

Real Acceleration of Communication Process in Distributed Algorithms with Compression*

Svetlana Tkachenko¹ Artem Andreev¹ Aleksandr Beznosikov^{1,2} Alexander Gasnikov^{1,2}

Institute for Information Transmission Problems, Moscow, Russia

Abstract. Modern applied optimization problems become more and more complex every day. Due to this fact, distributed algorithms that can speed up the process of solving an optimization problem through parallelization are of great importance. The main bottleneck of distributed algorithms is communications, which can slow down the method dramatically. One way to solve this issue is to use compression of transmitted information. In the current literature on theoretical distributed optimization, it is generally accepted that as much as we compress information, so much we reduce communication time. But in reality, the communication time depends not only on the size of the transmitted information, but also, for example, on the message startup time. In this paper, we study distributed optimization algorithms under the assumption of a more complex and closer-to-reality dependence of transmission time on compression. In particular, we describe the real speedup achieved by compression, analyze how much it makes sense to compress information, and present an adaptive way to select the power of compression depending on unknown or changing parameters of the communication process.

Keywords: distributed optimization · compression · acceleration

1 Introduction

Modern realities pose more and more complex optimization problems that need to be solved. For example, to improve the generalization of deployed models, machine learning engineers need to rely on training datasets of ever increasing sizes and on elaborate large-scale over-parametrized models [3]. Therefore, it is increasingly necessary to resort to the use of distributed approaches to solving the optimization problem. The essence of distributed optimization is to the process of streamlining the target function by using multiple computing resources that are scattered across different machines or servers. It enables optimization algorithms to run in parallel, which can greatly increase the speed and efficiency of finding the optimal solution. Therefore, distributed optimization is widely used in various domains, including machine learning, data science, and operations research [20].

*The research was supported by Russian Science Foundation (project No. 23-11-00229).

However, when utilizing parallel computation in a distributed optimization environment, a common challenge is the communication between the computational devices. Since the agents function independently, they must exchange information to harmonise their local solutions and revise the global solution. Meanwhile, communication time is a waste that prevents full parallelization. Therefore, to struggle for effective communication and to address the communication bottleneck issue is a key point in distributed optimization [13,19,9].

Employing compression of forwarded information is one of the viable solutions to decrease communication expenses [18,1]. It assists in reducing file size while preserving important information. With the use of effective compression algorithms, transmission time can be considerably reduced both in theory and in practice [10].

Several models describing the dependence of transmission time on message size can be found in literature. The most frequently utilized model in the theoretical optimization is $T = \beta s$, where T is the transmission time, β is the delay-size relationship, and s is the message size. Meanwhile, there is a more practical and widespread model that has stayed away from theoretical optimization. This model is $T = \beta s + \alpha$, where α is the server initialization time [8]. The simpler model indicate that transmission time can be reduced by a factor of n by transmitting n times less information. Nevertheless, the practical results contrast with the theoretical ones. In actuality, messaging involves initializing the channel, which refers to establishing a connection between the sender and the recipient. The second model accounts for this. This implies that there is minimal distinction when transmitting 1 or 2 bits, but once we send 100 Mb and 200 Mb, the variance is substantial. Therefore, we necessitate an accurate representation to characterise the communications.

Considering the issue of communication expenses and the proposed solution, the main questions of this study can be posed:

1. *Which model better describes the real world of messaging?*
2. *How does this change the theory of distributed optimization?*
3. *How can we determine the parameters of this model?*
4. *What is the most efficient method of calculating the parameters of this model in the event of frequent data updates?*

1.1 Contributions

More practical communication model: Instead of the classical delay versus size model of $T = \beta s$ (where β represents the relationship between delay and size), we have adopted a more realistic approach of $T = \beta s + \alpha$ (taking into account α – the server initialization time). When a worker sends a message to the server, the channel initialization time contributes significantly to the small size of message transmission. It is, therefore, essential to consider this factor.

Impact of model on communication complexities: We analyze how the more practical model from the previous paragraph affects the communication

costs of modern distributed algorithms with compression. We consider state-of-the-art methods that have the best theoretical guarantees for convex and non-convex problems.

