
Towards a Better Theoretical Understanding of Independent Subnetwork Training

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Modern advancements in large-scale machine learning would be impossible without
2 the paradigm of data-parallel distributed computing. Since distributed computing
3 with large-scale models imparts excessive pressure on communication channels, a
4 lot of recent research was directed towards co-designing communication compression
5 strategies and training algorithms with the goal of reducing communication
6 costs. While pure data parallelism allows better data scaling, it suffers from poor
7 model scaling properties. Indeed, compute nodes are severely limited by memory
8 constraints, preventing further increases in model size. For this reason, the latest
9 achievements in training giant neural network models rely on some form of model
10 parallelism as well. In this work, we take a closer theoretical look at Independent
11 Subnetwork Training (IST), which is a recently proposed and highly effective
12 technique for solving the aforementioned problems. We identify fundamental
13 differences between IST and alternative approaches, such as distributed methods
14 with compressed communication, and provide a precise analysis of its optimization
15 performance on a quadratic model.

16 1 Introduction

17 A huge part of today’s machine learning success drives from the possibility to build more and more
18 complex models and train them on increasingly larger datasets. This fast progress has become
19 feasible due to advancements in distributed optimization, which is necessary for proper scaling
20 when the training data sizes grow [50]. In a typical scenario data parallelism is used for efficiency
21 which consists of sharding the dataset across computing devices. This allowed very efficient scaling
22 and accelerating of training moderately sized models by using additional hardware [19]. Though,
23 such data parallel approach can suffer from communication bottleneck, which sparked a lot of
24 research on distributed optimization with compressed communication of the parameters between
25 nodes [3, 27, 38].

26 1.1 The need for model parallel

27 Despite the efficiency gains of data parallelism, it has some fundamental limitations when it comes to
28 scaling up the model size. As the model dimension grows, the amount of memory required to store
29 and update the parameters also increases, which becomes problematic due to resource constraints
30 on individual devices. This has led to the development of model parallelism [11, 37], which splits
31 a large model across multiple nodes, with each node responsible for computations of model parts
32 [15, 47]. However, naive model parallelism also poses challenges because each node can only update
33 its portion of the model based on the data it has access to. This creates a need for a very careful
34 management of communication between devices. Thus, a combination of both data and model
35 parallelism is often necessary to achieve efficient and scalable training of huge models.

Algorithm 1 Distributed Submodel (Stochastic) Gradient Descent

```
1: Parameters: learning rate  $\gamma > 0$ ; sketches  $\mathbf{C}_1, \dots, \mathbf{C}_n$ ; initial model  $x^0 \in \mathbb{R}^d$ 
2: for  $k = 0, 1, 2 \dots$  do
3:   Select submodels  $w_i^k = \mathbf{C}_i^k x^k$  for  $i \in [n]$  and broadcast to all computing nodes
4:   for  $i = 1, \dots, n$  in parallel do
5:     Compute local (stochastic) gradient w.r.t. submodel:  $\mathbf{C}_i^k \nabla f_i(w_i^k)$ 
6:     Take (maybe multiple) gradient descent step  $z_i^+ = w_i^k - \gamma \mathbf{C}_i^k \nabla f_i(w_i^k)$ 
7:     Send  $z_i^+$  to the server
8:   end for
9:   Aggregate/merge received submodels:  $x^{k+1} = \frac{1}{n} \sum_{i=1}^n z_i^+$ 
10: end for
```

36 **IST.** Independent Subnetwork Training (IST) is a technique which suggests dividing the neural
37 network into smaller independent sub-networks, training them in a distributed parallel fashion and then
38 aggregating the results to update the weights of the whole model. According to IST, every subnetwork
39 is operational on its own, has fewer parameters than the full model, and this not only reduces the load
40 on computing nodes but also results in faster synchronization. A generalized analog of the described
41 method is formalized as an iterative procedure in Algorithm 1. This paradigm was pioneered by
42 [45] for networks with fully-connected layers and was later extended to ResNets [14] and Graph
43 architectures [43]. Previous experimental studies have shown that IST is a very promising approach
44 for various applications as it allows to effectively combine data with model parallelism and train
45 larger models with limited compute. In addition, [28] performed theoretical analysis of IST for
46 overparameterized single hidden layer neural networks with ReLU activations. The idea of IST was
47 also recently extended to the federated setting via an asynchronous distributed dropout [13] technique.

