $\mathcal{O}\left(dL_0^2R_0^2/\varepsilon^2\right)$. To formally quantify communication costs, we introduce the following definition.

Definition 1. The worker-to-server (w2s, uplink) and server-to-worker (s2w, downlink) communication complexities of a method are the expected number of coordinates/floats that a worker sends to the server and that the server sends to a worker, respectively, to find an ε -solution.

Communication compression techniques, such as sparsification (Wang et al., 2018; Mishchenko et al., 2020; Alistarh et al., 2018b; Wangni et al., 2018; Konečný & Richtárik, 2018) and quantization (Alistarh et al., 2017; Wen et al., 2017; Zhang et al., 2016; Horváth et al., 2022; Wu et al., 2018; Mishchenko et al., 2019), are known to be immensely powerful for reducing the communication overhead of gradient-type methods. Existing literature primarily considers two main classes of compression operators: unbiased and biased (contractive) compressors.

Definition 2. (Unbiased compressor). A stochastic mapping $Q: \mathbb{R}^d \to \mathbb{R}^d$ is called an unbiased compressor/compression operator if there exists $\omega \geq 0$ such that for any $x \in \mathbb{R}^d$: $\mathbb{E}[\mathcal{Q}(x)] = x, \quad \mathbb{E}\left[\|\mathcal{Q}(x) - x\|_2^2\right] \leq \omega \|x\|_2^2.$

$$\mathbb{E}[\mathcal{Q}(x)] = x, \quad \mathbb{E}\left[\|\mathcal{Q}(x) - x\|_2^2\right] \le \omega \|x\|_2^2. \tag{5}$$

This definition encompasses a wide range of well-known compression techniques, including RandK sparsification (Stich et al., 2018), random dithering (Roberts, 1962; Goodall, 1951), and natural

¹In this paper, we use the non-normalized form (4) of the subgradient method studied in (Vorontsova et al., 2021; Bubeck, 2015; Beck, 2017; Duchi, 2018; Lan, 2020; Drusvyatskiy, 2020; Nemirovski et al., 2009). Earlier works (Shor et al., 1985; Polyak, 1987) typically employed SM in the form $x^{t+1} = x^t - \frac{\gamma_t}{\|\partial f(x^t)\|} \partial f(x^t)$, which includes an additional normalization term $\|\partial f(x^t)\|$.

²For constrained optimization problems, the subgradient method typically operates through projections onto a convex set X (see (Bubeck, 2015; Lacoste-Julien et al., 2012; Beck, 2017; Duchi, 2018)). However, when optimizing over an unbounded domain, i.e., $\mathcal{X} = \mathbb{R}^d$, projections are not needed.

compression (Horváth et al., 2022). Notable examples of methods employing compressor (5) are QSGD (Alistarh et al., 2017), DCGD (Khirirat et al., 2018), MARINA (Gorbunov et al., 2021), DIANA (Mishchenko et al., 2019), VR-DIANA (Horváth et al., 2019), DASHA (Tyurin & Richtárik, 2023), FedCOMGATE (Haddadpour et al., 2021), FedPAQ (Reisizadeh et al., 2020), FedSTEPH (Das et al., 2020), FedCOM (Haddadpour et al., 2021), ADIANA (Li et al., 2020), NEOLITHIC (Huang et al., 2022a), ACGD (Li et al., 2020), and CANITA (Li & Richtárik, 2021). However, Definition 2 does not cover another important class of practically more favorable compressors, called *contractive*, which are usually biased.

Definition 3. (Contractive compressor). A stochastic mapping $C : \mathbb{R}^d \to \mathbb{R}^d$ is called a contractive compressor/compression operator if there exists $\alpha \in (0,1]$ such that for any $x \in \mathbb{R}^d$:

$$\mathbb{E}\left[\|\mathcal{C}(x) - x\|_{2}^{2}\right] \le (1 - \alpha) \|x\|_{2}^{2}. \tag{6}$$

We denote the families of compressors satisfying Definitions 2 and 3 by $\mathbb{U}(\omega)$ and $\mathbb{B}(\alpha)$, respectively.³

Inequality (6) is satisfied by many compressors, including Top K (Ström, 2015; Dryden et al., 2016; Aji & Heafield, 2017; Alistarh et al., 2018b), quantization (Alistarh et al., 2017; Horváth et al., 2022), low-rank approximations (Vogels et al., 2019; 2020; Safaryan et al., 2021), and count-sketches (Ivkin et al., 2019; Rothchild et al., 2020). For broader surveys, see (Beznosikov et al., 2023; Demidovich et al., 2023; Safaryan et al., 2022). However, naive distributed SGD with biased compression (e.g., Top K) can diverge (Beznosikov et al., 2023). Error Feedback (EF14), introduced by Seide et al. (2014), emerged as a key technique to avert such divergence. Early theory of EF14 was confined to single-node settings (Stich et al., 2018; Alistarh et al., 2018a; Stich & Karimireddy, 2019), then expanded to distributed setting (Cordonnier, 2018; Beznosikov et al., 2023; Koloskova et al., 2020). Richtárik et al. (2021) reformulated EF14 into EF21, achieving optimal $\mathcal{O}\left(^1/T\right)$ convergence for smooth non-convex problems, surpassing the previous $\mathcal{O}\left(^1/T^{2/3}\right)$ rate (Koloskova et al., 2020).

The EF21 framework led to multiple variants (Richtárik et al., 2022; Fatkhullin et al., 2021), including bidirectional (s2w and w2s) biased compression. Gruntkowska et al. (2023) introduced EF21-P, combining biased s2w and unbiased w2s to improve complexity in smooth strongly convex settings. Later, Gruntkowska et al. (2024) proposed MARINA-P for smooth non-convex problems, delivering sharper bounds than both EF21 and EF21-P. Concurrently, Anonymous (2024) provided the first non-smooth convergence guarantees for EF21-P, albeit restricted to single-node scenarios.

In order to express communication complexities, we will further need the following quantities.

Definition 4 (Expected density). For the given compression operators $\mathcal{Q}(x)$ and $\mathcal{C}(x)$, we define the expected density as $\zeta_{\mathcal{Q}} = \sup_{x \in \mathbb{R}^d} \mathbb{E}\left[\|\mathcal{Q}(x)\|_0\right]$ and $\zeta_{\mathcal{C}} = \sup_{x \in \mathbb{R}^d} \mathbb{E}\left[\|\mathcal{C}(x)\|_0\right]$, 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 compression operator since $\|Q(x)\|_0 \le d$ and $\|C(x)\|_0 \le d$.

The landscape of communication-efficient federated methods for non-smooth optimization is largely unexplored, with most research targeting smooth objectives or single-node settings. Below, we highlight open challenges and gaps.

Numerous works study s2w compression (Zheng et al., 2019; Gruntkowska et al., 2023; Fatkhullin et al., 2021; Philippenko & Dieuleveut, 2021; Liu et al., 2020; Gorbunov et al., 2020; Safaryan et al., 2022; Huang et al., 2022b; Horváth et al., 2022; Tang et al., 2019; Tyurin & Richtarik, 2023; Gruntkowska et al., 2024), yet almost all focus on smooth objectives. To our knowledge, only Karimireddy et al. (2019) and Anonymous (2024) offer non-smooth convex guarantees with s2w compression, and both are limited to single-node settings with minimal relevance to federated learning.