Estimate of α , β : In order to calculate the two coefficients α and β based on the real data on the dependence of delay on message volume, we assume that $\alpha = \alpha_{const} + \delta_\alpha$ and $\beta = \beta_{const} + \delta_\beta$ where δ_α and δ_β follow independent normal distributions. α_{const} and β_{const} are the true values of these coefficients, δ_α and δ_β are errors of measurement or calculation.

Using this information, it is possible to calculate the coefficients α_{const} and β_{const} through statistical techniques, such as least squares. Rather than storing the complete dataset of message size and delay time, we can update the summation of variables. This approach enables us to update just four variables without needing to recalculate the coefficients using the least squares method.

The estimation process aims to find the values of α and β that best fit the observed data, thus providing insight into the server initialization time and the delay-volume relationship.

2 Problem statement

We consider the optimization problems of the form:

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

where x is the optimization variable. For example, in the context of ML, $x \in \mathbb{R}^d$ contains the parameters of the statistical model to be trained, n – number of employees/devices and functions $f_i(x) : \mathbb{R}^d \rightarrow \mathbb{R}$ – model data loss x , stored on the device i .

2.1 Distributed optimization with compression

We give an illustration of the traditional use of distributed optimization, utilizing the gradient descent algorithm as an instance – see Algorithm 1.

As noted above to handle large data sets, compression is necessary. In Algorithm 1 this can be represented by the compression operator $\mathcal{C} : \mathbb{R}^d \rightarrow \mathbb{R}^d$. In particular, $\nabla f_i(x^k)$ in line 4 should be replaced by $\mathcal{C}(\nabla f_i(x^k))$ and, accordingly, in line 6 we aggregate the compressed gradients $\mathcal{C}(\nabla f_i(x^k))$:

$$x^{k+1} = x^k - \gamma_k \cdot \frac{1}{n} \sum_{i=1}^n \mathcal{C}(\nabla f_i(x^k)).$$

This approach is basic, but does not give the best convergence results [11,6]. One can note that more advanced methods with compression use more tricky schemes, in particular, they are based on various variance reduction techniques,

Algorithm 1

- 1: **Initialization:** choose $x^0 \in \mathbb{R}^d$ and stepsizes $\{\gamma_k\}_{k=0}^K$
- 2: **for** $k = 0, 1, \dots, K$ **do**
- 3: Server sends x^k to all n nodes
- 4: Each i -th node, in parallel with the others, calculates the gradient of its corresponding function f_i :

$$\nabla f_i(x^k)$$

- 5: All nodes send $\nabla f_i(x^k)$ to the server
- 6: Server performs aggregation:

$$x^{k+1} = x^k - \gamma_k \cdot \frac{1}{n} \sum_{i=1}^n \nabla f_i(x^k)$$

- 7: **end for**
-

which prescribe to compress not the gradient itself, but the difference between the gradient and some reference value [15,14,10,17,7,4,5].

The theory of convergence of methods with compression is based on a formal definition of the properties of \mathcal{C} operators. In particular, two classes of operators: unbiased and biased, are often distinguished in the literature.

Definition 1. \mathcal{C} is an unbiased compression with $\zeta \geq 1$ if \mathcal{C} is unbiased ($\mathbb{E}[\mathcal{C}(x)] = x$) and $\mathbb{E}[\|\mathcal{C}(x)\|_2^2] \leq \zeta \|x\|_2^2$ for all $x \in \mathbb{R}^d$.

Definition 2. \mathcal{C} is a biased compression with $\delta \geq 1$ if $\mathbb{E}[\|\mathcal{C}(x) - x\|_2^2] \leq (1 - 1/\delta) \|x\|_2^2$ for all $x \in \mathbb{R}^d$.

Meanwhile, these definitions do not give a complete picture about compression operators. The definitions are interesting for proving convergence and obtaining iterative complexity of algorithms. But to obtain the communication cost in the amount of transmitted information, it is necessary to understand how much the operator reduces the transmitted information.

2.2 Degree of compression

In this subsection, we estimate the degree of compression $\omega_{inf} = \frac{\text{len}(x)}{\text{len}(\mathcal{C}(x))}$, where $\text{len}(x)$ is the number of bits of information to send $x \in \mathbb{R}^d$. We consider different classical compression operators.