48 **Federated Learning.** Another important setting when the data is distributed (due to privacy reasons)
49 is Federated Learning [22, 27, 31]. In this scenario computing devices are often heterogeneous and
50 more resource-constrained [5] (e.g. mobile phones) in comparison to data-center setting. Such
51 challenges prompted extensive research efforts into selecting smaller and more efficient submodels
52 for local on-device training [2, 6, 8, 12, 20, 21, 29, 35, 42, 44]. Many of these works propose
53 approaches to adapt submodels, often tailored to specific neural network architectures, based on
54 the capabilities of individual clients for various machine learning tasks. However, there is a lack of
55 comprehension regarding the theoretical properties of these methods.

56 1.2 Summary of contributions

57 When reviewing the literature, we have found that a rigorous understanding of IST convergence
58 virtually does not exist, which motivates our work. The main contributions of this paper include

- 59 • A novel approach to analyzing distributed methods that combine data and model parallelism
60 by operating with sparse submodels for a quadratic model.
- 61 • The first analysis of independent subnetwork training in homogeneous and heterogeneous
62 scenarios without restrictive assumptions on gradient estimators.
- 63 • Identification of settings when IST can optimize very efficiently or converge not to the
64 optimal solution but only to an irreducible neighborhood which is also tightly characterized.
- 65 • Experimental validation of the proposed theory through carefully designed illustrative
66 experiments. Due to space limitations, the results (and proofs) are provided in the Appendix.

67 2 Formalism and Setup

68 We consider the standard optimization formulation of distributed/federated learning problem [41],

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

69 where n is the number of clients/workers, each $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$ represents the loss of the model
70 parameterized by vector $x \in \mathbb{R}^d$ on the data of client i .

71 A typical Stochastic Gradient Descent (SGD) type method for solving this problem has the form

$$x^{k+1} = x^k - \gamma g^k, \quad g^k = \frac{1}{n} \sum_{i=1}^n g_i^k, \quad (2)$$

72 where $\gamma > 0$ is a stepsize and g_i^k is a suitably constructed estimator of $\nabla f_i(x^k)$. In the distributed
 73 setting, computation of gradient estimators g_i^k is typically performed by clients, sent to the server,
 74 which subsequently performs aggregation via averaging $g^k = \frac{1}{n} \sum_{i=1}^n g_i^k$. The result is then used to
 75 update the model x^{k+1} via a gradient-type method (2), and at the next iteration the model is broadcast
 76 back to the clients. The process is repeated iteratively until a model of suitable qualities is found.

77 One of the main techniques used to accelerate distributed training is lossy *communication compression*
 78 [3, 27, 38]. It suggests applying a (possibly randomized) lossy compression mapping \mathcal{C} to a
 79 vector/matrix/tensor x before it is transmitted. This saves bits sent per every communication round
 80 at the cost of transmitting a less accurate estimate $\mathcal{C}(x)$ of x . The error caused by this routine also
 81 causes convergence issues, and to the best of our knowledge, convergence of IST-based techniques is
 82 for this reason not yet understood.

83 **Definition 1** (Unbiased compressor). *A randomized mapping $\mathcal{C} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is an **unbiased compression**
 84 **operator** ($\mathcal{C} \in \mathbb{U}(\omega)$ for brevity) if for some $\omega \geq 0$ and $\forall x \in \mathbb{R}^d$*

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

85 A notable example of a mapping from this class is the *random sparsification* (Rand-q for $q \in$
 86 $\{1, \dots, d\}$) operator defined by

$$\mathcal{C}_{\text{Rand-q}}(x) := \mathbf{C}_q x = \frac{d}{q} \sum_{i \in S} e_i e_i^\top x, \quad (4)$$

87 where $e_1, \dots, e_d \in \mathbb{R}^d$ are standard unit basis vectors in \mathbb{R}^d , and S is a random subset of $[d] :=$
 88 $\{1, \dots, d\}$ sampled from the uniform distribution on the all subsets of $[d]$ with cardinality q . Rand-q
 89 belongs to $\mathbb{U}(d/q - 1)$, which means that the more elements are “dropped” (lower q), the higher is
 90 the variance ω of the compressor.