Distributed subgradient methods are well-studied, but either lack compression (Nedic & Ozdaglar, 2009; Kiwiel & Lindberg, 2001; Hishinuma & Iiduka, 2015; Zheng et al., 2022) or focus on specific operators like quantization (Xia et al., 2023; Doan et al., 2020; 2018; Xia et al., 2022; Emiola & Enyioha, 2022), covering only the w2s direction. No comprehensive treatments address s2w compression in non-smooth distributed optimization.

³Notably, it can easily be verified (see Lemma 8 in (Richtárik et al., 2021)) that if $Q \in \mathbb{U}(\omega)$, then $(\omega+1)^{-1}Q \in \mathbb{B}((\omega+1)^{-1})$, indicating that the family of biased compressors is wider.

Method	Non-smooth	Distributed	Compressed communications	Compression type	Adaptive stepsizes
EF14 (Karimireddy et al., 2019)	✓	X	✓	w2s	×
EF21-P (Anonymous, 2024)	✓	Х	✓	s2w	✓
MARINA-P (Gruntkowska et al., 2024)	×	✓	✓	s2w	×
SM with Polyak Stepsize (Hazan & Kakade, 2019)	✓	Х	X	-	✓
SM with Quantization (Xia et al., 2023)	✓	✓	✓	w2s	×
EF21-P [OURS]	1	✓	✓	s2w	✓
MARINA-P [OURS]	1	✓	√	s2w	/

Table 1: Summary of optimization methods employing worker-to-server (w2s) or server-to-worker (s2w) compression schemes.

Recent works on adaptive stepsizes (Khaled et al., 2023; Defazio et al., 2023; 2024; Mishchenko & Defazio; Defazio & Mishchenko, 2023) primarily handle single-node problems. Polyak stepsizes (Polyak, 1987; Hazan & Kakade, 2019) remain popular, but current studies (Loizou et al., 2021; Oikonomou & Loizou, 2024; Jiang & Stich, 2024) often assume smoothness or lack distributed analysis. Even existing non-smooth convex results (Hazan & Kakade, 2019; Schaipp et al., 2023) remain restricted to single-node contexts.

In summary, the intersection of non-smooth optimization, communication efficiency, and federated learning remains underexplored. Our work addresses this gap by providing the first comprehensive study of distributed non-smooth optimization with s2w compression and adaptive stepsizes, while maintaining optimal convergence rates.

2 Contributions

We now summarize our main contributions: • Extension of EF21-P to distributed non-smooth settings. We generalize EF21-P (Anonymous, 2024) from single-node to distributed architectures, proving optimal $\mathcal{O}(1/\sqrt{T})$ rates for both Polyak and constant stepsizes, and a suboptimal $\mathcal{O}(\log T/\sqrt{T})$ rate for decreasing stepsizes. Our communication complexity matches classical distributed subgradient methods, addressing a longstanding gap in Error Feedback theory for non-smooth problems.

- Introduction of MARINA-P for non-smooth objectives. Building on Gruntkowska et al. (2024), we extend MARINA-P from smooth non-convex to non-smooth convex settings, establishing optimal $\mathcal{O}(^{1}/\sqrt{T})$ rates for constant and Polyak stepsizes, along with a suboptimal $\mathcal{O}(^{\log T}/\sqrt{T})$ rate under decreasing steps.
- Superior performance of MARINA-P with correlated compressors. Empirical results show that MARINA-P, when paired with correlated compressors, surpasses EF21-P in non-smooth settings. This extends the known benefits of correlated compressors, previously shown for smooth non-convex objectives, to non-smooth convex federated tasks.
- Support for diverse stepsize schedules. We provide theoretical guarantees for both algorithms under constant, decreasing, and Polyak stepsizes, bridging the gap between foundational theory and practical deep learning use cases.

To our knowledge, these are the first theoretical results for distributed non-smooth optimization incorporating s2w compression and adaptive stepsizes, while achieving provably optimal convergence rates.

Algorithm 1 EF21-P (distributed version)

```
271
              1: Input: initial points w^0 = x^0 \in \mathbb{R}^d, stepsize \gamma_0 > 0
272
             2: for t = 0, 1, 2, \dots, T do
273
             3:
                       for i = 1, \dots, n on Workers in parallel do
274
             4:
                            Receive compressed difference \Delta^t from server
275
                             Compute local subgradient g_i^t = \partial f_i(w^t) and send it to server
             5:
276
             6:
277
             7:
                       On Server:
278
             8:
                       Receive g_i^t from workers
                       Choose stepsize \gamma_t (can be set according to (9), (10), or (11)) x^{t+1} = x^t - \gamma_t \frac{1}{n} \sum_{i=1}^n g_i^t
Compute \Delta^{t+1} = \mathcal{C}(x^{t+1} - w^t) and broadcast it to workers
279
             9:
            10:
281
            11:
                       w^{t+1} = w^t + \Delta^{t+1}
            12:
                       for i=1,\dots,n on Workers in parallel do w^{t+1}=w^t+\Delta^{t+1}
            13:
283
            14:
284
            15:
                       end for
285
            16: end for
286
            17: Output: x^T
287
```

3 EF21-P

270

289

290291292

293

294

295

296

297

298

299

300

301

302

303

304

305 306

307

308

309 310

311

312 313

314 315

316

317 318

319320321322323

We now present the first major contribution of our paper: a distributed version of EF21-P for the non-smooth setting.

Let us first recap the standard single-node EF21-P algorithm, which aims to solve (1) via the iterative process:

$$x^{t+1} = x^{t} - \gamma_{t} \nabla f(w^{t})$$

$$w^{t+1} = w^{t} + C^{t} (x^{t+1} - w^{t}),$$
(7)

where $\gamma_t>0$ is a stepsize, $x^0\in\mathbb{R}^d$ is the initial iterate, $w^0=x^0\in\mathbb{R}^d$ is the initial iterate shift, and \mathcal{C}^t is an instantiation of a randomized contractive compressor \mathcal{C} sampled at time t. This method was proposed as a primal counterpart to the standard EF21. It has proven particularly useful in bidirectional settings where primal compression is performed on the server side, allowing for the decoupling of primal and dual compression parameter constants. For more details, we refer the reader to the original paper (Gruntkowska et al., 2023). Anonymous (2024) first extended EF21-P to the non-smooth setting. Their key modification was replacing the "smooth" update step with a "non-smooth" one: $x^{t+1}=x^t-\gamma_t\partial f(w^t)$.

They proved an optimal rates of $\mathcal{O}(1/\sqrt{T})$ for Polyak and constant stepsizes, and a suboptimal rate of $\mathcal{O}(\log T/\sqrt{T})$ for decreasing stepsizes, but only for the single-node regime. In Algorithm 1, we extend these results to the distributed setting, allowing for parallel computation of subgradients $\partial f(w^t)$.

At each iteration of distributed EF21-P, the workers calculate $\partial f_i(w^t)$ and transmit it to the server. The server then averages the subgradients and updates the global model x^t . Subsequently, it computes the compressed difference $\Delta^{t+1} = \mathcal{C}_i^t(x^{t+1} - w^t)$ and broadcasts the same vector Δ^{t+1} to all workers. Both the server and workers then use Δ^{t+1} to update their internal states w^t . Note that this procedure ensures that the states w^t remain synchronized between workers and the server.