Definition 3. For $k \in [d] := \{1, \dots, d\}$, the unbiased random (aka Rand- k) sparsification operator is defined via

$$\mathcal{C}(x) := \frac{d}{k} \sum_{i \in S} x_i e_i,$$

where $S \subseteq [d]$ is the k -nice sampling; i.e., a subset of $[d]$ of cardinality k chosen uniformly at random, and e_1, \dots, e_d are the standard unit basis vectors in \mathbb{R}^d .

Lemma 1. For the unbiased random sparsification $\omega_{inf} = \frac{d}{k}$.

Proof. Initial vector x contains d non-zero coordinates, and the compressed one contains k . Then $\omega_{inf} = \frac{d}{k}$. Here it is important to clarify that in the general case it is necessary to forward numbers of non-zero coordinates as well. But if the same random generator with the same seed is installed on the sending and receiving devices, it is possible to synchronize the randomness for free, and then there is no need to send additional information.

Definition 4 (see [2]). Top- k sparsification operator is defined via

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

where coordinates are ordered by their magnitudes so that $|x_{(1)}| \leq |x_{(2)}| \leq \dots \leq |x_{(d)}|$.

Top- k is a greedy version of unbiased random sparsification.

Lemma 2. For Top- k sparsification $\omega_{inf} = \frac{d \cdot \text{len}(x)}{k \cdot \text{len}(x) + k \cdot \lceil \log_2 d \rceil}$.

Proof. Similar to Unbiased random sparsification initial vector x contains d non-zero coordinates, and the compressed one contains k . But here, unlike random sparsification, we have to pass the numbers of selected non-zero coordinates. To encode the numbers from 1 to d , $\lceil \log_2 d \rceil$ bits are needed. The total number of transmitted bits is $k \cdot \text{len}(x) + k \cdot \lceil \log_2 d \rceil$. Then $\omega_{inf} = \frac{d \cdot \text{len}(x)}{k \cdot \text{len}(x) + k \cdot \lceil \log_2 d \rceil}$.

Definition 5 (see [12]). Natural compression operator \mathcal{C}_{nat} is defined as follows:

$$\mathcal{C}(x) = \begin{cases} \text{sign}(x) \cdot 2^{\lceil \log_2 |x| \rceil}, & \text{with } p(x), \\ \text{sign}(x) \cdot 2^{\lceil \log_2 |x| \rceil}, & \text{with } 1 - p(x), \end{cases}$$

where probability $p(x) := \frac{2^{\lceil \log_2 |x| \rceil} - |x|}{2^{\lceil \log_2 |x| \rceil}}$.

The essence of this compression is random rounding to the nearest power of two. In terms of computing on a computer with 32bit float type, this is simply equivalent to using only the sign bit and 8 bits from the exponent.

Lemma 3. For Natural compression $\omega_{inf} = \frac{32}{9}$.

Proof. The statement follows directly from the use of such compression with 32bit float. Instead of 32 bits we send 9.

Definition 6 (see [21]). Rank- r Power compression introduced by [21] is a compressed-decompressed approach based on the low-rank approximate decomposition of the matrix $X \in \mathbb{R}^{n \times m}$ (transformed version of the original parcel vector x).

Lemma 4. For Rank- r PowerSGD compression $\omega_{inf} = \frac{nm}{r(n+m)}$.

Proof. The product of matrices PQ^T , $P \in \mathbb{R}^{n \times r}$, $Q \in \mathbb{R}^{m \times r}$ approximates the matrix $X \in \mathbb{R}^{n \times m}$. Thus, instead of storing $n \cdot m$ numbers must be $r \cdot n + r \cdot m$. Then $\omega_{inf} = \frac{nm}{r(n+m)}$.

The results obtained above are summarized in Table 1.

Compression operator	ω_{inf}
Unbiased random sparsification	$\frac{d}{k}$
Top- k sparsification [2]	$\frac{d \cdot \text{len}(x)}{k \cdot \text{len}(x) + k \cdot \lceil \log_2 d \rceil}$
Natural compression [12]	$\frac{32}{9}$
Rank- r Power compression [21]	$\frac{nm}{r(n+m)}$

Table 1: ω_{inf} for different compression operators.

As previously stated, the estimation of communication time necessitates ω_{inf} . This coefficient serves as a reliable indicator of the effectiveness of the compression operator in each instance, particularly if one knows the optimal frequency of message compression (which is the objective of this research paper).

3 Main part

3.1 Transmission time model and convergence complexities

We consider the following model of transmission time:

$$T(s) = \alpha + \beta \cdot s,$$

where s is the size of the packages, β represents the time to transmit one unit of the information, α represents the time to initialize the channel, which is the delay that occurs before any message transmission occurs. This delay may involve activities such as creating a connection, verifying the user's identity, or loading essential resources. As mentioned earlier, in papers on theoretical

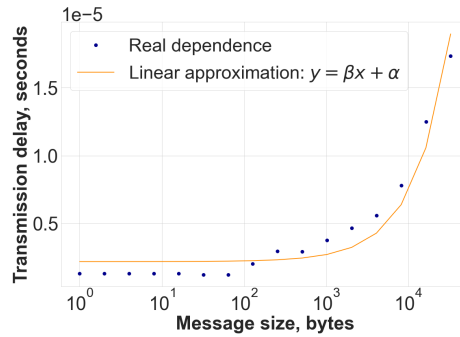


Fig. 1: Dependence of communication time on the size of the transmitted messages for MSU supercomputer "Lomonosov": blue dots – real values, orange line – approximation.

optimization and convergence estimates, it is assumed that $\alpha = 0$. But let us examine the plot presented in Figure 1, which portrays the relationship between the transmission delay and message size in a live network. Using this plot we can see the effect of α on communication time.

Note that the theoretical results on the communication cost of distributed algorithms with compression depend on the parameter

$$\eta = \frac{T(\text{len}(\nabla f_i(x)))}{T(\text{len}(\nabla f_i(x))/w_{inf})}. \quad (1)$$

In particular, the best results for method with unbiased compression for convex [14] and non-convex [10] problems linear depends on $\left(\frac{1}{\eta} + \frac{\zeta}{\eta\sqrt{n}}\right)$. The state-of-the-art results for biased compression [16,17] linear depends on $\left(\frac{1}{\eta} + \frac{\delta}{\eta}\right)$. It is easy to see that if $\alpha = 0$, the expression (1) gives $\eta = w_{inf}$. But if $\alpha \neq 0$, it is possible that $\eta = 1$, thus the impact of even a large w_{inf} can be almost completely canceled.

3.2 Division into areas

Let us examine the original plot (Figure 1) and divide it into three conditional ranges (Figure 2). In the first range, the coefficient α is the most significant. This means that if the size of the message s falls within this range, then α greatly exceeds βs . In the second range, both coefficients α and β hold value, with α and βs being close in value. In the third range, β carries the most significance, with βs greatly exceeding α .

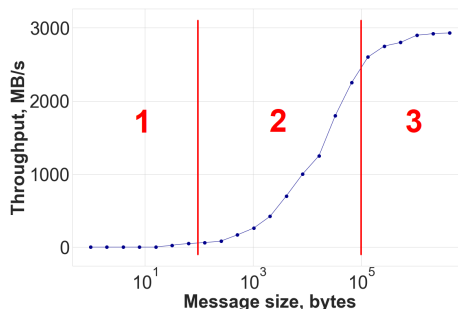


Fig. 2: Division into fields according to the importance of the summands α and βs

Let us examine how the transmission delay varies with changes in message size (see Figure 3, Figure 4, and Figure 5). We determine the number of times the message size alters during compression, and subsequently how many times the transmission time changes:

- When transitioning between areas 3 to 3 and 3 to 2, the message is compressed by a factor of n , while communication time is reduced by $0,95 \cdot n$.
- When transitioning between areas 2 to 2 and 2 to 1, compression is approximately 40 times greater than the reduction in communication cost.
- When transitioning between areas 3 to 1, the time saved is insignificant compared to the compression size (with a compression of 5000 times, the transmission time is only reduced by approximately 100 times).