We now present the convergence result for our distributed EF21-P algorithm.

Theorem 1. Let Assumptions 1, 2 and 3 hold. Define a Lyapunov function $V^t := \|x^t - x^*\|_2^2 + \frac{1}{\lambda_* \theta} \|w^t - x^t\|_2^2$, where $\lambda_* := \frac{\sqrt{1-\alpha}}{1-\sqrt{1-\alpha}}$ and $\theta := 1 - \sqrt{1-\alpha}$. Define also a constant $B_* := 1 + 2\frac{\sqrt{1-\alpha}}{1-\sqrt{1-\alpha}}$. Let $\{w^t\}$ be the sequence produced by EF21-P (Algorithm 1). Define $\overline{w}^T := \frac{1}{T} \sum_{t=0}^{T-1} w^t$ and $\widehat{w}^T := \frac{1}{\sum_{t=0}^{T-1} \gamma_t} \sum_{t=0}^{T-1} \gamma_t w^t$.

⁴Since it operates in the primal space of model parameters

1. (Constant stepsize). If $\gamma_t := \gamma > 0$, then

$$\mathbb{E}\left[f(\overline{w}^T) - f(x^*)\right] \le \frac{V^0}{2\gamma T} + \frac{B_* L_0^2 \gamma}{2}.$$
 (8)

If, moreover, optimal γ is chosen i.e.

324

325

326

327

328 329

330 331

332

333 334

335 336

337

338

339 340 341

342

348

349

350

351

352

353

354

355

356

357

358

359 360

361

362

363

364 365 366

367

368

369

370

371

372

373 374

375

376

377

$$\gamma := \frac{1}{\sqrt{T}} \sqrt{\frac{V^0}{B_* L_0^2}},\tag{9}$$

then $\mathbb{E}\left[f(\overline{w}^T) - f(x^*)\right] \leq \frac{\sqrt{B_*L_0^2V^0}}{\sqrt{T}}.$

2. Polyak stepsize. If γ_t is chosen as

$$\gamma_t := \frac{f(w^t) - f(x^*)}{B_* \|\partial f(w^t)\|_2^2}, \tag{10}$$

then $\mathbb{E}\left[f(\overline{w}^T) - f(x^*)\right] \leq \frac{\sqrt{B_*L_0^2V^0}}{\sqrt{T}}$.

3. (Decreasing stepsize). If γ_t is chosen as

$$\gamma_t := \frac{\gamma_0}{\sqrt{t+1}},\tag{11}$$

$$\gamma_t := \frac{\gamma_0}{\sqrt{t+1}},$$
 then $\mathbb{E}\left[f(\widehat{w}^T) - f(x^*)\right] \leq \frac{V^0 + 2\gamma_0^2 B_* L_0^2 \log(T+1)}{\gamma_0 \sqrt{T}}.$ If, moreover, optimal γ_0 is chosen i.e.

If, moreover, optimal γ_0 *is chosen i.e.*

If, moreover, optimal
$$\gamma_0$$
 is chosen i.e.
$$\gamma_0 := \sqrt{\frac{V_0}{2B_* L_0^2 \log(T+1)}},$$
then $\mathbb{E}\left[f(\widehat{w}^T) - f(x^*)\right] \le 2\sqrt{2B_* L_0^2 V_0} \sqrt{\frac{\log(T+1)}{T}}.$

Let us analyze the obtained results. The constant $B_*:=1+2\frac{\sqrt{1-\alpha}}{1-\sqrt{1-\alpha}}\leq \frac{4}{\alpha}-1$ is a decreasing function in α , which aligns with intuition since larger values of α correspond to less aggressive compression regimes. For both constant (9) and Polyak (10) stepsizes, we achieve the optimal rate of $\mathcal{O}(1/\sqrt{T})$ known for uncompressed subgradient methods (Nesterov, 2013; Arjevani & Shamir, 2015). However, achieving this rate requires either knowing the total number of iterations T in advance (for constant stepsize) or knowing the optimal value $f(x^*)$ (for Polyak stepsize), which may be impractical in many applications. When neither T nor $f(x^*)$ is known, one can employ the decreasing stepsize strategy (11). This approach leads to a suboptimal convergence rate of $\mathcal{O}(\log T/\sqrt{T})$ – a well-known limitation of subgradient methods (Nesterov, 2013; Anonymous, 2024).

For both constant and Polyak stepsizes, the following corollary provides explicit complexity bounds, characterizing both the number of iterations and the total communication cost needed to obtain an ε -approximate solution.

Corollary 1. Let the conditions of the Theorem 1 are met. If γ_t is set according to (9) or (10) (constant or Polyak stepsizes) then EF21-P (Algorithm 1) requires $T = \mathcal{O}\left(\frac{L_0^2 R_0^2}{\alpha \varepsilon^2}\right)$ iterations/communication rounds in order to achieve $\mathbb{E}\left[f(\overline{w}^T) - f(x^*)\right] \leq \varepsilon$, and the expected total communication cost per worker is $\mathcal{O}(d + \zeta_{\mathcal{C}}T)$.

Let us analyze the implications of Corollary 1. In the uncompressed case ($\alpha = 1$), our algorithm achieves the optimal rate of standard Subgradient Methods (SM) (Nesterov, 2013) for first-order nonsmooth optimization. With Top K compression ($\zeta_C = K$), the communication complexity becomes $\mathcal{O}\left(dL_0^2R_0^2/\varepsilon^2\right)$, matching the worst-case complexity of distributed SM. This indicates that EF21-P with Top K compression cannot improve upon SM's complexity regardless of the compression parameter α – a fundamental limitation in communication-compressed non-smooth optimization. Our findings align with Balkanski & Singer (2018), who demonstrated that parallelization provides no worst-case benefits for non-smooth optimization.

From a practical perspective, EF21-P's main limitation stems from broadcasting identical compressed differences Δ_t to all workers, potentially leading to poor approximations of x^{t+1} by $w^t + \Delta_t$. The MARINA-P algorithm (Gruntkowska et al., 2024), originally developed for smooth non-convex problems, addresses this limitation. In the following section, we extend MARINA-P to the non-smooth setting.