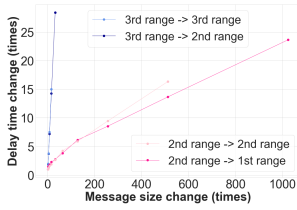


Fig. 3: Transitions from area 2 to areas 1, 2

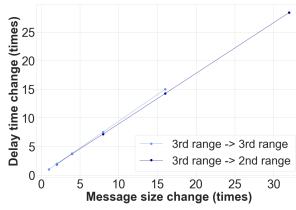


Fig. 4: Transitions from area 3 to areas 2, 3

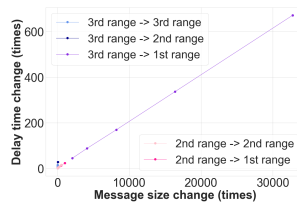


Fig. 5: Transitions from area 3 to area 1

Conclusion: It is most feasible to travel from area 3 to 3 or 2, and it is not financially viable to travel from area 3 to 1.

Let us consider an example how compression of a message affects transmission time. We use the distributed system from Figure 1. The size of an uncompressed message is 10^{22} . Figure 6 demonstrates that the variation in time is almost linear at first, but then the compression loses its effectiveness. Hence, it can be concluded that the compression of a message has an impact on the transmission time.

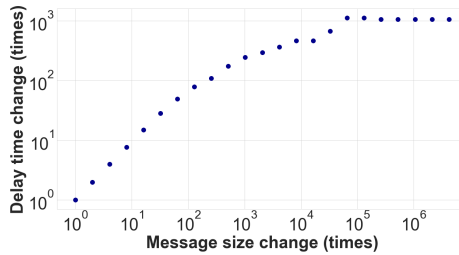


Fig. 6: Transmission delay reduction on message size reduction

Conclusion: High compression does not result in a significant time savings, thus extensively compressing a message is not an efficient approach.

3.3 A way to find α and β for an unknown network parameters

Here is a method for determining α and β when the network parameters are unknown, when we can specify areas as in Figure 2. We formalize the problem of finding or estimating α and β as follows.

Condition: It is possible to transmit messages of varying sizes ranging from 0 to P_{max} , which represents the maximum message size that we can send. It is imperative to consider technological constraints when evaluating the feasibility of message transmission. For instance, for each message, the values α and β are stochastic and have the laws: $\alpha(t) = \alpha_{const} + \delta\alpha$ and $\beta(t) = \beta_{const} + \delta\beta$. That is, $\alpha(t)$ and $\beta(t)$ vary among samples and consist of a constant value plus stochastic noise δ (which follows a normal distribution with mean 0 and variance σ^2 , where $\sigma = \alpha_m \cdot \alpha_{const}$ or $\sigma = \beta_m \cdot \beta_{const}$ depending on the nature of δ).

Suggested solution: Let us apply the formulas of the method of least squares to recalculate α and β . During the operation of the main optimization algorithm we vary message sizes. Firstly, we calculate the delay at 2 points. Then,

for each subsequent step, we select (deterministically or randomly) the next point from within the interval $(0, P_{max})$ and calculate the delay. Here is an example of code that executes this algorithm:

Algorithm 2

- 1: **Parameters:** largest message size P_{max} ;
 - 2: **Initialization:** compute times y_1, y_2 for message sizes x_1, x_2 respectively ($y = \beta \cdot x + \alpha$). Set $s_x^2 = x_1 + x_2$, $s_y^2 = y_1 + y_2$, $s_{xy}^2 = y_1 \cdot x_1 + y_2 \cdot x_2$, $s_{xx}^2 = x_1 \cdot x_1 + x_2 \cdot x_2$;
 - 3: **for** $k = 3, 4, \dots$ **do**
 - 4: **for** new $x_k \in [0, P_{max}]$, compute y_k
 - 5: $s_x^k = s_x^{k-1} + x_k$
 - 6: $s_y^k = s_y^{k-1} + y_k$
 - 7: $s_{xy}^k = s_{xy}^{k-1} + x_k \cdot y_k$
 - 8: $s_{xx}^k = s_{xx}^{k-1} + x_k \cdot x_k$
 - 9: $\beta_k = \frac{k \cdot s_{xy}^k - s_x^k \cdot s_y^k}{k \cdot s_{xx}^k - (s_x^k)^2}$
 - 10: $\alpha_k = \frac{s_y^k - \beta_k \cdot s_x^k}{k}$
 - 11: **end for**
-

The algorithm recalculates the α and β coefficients using the least squares formulas. It is worth pointing out that it is very expensive to recalculate the parameters α and β using the least squares method and to store all data of message sizes and times of transmission $\{x_i, y_i\}$. But Algorithm 2 can work online. We need only 4 variables: the sum of message sizes s_x , the sum of delays s_y , and the sums needed for the least squares calculation s_{xy} and s_{xx} .

Proposition 1. β_k, α_k from Algorithm 2 are unbiased estimations of β and α , namely $\mathbb{E}[\beta_k] = \beta_{const}$ and $\mathbb{E}[\alpha_k] = \alpha_{const}$.

Proof. We start from β_k :

$$\mathbb{E}[\beta_k] = \mathbb{E} \left[\frac{k \cdot s_{xy}^k - s_x^k \cdot s_y^k}{k \cdot s_{xx}^k - (s_x^k)^2} \right] = \frac{k \cdot \mathbb{E}[s_{xy}^k] - s_x^k \cdot \mathbb{E}[s_y^k]}{k \cdot s_{xx}^k - (s_x^k)^2}. \quad (2)$$

Next, we estimate $\mathbb{E}[s_{xy}^k]$ and $s_x^k \cdot \mathbb{E}[s_y^k]$

$$\begin{aligned} \mathbb{E}[s_{xy}^k] &= \mathbb{E} \left[\sum_{i=1}^k (x_i y_i) \right] = \mathbb{E} \left[\sum_{i=1}^k (x_i \mathbb{E}[y_i]) \right] = \sum_{i=1}^k x_i (\beta_{const} x_i + \alpha_{const}) \\ &= \beta_{const} \sum_{i=1}^k x_i^2 + \alpha_{const} \sum_{i=1}^k x_i, \end{aligned} \quad (3)$$

$$\begin{aligned} s_x^k \cdot \mathbb{E}[s_y^k] &= \left(\sum_{i=1}^k x_i \right) \cdot \mathbb{E} \left[\sum_{i=1}^k y_i \right] = \left(\sum_{i=1}^k x_i \right) \cdot \sum_{i=1}^k (\beta_{const} x_i + \alpha_{const}) = \\ &= \beta_{const} \left(\sum_{i=1}^k x_i \right)^2 + k \cdot \alpha_{const} \sum_{i=1}^k x_i. \end{aligned} \quad (4)$$

Substituting (4) and (3) to (2), we get

$$\mathbb{E}[\beta_k] = \frac{k \cdot (\beta_{const} \cdot s_{xx}^k + \alpha_{const} \cdot s_x^k) - \beta_{const} \cdot (s_x^k)^2 - k \cdot \alpha_{const} \cdot s_x^k}{k \cdot s_{xx}^k - (s_x^k)^2} = \beta_{const}.$$

Finally, for $\mathbb{E}[\alpha_k]$ we obtain

$$\begin{aligned} \mathbb{E}[\alpha_k] &= \frac{\mathbb{E}[s_y^k - \beta_k \cdot s_x^k]}{k} = \frac{\mathbb{E}[s_y^k] - \mathbb{E}[\beta_k] \cdot s_x^k}{k} \\ &= \frac{\sum_{i=1}^k (\beta_{const} x_i + \alpha_{const}) - \beta_{const} \cdot s_x^k}{k} = \frac{k \cdot \alpha_{const}}{k} = \alpha_{const}. \end{aligned}$$

4 Conclusions

In this paper, we considered a realistic communication cost model $T(s) = \beta s + \alpha$, which takes into account α – the server initialization time. We tried to discuss how it affects to communication time complexities of algorithms. We also provided the algorithm for determining the coefficients α and β utilizing statistical techniques such as the least squares method, alongside estimated uncertainties related to this approach. Rather than storing the complete message size and delay time sets, it is viable to update some combinations of the variables.

One can note that the model considered in this paper can also be improved. For example, we can also include the time required for the communication operator counting. In some cases, this can be quite expensive, which slows down the computational process. Taking this time into account is an important detail for future research.