Algorithm 2 MARINA-P

378

401 402

403 404

405

406

407

408

409

410

411 412

413

414

415 416

417

418

423 424

425

430

431

```
379
             1: Input: initial point x^0 \in \mathbb{R}^d, initial model shifts w_i^0 = x^0 for all i \in [n], stepsize \gamma_0 > 0,
380
                 probability 0 , compressors <math>\mathcal{Q}_i^t \in \mathbb{U}(\omega) for all i \in [n]
             2: for t = 0, 1, ..., T do
382
             3:
                      for i = 1, \dots, n on Workers in parallel do
             4:
                           Compute local subgradient g_i^t = \partial f_i(w_i^t) and send it to server
384
             5:
                      end for
             6:
                      On Server:
386
             7:
                      Receive g_i^t from workers
                      Choose stepsize \gamma_t (can be set according to (13), (14), or (15)) x^{t+1} = x^t_i - \gamma_t \frac{1}{n} \sum_{i=1}^n g_i^t
387
             8:
             9:
                      Sample c^t \sim \text{Bernoulli}(p)
389
            10:
                      if c^t = 0 then
390
            11:
                           Send Q_i^t(x^{t+1} - x^t) to worker i for i \in [n]
            12:
391
            13:
                      else
392
                           Send x^{t+1} to all workers
            14:
393
            15:
394
            16:
                      for i = 1, ..., n on Workers in parallel do
395
                           w_i^{t+1} = \begin{cases} x^{t+1} & \text{if } c^t = 1\\ w_i^t + \mathcal{Q}_i^t (x^{t+1} - x^t) & \text{if } c^t = 0 \end{cases}
            17:
397
                      end for
            18:
398
            19: end for
399
            20: Output: x^T
400
```

MARINA-P

Building upon the foundations of the standard MARINA algorithm (Gorbunov et al., 2021; Szlendak et al., 2022), Gruntkowska et al. (2024) introduced MARINA-P, a primal counterpart designed to operate in the model parameter space. This section presents an extension of MARINA-P to the non-smooth convex setting. At each iteration, workers compute subgradients $\partial f_i(w_i^t)$ and transmit them to the server. The server aggregates these subgradients and updates the global model x^t . With probability p (typically small), the server sends the uncompressed updated model x^{t+1} to all workers. Otherwise, each worker i receives a compressed vector $\mathcal{Q}_i^t(x^{t+1}-x^t)$. Workers then update their local models w_i^{t+1} based on the received information. A key feature of MARINA-P is that the compressed vectors $\mathcal{Q}_1^t(x^{t+1}-x^t),\dots,\mathcal{Q}_n^t(x^{t+1}-x^t)$ can differ across workers. This distinction is crucial for the algorithm's practical superiority, as it allows for potentially better approximations of x^{t+1} compared to methods like EF21-P, especially when the compressors $\mathcal{Q}_1, \dots, \mathcal{Q}_n$ are correlated.

We now present the main convergence results for MARINA-P in the non-smooth convex setting.

Theorem 2. Let Assumptions 1, 2 and 3 hold. Define a Lyapunov function $V^t := \|x^t - x^*\|_2^2 +$ $\frac{1}{\lambda_* p} \frac{1}{n} \sum_{i=1}^n \|w_i^t - x^t\|_2^2, \text{ where } \lambda_* := \frac{\overline{L}_0}{\widetilde{L}_0} \sqrt{\frac{(1-p)\omega}{p}}. \text{ Define also a constant } \widetilde{B}_* := \overline{L}_0^2 + \overline{L}_0^2$ $2\overline{L}_0\widetilde{L}_0\sqrt{\frac{(1-p)\omega}{p}}$. Let $\{w_i^t\}$ be the sequence produced by MARINA-P (Algorithm 2). Define $\overline{w}_i^T:=rac{1}{T}\sum_{t=0}^{T-1}w_i^t$ and $\widehat{w}_i^T:=rac{1}{\sum_{t=0}^{T-1}\gamma_t}\sum_{t=0}^{T-1}\gamma_tw_i^t$ for all $i\in[n]$.

1. (Constant stepsize). If $\gamma_t := \gamma > 0$, then $\mathbb{E}\left[\frac{1}{n}\sum_{i=1}^n f_i(\overline{w}_i^T) - f(x^*)\right] \leq \frac{V^0}{2\sqrt{T}} + \frac{\widetilde{B}_*\gamma}{2}$.

If, moreover, optimal γ is chosen i.e.

If, moreover, optimal
$$\gamma$$
 is chosen i.e.
$$\gamma := \frac{1}{\sqrt{T}} \sqrt{\frac{V^0}{\tilde{B}_*}},$$
 then $\mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n f_i(\overline{w}_i^T) - f(x^*)\right] \leq \frac{\sqrt{\tilde{B}_* V^0}}{\sqrt{T}}.$
2. Polyak stepsize. If γ_t is chosen as

2. Polyak stepsize. If γ_t is chosen as

$$\gamma_{t} := \frac{\frac{1}{n} \sum_{i=1}^{n} f_{i}(w_{i}^{t}) - f(x^{*})}{\left\|\frac{1}{n} \sum_{i=1}^{n} \partial f_{i}(w_{i}^{t})\right\|_{2}^{2} \left(1 + 2 \frac{\sqrt{\frac{1}{n} \sum_{i=1}^{n} \left\|\partial f_{i}(w_{i}^{t})\right\|_{2}^{2}}{\left\|\frac{1}{n} \sum_{i=1}^{n} \partial f_{i}(w_{i}^{t})\right\|_{2}^{2}} \sqrt{\frac{(1-p)\omega}{p}}\right)},$$
(14)

then
$$\mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}f_{i}(\overline{w}_{i}^{T})-f(x^{*})\right] \leq \frac{\sqrt{\widetilde{B}_{\star}V^{0}}}{\sqrt{T}}$$
.

3. (Decreasing stepsize). If γ_t is chosen as

$$\gamma_t := \frac{\gamma_0}{\sqrt{t+1}},$$

$$then \mathbb{E}\left[\frac{1}{n}\sum_{i=1}^n f_i(\widehat{w}_i^T) - f(x^*)\right] \le \frac{V^0 + 2\gamma_0^2 \widetilde{B}_* \log(T+1)}{\gamma_0 \sqrt{T}}.$$
(15)

If, moreover, optimal γ_0 *is chosen i.e.*

$$\gamma_0 := \sqrt{\frac{V_0}{2\widetilde{B}_* \log(T+1)}},$$

$$then \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n f_i(\widehat{w}_i^T) - f(x^*)\right] \le 2\sqrt{2\widetilde{B}_* V_0} \sqrt{\frac{\log(T+1)}{T}}.$$

$$(16)$$

Remark 1. For both EF21-P and MARINA-P, the Polyak stepsize can be efficiently implemented in the distributed setting without additional per-iteration communication overhead. This is because the subgradient values $\partial f_i(w^t)$ (for EF21-P) and $\partial f_i(w^t_i)$ (for MARINA-P) are already computed by the clients and transmitted to the server as part of the algorithm's regular operations.

Let us analyze these results. The constant $\widetilde{B}_* := \overline{L}_0^2 + 2\overline{L}_0\widetilde{L}_0\sqrt{\frac{(1-p)\omega}{p}}$ depends on both compression parameters ω and p. Smaller values of ω correspond to less aggressive compression, while larger values of p indicate more frequent uncompressed communication – both cases lead to smaller \widetilde{B}_* and consequently faster convergence. For both constant (13) and Polyak (14) stepsizes, we obtain the optimal rate of $\mathcal{O}\left(1/\sqrt{T}\right)$ (Nesterov, 2013; Arjevani & Shamir, 2015). As with EF21-P, achieving this rate requires either knowing the total iterations T (for constant stepsize) or the optimal value $f(x^*)$ (for Polyak stepsize) in advance. When such knowledge is unavailable, the decreasing stepsize strategy offers a practical alternative, though it results in a suboptimal $\mathcal{O}\left(\log T/\sqrt{T}\right)$ convergence rate – a characteristic limitation of subgradient methods (Nesterov, 2013). It is worth noting that implementing the Polyak stepsize only requires an estimate of $f(x^*)$, rather than knowledge of the Lipschitz constant L_0 . This characteristic is common among Polyak stepsizes (Loizou et al., 2021).

For the constant and Polyak stepsize regimes, the following corollary establishes complexity bounds and characterizes the communication costs required to achieve an ε -approximate solution.

Corollary 2. Let the conditions of the Theorem 2 are met and $p = \zeta Q/d$. If γ_t is set according to (13) or (14) (constant or Polyak stepsizes) then MARINA-P (Algorithm 2) requires $T = \mathcal{O}\left(\frac{R_0^2}{\varepsilon^2}\left(\overline{L}_0^2 + \overline{L}_0\widetilde{L}_0\sqrt{\omega\left(\frac{d}{\zeta_Q} - 1\right)}\right)\right)$ iterations/communication rounds in order to achieve $\mathbb{E}\left[\frac{1}{n}\sum_{i=1}^n f_i(\overline{w}_i^T) - f(x^*)\right] \leq \varepsilon$, and the expected total communication cost per worker is $\mathcal{O}\left(d + \zeta_Q T\right)$.

This corollary reveals several important properties. With Rand K compression ($\zeta_{\mathcal{Q}} = K$, $\omega = d/K - 1$) (Beznosikov et al., 2023), MARINA-P achieves communication complexity $\mathcal{O}\left(d\widetilde{L}_0^2R_0^2/\varepsilon^2\right)$. Under the condition $\widetilde{L}_0^2 = \mathcal{O}\left(L_0\right)$, this matches the optimal per-worker complexity of standard SM, up to constant factors (Nesterov, 2013). A notable feature of our complexity result is its independence from the number of workers n in the non-smooth setting – a known phenomenon in subgradient methods (Arjevani & Shamir, 2015; Balkanski & Singer, 2018). This contrasts with MARINA-P's behavior in smooth non-convex settings (Gruntkowska et al., 2024), where complexity scales as $\mathcal{O}(1/n)$. The absence of theoretical bounds predicting such scaling behavior in non-smooth distributed settings presents an interesting direction for future research.

MARINA-P's primary advantage over EF21-P lies in its ability to employ worker-specific compression operators Q_i , enabling more accurate approximations of the global model, particularly when using correlated compressors. The following section examines various constructions of Q_i that leverage this flexibility to enhance practical performance.

5 IMPACT STATEMENT

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

REFERENCES

- Alham Fikri Aji and Kenneth Heafield. Sparse communication for distributed gradient descent. *arXiv* preprint arXiv:1704.05021, 2017.
- Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. QSGD: Communication-efficient SGD via gradient quantization and encoding. *Advances in neural information processing systems*, 30, 2017.
- Dan Alistarh, Torsten Hoefler, Mikael Johansson, Sarit Khirirat, Nikola Konstantinov, and Cédric Renggli. The convergence of sparsified gradient methods. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018a.
 - Dan Alistarh, Torsten Hoefler, Mikael Johansson, Nikola Konstantinov, Sarit Khirirat, and Cédric Renggli. The convergence of sparsified gradient methods. *Advances in Neural Information Processing Systems*, 31, 2018b.
 - Anonymous. Error feedback for smooth and nonsmooth convex optimization with constant, decreasing and polyak stepsizes. In *Submitted to The Thirteenth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=Qv9TG9yDG0. under review.
 - Yossi Arjevani and Ohad Shamir. Communication complexity of distributed convex learning and optimization. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015. URL https://proceedings.neurips.cc/paper_files/paper/2015/file/7fec306dle665bc9c748b5d2b99a6e97-[]Paper.pdf.
 - Francis Bach. Learning with submodular functions: A convex optimization perspective. *Foundations and Trends*® *in machine learning*, 6(2-3):145–373, 2013.
 - Eric Balkanski and Yaron Singer. Parallelization does not accelerate convex optimization: Adaptivity lower bounds for non-smooth convex minimization. *arXiv preprint arXiv:1808.03880*, 2018.
 - Amir Beck. First-order methods in optimization. SIAM, 2017.
 - Aleksandr Beznosikov, Samuel Horváth, Peter Richtárik, and Mher Safaryan. On biased compression for distributed learning. *Journal of Machine Learning Research*, 24(276):1–50, 2023.
 - Leon Bottou. Stochastic Gradient Descent Tricks, volume 7700 of Lecture Notes in Computer Science (LNCS), pp. 430–445. Springer, neural networks, tricks of the trade, reloaded edition, January 2012. URL https://www.microsoft.com/en-[]us/research/publication/stochastic-[]gradient-[]tricks/.
 - Léon Bottou, Frank Curtis, and Jorge Nocedal. Optimization methods for large-scale machine learning. *SIAM Review*, 60:223–311, 2018.
 - Sébastien Bubeck. Convex optimization: Algorithms and complexity. *Foundations and Trends*® *in Machine Learning*, 8(3-4):231–357, 2015.
 - Antonin Chambolle. An algorithm for total variation minimization and applications. *Journal of Mathematical imaging and vision*, 20:89–97, 2004.
 - Jean-Baptiste Cordonnier. Convex optimization using sparsified stochastic gradient descent with memory. 2018.
 - Corinna Cortes. Support-vector networks. *Machine Learning*, 1995.
- Rudrajit Das, Abolfazl Hashemi, Sujay Sanghavi, and Inderjit S Dhillon. Improved convergence rates for non-convex federated learning with compression. *arXiv preprint arXiv:2012.04061*, 2020.
- Jeffrey Dean, Greg Corrado, Rajat Monga, Kai Chen, Matthieu Devin, Mark Mao, Andrew Senior, Paul Tucker, Ke Yang, Quoc V Le, and et al. Large scale distributed deep networks. In *Advances in Neural Information Processing Systems*, pp. 1223–1231, 2012.