References

1. Alistarh, D., Grubic, D., Li, J., Tomioka, R., Vojnovic, M.: QSGD: Communication-efficient SGD via gradient quantization and encoding. In: *Advances in Neural Information Processing Systems*. pp. 1709–1720 (2017)
2. Alistarh, D., Hoefler, T., Johansson, M., Khirirat, S., Konstantinov, N., Renggli, C.: The convergence of sparsified gradient methods (2018)
3. Arora, S., Cohen, N., Hazan, E.: On the optimization of deep networks: Implicit acceleration by overparameterization. In: *International Conference on Machine Learning*. pp. 244–253. PMLR (2018)
4. Beznosikov, A., Gasnikov, A.: Compression and data similarity: Combination of two techniques for communication-efficient solving of distributed variational inequalities. *arXiv preprint arXiv:2206.09446* (2022)
5. Beznosikov, A., Gasnikov, A.: Similarity, compression and local steps: Three pillars of efficient communications for distributed variational inequalities. *arXiv preprint arXiv:2302.07615* (2023)
6. Beznosikov, A., Horváth, S., Richtárik, P., Safaryan, M.: On biased compression for distributed learning. *arXiv preprint arXiv:2002.12410* (2020)

7. Beznosikov, A., Richtárik, P., Diskin, M., Ryabinin, M., Gasnikov, A.: Distributed methods with compressed communication for solving variational inequalities, with theoretical guarantees. *Advances in Neural Information Processing Systems* **35**, 14013–14029 (2022)
8. Chan, E., Heimlich, M., Purkayastha, A., Van De Geijn, R.: Collective communication: theory, practice, and experience. *Concurrency and Computation: Practice and Experience* **19**(13), 1749–1783 (2007)
9. Ghosh, A., Maity, R.K., Mazumdar, A., Ramchandran, K.: Communication efficient distributed approximate newton method. In: 2020 IEEE International Symposium on Information Theory (ISIT). pp. 2539–2544. IEEE (2020)
10. Gorbunov, E., Burlachenko, K.P., Li, Z., Richtárik, P.: Marina: Faster non-convex distributed learning with compression. In: International Conference on Machine Learning. pp. 3788–3798. PMLR (2021)
11. Gorbunov, E., Hanzely, F., Richtárik, P.: A unified theory of sgd: Variance reduction, sampling, quantization and coordinate descent. In: International Conference on Artificial Intelligence and Statistics. pp. 680–690. PMLR (2020)
12. Horvath, S., Ho, C.Y., Horvath, L., Sahu, A.N., Canini, M., Richtarik, P.: Natural compression for distributed deep learning (2022)
13. Konečný, J., McMahan, H.B., Yu, F.X., Richtárik, P., Suresh, A.T., Bacon, D.: Federated learning: Strategies for improving communication efficiency. arXiv preprint arXiv:1610.05492 (2016)
14. Li, Z., Kovalev, D., Qian, X., Richtárik, P.: Acceleration for compressed gradient descent in distributed and federated optimization. arXiv preprint arXiv:2002.11364 (2020)
15. Mishchenko, K., Gorbunov, E., Takáč, M., Richtárik, P.: Distributed learning with compressed gradient differences. arXiv preprint arXiv:1901.09269 (2019)
16. Qian, X., Richtárik, P., Zhang, T.: Error compensated distributed sgd can be accelerated. *Advances in Neural Information Processing Systems* **34**, 30401–30413 (2021)
17. Richtárik, P., Sokolov, I., Fatkhullin, I.: EF21: A new, simpler, theoretically better, and practically faster error feedback. arXiv preprint arXiv:2106.05203 (2021)
18. Seide, F., Fu, H., Droppo, J., Li, G., 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)
19. Smith, V., Forte, S., Chenxin, M., Takáč, M., Jordan, M.I., Jaggi, M.: Cocoa: A general framework for communication-efficient distributed optimization. *Journal of Machine Learning Research* **18**, 230 (2018)
20. Verbraeken, J., Wolting, M., Katzy, J., Kloppenburg, J., Verbelen, T., Rellermeyer, J.S.: A survey on distributed machine learning. *Acm computing surveys (csur)* **53**(2), 1–33 (2020)
21. Vogels, T., Karimireddy, S.P., Jaggi, M.: Powersgd: Practical low-rank gradient compression for distributed optimization. *Advances in Neural Information Processing Systems* **32** (2019)