- Aaron Defazio and Konstantin Mishchenko. Learning-rate-free learning by d-adaptation. In *International Conference on Machine Learning*, pp. 7449–7479. PMLR, 2023.
 - Aaron Defazio, Ashok Cutkosky, Harsh Mehta, and Konstantin Mishchenko. When, why and how much? adaptive learning rate scheduling by refinement. *arXiv preprint arXiv:2310.07831*, 2023.
 - Aaron Defazio, Xingyu Alice Yang, Harsh Mehta, Konstantin Mishchenko, Ahmed Khaled, and Ashok Cutkosky. The road less scheduled. *arXiv preprint arXiv:2405.15682*, 2024.
 - Yury Demidovich, Grigory Malinovsky, Igor Sokolov, and Peter Richtárik. A guide through the zoo of biased SGD. *Advances in Neural Information Processing Systems*, 36:23158–23171, 2023.
 - Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
 - Thinh T Doan, Siva Theja Maguluri, and Justin Romberg. On the convergence of distributed subgradient methods under quantization. In 2018 56th Annual Allerton Conference on Communication, Control, and Computing (Allerton), pp. 567–574. IEEE, 2018.
 - Thinh T Doan, Siva Theja Maguluri, and Justin Romberg. Fast convergence rates of distributed subgradient methods with adaptive quantization. *IEEE Transactions on Automatic Control*, 66(5): 2191–2205, 2020.
 - Dmitriy Drusvyatskiy. Convex analysis and nonsmooth optimization. University Lecture, 2020.
 - Nikoli Dryden, Tim Moon, Sam Ade Jacobs, and Brian Van Essen. Communication quantization for data-parallel training of deep neural networks. In 2016 2nd Workshop on Machine Learning in HPC Environments (MLHPC), pp. 1–8, 2016. doi: 10.1109/MLHPC.2016.004.
 - John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(7), 2011.
 - John C Duchi. Introductory lectures on stochastic optimization. *The mathematics of data*, 25:99–186, 2018.
 - Iyanuoluwa Emiola and Chinwendu Enyioha. Quantized and distributed subgradient optimization method with malicious attack. *IEEE Control Systems Letters*, 7:181–186, 2022.
 - Ilyas Fatkhullin, Igor Sokolov, Eduard Gorbunov, Zhize Li, and Peter Richtárik. EF21 with bells & whistles: Practical algorithmic extensions of modern error feedback. *arXiv preprint arXiv:2110.03294*, 2021.
 - Xavier Glorot, Antoine Bordes, and Yoshua Bengio. Deep sparse rectifier neural networks. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pp. 315–323. JMLR Workshop and Conference Proceedings, 2011.
 - WM Goodall. Television by pulse code modulation. *Bell System Technical Journal*, 30(1):33–49, 1951.
 - Eduard Gorbunov, Dmitry Kovalev, Dmitry Makarenko, and Peter Richtárik. Linearly converging error compensated SGD. *Advances in Neural Information Processing Systems*, 33:20889–20900, 2020.
 - Eduard Gorbunov, Konstantin P Burlachenko, Zhize Li, and Peter Richtárik. MARINA: Faster non-convex distributed learning with compression. In *International Conference on Machine Learning*, pp. 3788–3798. PMLR, 2021.
 - Kaja Gruntkowska, Alexander Tyurin, and Peter Richtárik. EF21-P and friends: Improved theoretical communication complexity for distributed optimization with bidirectional compression. In *International Conference on Machine Learning*, pp. 11761–11807. PMLR, 2023.
 - Kaja Gruntkowska, Alexander Tyurin, and Peter Richtárik. Improving the worst-case bidirectional communication complexity for nonconvex distributed optimization under function similarity. *arXiv* preprint arXiv:2402.06412, 2024.

- Farzin Haddadpour, Mohammad Mahdi Kamani, Aryan Mokhtari, and Mehrdad Mahdavi. Federated learning with compression: Unified analysis and sharp guarantees. In *International Conference on Artificial Intelligence and Statistics*, pp. 2350–2358. PMLR, 2021.
 - David Hallac, Jure Leskovec, and Stephen Boyd. Network lasso: Clustering and optimization in large graphs. In *Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining*, pp. 387–396, 2015.
 - Trevor Hastie, Robert Tibshirani, Jerome H Friedman, and Jerome H Friedman. *The elements of statistical learning: data mining, inference, and prediction*, volume 2. Springer, 2009.
 - Elad Hazan and Sham Kakade. Revisiting the Polyak step size. *arXiv preprint arXiv:1905.00313*, 2019.
 - Elad Hazan and Satyen Kale. Beyond the regret minimization barrier: optimal algorithms for stochastic strongly-convex optimization. *The Journal of Machine Learning Research*, 15(1): 2489–2512, 2014.
 - Elad Hazan, Amit Agarwal, and Satyen Kale. Logarithmic regret algorithms for online convex optimization. *Machine Learning*, 69(2):169–192, 2007.
 - Kazuhiro Hishinuma and Hideaki Iiduka. Parallel subgradient method for nonsmooth convex optimization with a simple constraint. *Linear Nonlinear Anal*, 1:67–77, 2015.
 - Samuel Horváth, Dmitry Kovalev, Konstantin Mishchenko, Sebastian Stich, and Peter Richtárik. Stochastic distributed learning with gradient quantization and variance reduction. *arXiv* preprint *arXiv*:1904.05115, 2019.
 - Samuel Horváth, Chen-Yu Ho, Ludovit Horváth, Atal Narayan Sahu, Marco Canini, and Peter Richtárik. Natural compression for distributed deep learning. In *Mathematical and Scientific Machine Learning*, pp. 129–141. PMLR, 2022.
 - Junxian Huang, Feng Qian, Alexandre Gerber, Z Morley Mao, Subhabrata Sen, and Oliver Spatscheck. A close examination of performance and power characteristics of 4g lte networks. In *Proceedings of the 10th international conference on Mobile systems, applications, and services*, pp. 225–238, 2012.
 - Xinmeng Huang, Yiming Chen, Wotao Yin, and Kun Yuan. Lower bounds and nearly optimal algorithms in distributed learning with communication compression. *Advances in Neural Information Processing Systems*, 35:18955–18969, 2022a.
 - Xinmeng Huang, Yiming Chen, Wotao Yin, and Kun Yuan. Lower bounds and nearly optimal algorithms in distributed learning with communication compression. *Advances in Neural Information Processing Systems*, 35:18955–18969, 2022b.
 - Peter J Huber. Robust estimation of a location parameter. *The Annals of Mathematical Statistics*, pp. 73–101, 1964.
 - Nikita Ivkin, Daniel Rothchild, Enayat Ullah, Vladimir braverman, Ion Stoica, and Raman Arora. Communication-efficient distributed SGD with sketching. *Advances in Neural Information Processing Systems*, 32, 2019.
 - Prateek Jain, Dheeraj Nagaraj, and Praneeth Netrapalli. Making the last iterate of SGD information theoretically optimal. In *Conference on Learning Theory*, pp. 1752–1755. PMLR, 2019.
 - Xiaowen Jiang and Sebastian U Stich. Adaptive SGD with polyak stepsize and line-search: Robust convergence and variance reduction. *Advances in Neural Information Processing Systems*, 36, 2024.
 - Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. Advances and open problems in federated learning. *Foundations and trends® in machine learning*, 14(1–2):1–210, 2021.

- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott
 Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models.
 arXiv preprint arXiv:2001.08361, 2020.
 - Sai Praneeth Karimireddy, Quentin Rebjock, Sebastian Stich, and Martin Jaggi. Error feedback fixes signsgd and other gradient compression schemes. In *International Conference on Machine Learning*, pp. 3252–3261. PMLR, 2019.
 - Ahmed Khaled, Konstantin Mishchenko, and Peter Richtárik. Tighter theory for local SGD on identical and heterogeneous data. In *The 23rd International Conference on Artificial Intelligence and Statistics (AISTATS 2020)*, 2020.
 - Ahmed Khaled, Konstantin Mishchenko, and Chi Jin. Dowg unleashed: An efficient universal parameter-free gradient descent method. *Advances in Neural Information Processing Systems*, 36: 6748–6769, 2023.
 - Sarit Khirirat, Hamid Reza Feyzmahdavian, and Mikael Johansson. Distributed learning with compressed gradients. *arXiv preprint arXiv:1806.06573*, 2018.
 - Diederik P. Kingma and Jimmy Ba. Adam: a method for stochastic optimization. In *The 3rd International Conference on Learning Representations*, 2014. URL https://arxiv.org/pdf/1412.6980.pdf.
 - KC Kiwiel and PO Lindberg. Parallel subgradient methods for convex optimization. In *Studies in Computational Mathematics*, volume 8, pp. 335–344. Elsevier, 2001.
 - Roger Koenker and Gilbert Bassett Jr. Regression quantiles. *Econometrica: journal of the Econometric Society*, pp. 33–50, 1978.
 - Anastasia Koloskova, Nicolas Loizou, Sadra Boreiri, Martin Jaggi, and Sebastian Stich. A unified theory of decentralized SGD with changing topology and local updates. In *International Conference on Machine Learning*, pp. 5381–5393. PMLR, 2020.
 - Jakub Konečný and Peter Richtárik. Randomized distributed mean estimation: Accuracy vs. communication. *Frontiers in Applied Mathematics and Statistics*, 4:62, 2018.
 - Jakub Konečný, H. Brendan McMahan, Daniel Ramage, and Peter Richtárik. Federated optimization: distributed machine learning for on-device intelligence. *arXiv:1610.02527*, 2016a.
 - Jakub Konečný, H. Brendan McMahan, Felix Yu, Peter Richtárik, Ananda Theertha Suresh, and Dave Bacon. Federated learning: strategies for improving communication efficiency. In NIPS Private Multi-Party Machine Learning Workshop, 2016b.
 - Simon Lacoste-Julien, Mark Schmidt, and Francis Bach. A simpler approach to obtaining an o (1/t) convergence rate for the projected stochastic subgradient method. *arXiv preprint arXiv:1212.2002*, 2012.
 - Guanghui Lan. First-order and stochastic optimization methods for machine learning, volume 1. Springer, 2020.
 - Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444, 2015.
 - Zhize Li and Peter Richtárik. CANITA: Faster rates for distributed convex optimization with communication compression. *Advances in Neural Information Processing Systems*, 34:13770–13781, 2021.
 - Zhize Li, Dmitry Kovalev, Xun Qian, and Peter Richtárik. Acceleration for compressed gradient descent in distributed and federated optimization. In *International Conference on Machine Learning (ICML)*, pp. 5895–5904. PMLR, 2020.
 - Yujun Lin, Song Han, Huizi Mao, Yu Wang, and Bill Dally. Deep gradient compression: Reducing the communication bandwidth for distributed training. In *International Conference on Learning Representations*, 2018.

- Xiaorui Liu, Yao Li, Jiliang Tang, and Ming Yan. A double residual compression algorithm for efficient distributed learning. In *International Conference on Artificial Intelligence and Statistics*, pp. 133–143. PMLR, 2020.
 - Nicolas Loizou, Sharan Vaswani, Issam Hadj Laradji, and Simon Lacoste-Julien. Stochastic Polyak step-size for SGD: An adaptive learning rate for fast convergence. In *International Conference on Artificial Intelligence and Statistics*, pp. 1306–1314. PMLR, 2021.
 - Brendan McMahan, Eider Moore, Daniel Ramage, and Blaise Agüera y Arcas. Federated learning of deep networks using model averaging. *arXiv preprint arXiv:1602.05629*, 2016.
 - H Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2017.
 - Konstantin Mishchenko and Aaron Defazio. Prodigy: An expeditiously adaptive parameter-free learner. In *Forty-first International Conference on Machine Learning*.
 - Konstantin Mishchenko, Eduard Gorbunov, Martin Takáč, and Peter Richtárik. Distributed learning with compressed gradient differences. *arXiv preprint arXiv:1901.09269*, 2019.
 - Konstantin Mishchenko, Filip Hanzely, and Peter Richtárik. 99% of worker-master communication in distributed optimization is not needed. In *Conference on Uncertainty in Artificial Intelligence*, pp. 979–988. PMLR, 2020.
 - Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th international conference on machine learning (ICML-10)*, pp. 807–814, 2010.
 - Arvind Narayanan, Xumiao Zhang, Ruiyang Zhu, Ahmad Hassan, Shuowei Jin, Xiao Zhu, Xiaoxuan Zhang, Denis Rybkin, Zhengxuan Yang, Zhuoqing Morley Mao, et al. A variegated look at 5g in the wild: performance, power, and qoe implications. In *Proceedings of the 2021 ACM SIGCOMM 2021 Conference*, pp. 610–625, 2021.
 - Angelia Nedic and Asuman Ozdaglar. Distributed subgradient methods for multi-agent optimization. *IEEE Transactions on Automatic Control*, 54(1):48–61, 2009.
 - Arkadi Nemirovski, Anatoli Juditsky, Guanghui Lan, and Alexander Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on optimization*, 19(4):1574–1609, 2009.
 - Yurii Nesterov. *Introductory lectures on convex optimization: A basic course*, volume 87. Springer Science & Business Media, 2013.
 - Dimitris Oikonomou and Nicolas Loizou. Stochastic polyak step-sizes and momentum: Convergence guarantees and practical performance. *arXiv preprint arXiv:2406.04142*, 2024.
 - Andrei Panferov, Yury Demidovich, Ahmad Rammal, and Peter Richtárik. Correlated quantization for faster nonconvex distributed optimization. *arXiv preprint arXiv:2401.05518*, 2024.
 - Constantin Philippenko and Aymeric Dieuleveut. Preserved central model for faster bidirectional compression in distributed settings. *Advances in Neural Information Processing Systems*, 34: 2387–2399, 2021.
- Boris T Polyak. Introduction to optimization. 1987.
- Alexander Rakhlin, Ohad Shamir, and Karthik Sridharan. Making gradient descent optimal for strongly convex stochastic optimization. *arXiv preprint arXiv:1109.5647*, 2011.
 - Nathan D Ratliff, J Andrew Bagnell, and Martin A Zinkevich. (approximate) subgradient methods for structured prediction. In *Artificial Intelligence and Statistics*, pp. 380–387. PMLR, 2007.
 - Sashank J Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. In *International Conference on Learning Representations*, 2018.

- Amirhossein Reisizadeh, Aryan Mokhtari, Hamed Hassani, Ali Jadbabaie, and Ramtin Pedarsani. Fedpaq: A communication-efficient federated learning method with periodic averaging and quantization. In *International conference on artificial intelligence and statistics*, pp. 2021–2031. PMLR, 2020.
 - Peter Richtárik, Igor Sokolov, and Ilyas Fatkhullin. EF21: A new, simpler, theoretically better, and practically faster error feedback. *Advances in Neural Information Processing Systems*, 34: 4384–4396, 2021.
 - Peter Richtárik, Igor Sokolov, Elnur Gasanov, Ilyas Fatkhullin, Zhize Li, and Eduard Gorbunov. 3PC: Three point compressors for communication-efficient distributed training and a better theory for lazy aggregation. In *International Conference on Machine Learning*, pp. 18596–18648. PMLR, 2022.
 - Peter Richtárik, Abdurakhmon Sadiev, and Yury Demidovich. A unified theory of stochastic proximal point methods without smoothness. *arXiv preprint arXiv:2405.15941*, 2024.
 - Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of mathematical statistics*, pp. 400–407, 1951.
 - Lawrence Roberts. Picture coding using pseudo-random noise. *IRE Transactions on Information Theory*, 8(2):145–154, 1962.
 - Daniel Rothchild, Ashwinee Panda, Enayat Ullah, Nikita Ivkin, Ion Stoica, Vladimir Braverman, Joseph Gonzalez, and Raman Arora. FetchSGD: Communication-efficient federated learning with sketching. In *International Conference on Machine Learning*, pp. 8253–8265. PMLR, 2020.
 - Leonid I Rudin, Stanley Osher, and Emad Fatemi. Nonlinear total variation based noise removal algorithms. *Physica D: nonlinear phenomena*, 60(1-4):259–268, 1992.
 - Abdurakhmon Sadiev, Grigory Malinovsky, Eduard Gorbunov, Igor Sokolov, Ahmed Khaled, Konstantin Burlachenko, and Peter Richtárik. Federated optimization algorithms with random reshuffling and gradient compression. *arXiv preprint arXiv:2206.07021*, 2022.
 - Mher Safaryan, Rustem Islamov, Xun Qian, and Peter Richtárik. FedNL: Making Newton-type methods applicable to federated learning. *arXiv preprint arXiv:2106.02969*, 2021.
 - Mher Safaryan, Egor Shulgin, and Peter Richtárik. Uncertainty principle for communication compression in distributed and federated learning and the search for an optimal compressor. *Information and Inference: A Journal of the IMA*, 11(2):557–580, 2022.
 - Fabian Schaipp, Robert M. Gower, and Michael Ulbrich. A stochastic proximal polyak step size. *Transactions on Machine Learning Research*, 2023.
 - Dominik Scherer, Andreas Müller, and Sven Behnke. Evaluation of pooling operations in convolutional architectures for object recognition. In *International conference on artificial neural networks*, pp. 92–101. Springer, 2010.
 - Frank Seide, Hao Fu, Jasha Droppo, Gang Li, and Dong Yu. 1-bit stochastic gradient descent and its application to data-parallel distributed training of speech DNNs. In *Fifteenth Annual Conference of the International Speech Communication Association*, 2014.
 - Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.
 - Shai Shalev-Shwartz, Yoram Singer, and Nathan Srebro. Pegasos: Primal estimated sub-gradient solver for svm. In *Proceedings of the 24th international conference on Machine learning*, pp. 807–814, 2007.
 - N. Z. Shor, Krzysztof C. Kiwiel, and Andrzej Ruszczyński. Minimization methods for nondifferentiable functions, 1985.
 - Sebastian U. Stich. Local SGD converges fast and communicates little. In *International Conference on Learning Representations*, 2020.

- Sebastian U Stich and Sai Praneeth Karimireddy. The error-feedback framework: Better rates for SGD with delayed gradients and compressed communication. *arXiv preprint arXiv:1909.05350*, 2019.
 - Sebastian U Stich, Jean-Baptiste Cordonnier, and Martin Jaggi. Sparsified SGD with memory. *Advances in neural information processing systems*, 31, 2018.
 - Nikko Ström. Scalable distributed DNN training using commodity GPU cloud computing. In *Interspeech 2015*, 2015.
 - Rafał Szlendak, Alexander Tyurin, and Peter Richtárik. Permutation compressors for provably faster distributed nonconvex optimization. In *International Conference on Learning Representations*, 2022.
 - Hanlin Tang, Chen Yu, Xiangru Lian, Tong Zhang, and Ji Liu. Doublesqueeze: Parallel stochastic gradient descent with double-pass error-compensated compression. In *International Conference on Machine Learning*, pp. 6155–6165. PMLR, 2019.
 - Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 58(1):267–288, 1996.
 - Alexander Tyurin and Peter Richtárik. DASHA: Distributed nonconvex optimization with communication compression and optimal oracle complexity. In *The Eleventh International Conference on Learning Representations*, 2023.
 - Alexander Tyurin and Peter Richtarik. 2Direction: Theoretically faster distributed training with bidirectional communication compression. *Advances in Neural Information Processing Systems*, 36:11737–11808, 2023.
 - Vladimir Vapnik. The nature of statistical learning theory. Springer science & business media, 2013.
 - Thijs Vogels, Sai Praneeth Karimireddy, and Martin Jaggi. PowerSGD: Practical low-rank gradient compression for distributed optimization. *Advances in Neural Information Processing Systems*, 32, 2019.
 - Thijs Vogels, Sai Praneeth Karimireddy, and Martin Jaggi. Practical low-rank communication compression in decentralized deep learning. *Advances in Neural Information Processing Systems*, 33:14171–14181, 2020.
 - Evgeniya Vorontsova, Roland Hildebrand, Alexander Gasnikov, and Fedor Stonyakin. Convex optimization. *arXiv preprint arXiv:2106.01946*, 2021.
 - Hongyi Wang, Scott Sievert, Shengchao Liu, Zachary Charles, Dimitris Papailiopoulos, and Stephen Wright. Atomo: Communication-efficient learning via atomic sparsification. *Advances in neural information processing systems*, 31, 2018.
 - Jianqiao Wangni, Jialei Wang, Ji Liu, and Tong Zhang. Gradient sparsification for communication-efficient distributed optimization. *Advances in Neural Information Processing Systems*, 31, 2018.
 - Wei Wen, Cong Xu, Feng Yan, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. Terngrad: Ternary gradients to reduce communication in distributed deep learning. *Advances in neural information processing systems*, 30, 2017.
 - Blake Woodworth, Kumar Kshitij Patel, Sebastian Stich, Zhen Dai, Brian Bullins, Brendan Mcmahan, Ohad Shamir, and Nathan Srebro. Is local SGD better than minibatch SGD? In *International Conference on Machine Learning*, pp. 10334–10343. PMLR, 2020.
 - Jiaxiang Wu, Weidong Huang, Junzhou Huang, and Tong Zhang. Error compensated quantized SGD and its applications to large-scale distributed optimization. In *International conference on machine learning*, pp. 5325–5333. PMLR, 2018.
 - Zhaoyue Xia, Jun Du, and Yong Ren. Convergence theory of generalized distributed subgradient method with random quantization. *arXiv preprint arXiv:2207.10969*, 2022.

- Zhaoyue Xia, Jun Du, Chunxiao Jiang, H Vincent Poor, Zhu Han, and Yong Ren. Distributed subgradient method with random quantization and flexible weights: Convergence analysis. *IEEE Transactions on Cybernetics*, 2023.
- Kai Yi, Timur Kharisov, Igor Sokolov, and Peter Richtárik. Cohort squeeze: Beyond a single communication round per cohort in cross-device federated learning. *arXiv preprint arXiv:2406.01115*, 2024.
- Moslem Zamani and François Glineur. Exact convergence rate of the last iterate in subgradient methods. *arXiv preprint arXiv:2307.11134*, 2023.
- Hantian Zhang, Jerry Li, Kaan Kara, Dan Alistarh, Ji Liu, and Ce Zhang. The ZipML framework for training models with end-to-end low precision: The cans, the cannots, and a little bit of deep learning. *arXiv preprint arXiv:1611.05402*, 2016.
- Shuai Zheng, Ziyue Huang, and James Kwok. Communication-efficient distributed blockwise momentum SGD with error-feedback. *Advances in Neural Information Processing Systems*, 32, 2019.
- Yuchen Zheng, Yujia Xie, Ilbin Lee, Amin Dehghanian, and Nicoleta Serban. Parallel subgradient algorithm with block dual decomposition for large-scale optimization. *European journal of operational research*, 299(1):60–74, 2022.
- Zhihan Zhu, Yanhao Zhang, and Yong Xia. Convergence rate of projected subgradient method with time-varying step-sizes. *Optimization Letters*, pp. 1–5, 2024.
- Hui Zou and Trevor Hastie. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 67(2):301–320, 2005.