91 In this work, we are mainly interested in a somewhat more general class of operators than mere
 92 sparsifiers. In particular, we are interested in compressing via the application of random matrices, i.e.,
 93 via *sketching*. A sketch $\mathbf{C}_i^k \in \mathbb{R}^{d \times d}$ can be used to represent submodel computations in the following
 94 way:

$$g_i^k := \mathbf{C}_i^k \nabla f_i(\mathbf{C}_i^k x^k), \quad (5)$$

95 where we require \mathbf{C}_i^k to be a symmetric positive semidefinite matrix. Such gradient estimate
 96 corresponds to computing the local gradient with respect to a sparse submodel model $\mathbf{C}_i^k x^k$, and
 97 additionally sketching the resulting gradient with the same matrix \mathbf{C}_i^k to guarantee that the resulting
 98 update lies in the lower-dimensional subspace.

99 Using this notion, Algorithm 1 (with one local gradient step) can be represented in the following form

$$x^{k+1} = \frac{1}{n} \sum_{i=1}^n [\mathbf{C}_i^k x^k - \gamma \mathbf{C}_i^k \nabla f_i(\mathbf{C}_i^k x^k)], \quad (6)$$

which is equivalent to the SGD-type update (2) when **perfect reconstruction** property holds

$$\mathbf{C}^k := \frac{1}{n} \sum_{i=1}^n \mathbf{C}_i^k = \mathbf{I},$$

100 where \mathbf{I} is the identity matrix (with probability one). This property holds for a specific class of
 101 compressors that are particularly useful for capturing the concept of an *independent* subnetwork
 102 partition.

103 **Definition 2** (Permutation sketch). *Assume that model size is greater than number of clients $d \geq n$
 104 and $d = qn$, where $q \geq 1$ is an integer¹. Let $\pi = (\pi_1, \dots, \pi_d)$ be a random permutation of $[d]$. Then
 105 for all $x \in \mathbb{R}^d$ and each $i \in [n]$ we define *Perm-q operator**

$$\mathbf{C}_i := n \cdot \sum_{j=q(i-1)+1}^{qi} e_{\pi_j} e_{\pi_j}^\top. \quad (7)$$

¹While this condition may look restrictive it naturally holds for distributed learning in a data-center setting. For other scenarios [40] generalized it for $n \geq d$ and block permutation case.

106 Perm-q is unbiased and can be conveniently used for representing (non-overlapping) structured
 107 decomposition of the model such that every client i is responsible for computations over a submodel
 108 $\mathbf{C}_i x^k$.

109 Our convergence analysis relies on assumption previously used for coordinate descent type methods.

110 **Assumption 1** (Matrix smoothness). *A differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is \mathbf{L} -smooth, if there
 111 exists a positive semi-definite matrix $\mathbf{L} \in \mathbb{R}^{d \times d}$ such that*

$$f(x+h) \leq f(x) + \langle \nabla f(x), h \rangle + \frac{1}{2} \langle \mathbf{L}h, h \rangle, \quad \forall x, h \in \mathbb{R}^d. \quad (8)$$

112 Standard L -smoothness condition is obtained as a special case of (8) for $\mathbf{L} = L \cdot \mathbf{I}$.

113 2.1 Issues with existing approaches

114 Consider the simplest gradient type method with compressed model in the single node setting

$$x^{k+1} = x^k - \gamma \nabla f(\mathcal{C}(x^k)). \quad (9)$$

115 Algorithms belonging to this family require a different analysis in comparison to SGD [16, 18],
 116 Distributed Compressed Gradient Descent [3, 26] and Randomized Coordinate Descent [34, 36] type
 117 methods because the gradient estimator is no longer unbiased

$$\mathbb{E} [\nabla f(\mathcal{C}(x))] \neq \nabla f(x) = \mathbb{E} [\mathcal{C}(\nabla f(x))]. \quad (10)$$

118 That is why such kind of algorithms are harder to analyze. So, prior results for *unbiased* SGD [25]
 119 can not be directly reused. Furthermore, the nature of the bias in this type of gradient estimator does
 120 not exhibit additive (zero-mean) noise, thereby preventing the application of previous analyses for
 121 biased SGD [1].

122 An assumption like bounded stochastic gradient norm extensively used in previous works [30, 48]
 123 hinders an accurate understanding of such methods. This assumption hides the fundamental difficulty
 124 of analyzing biased gradient estimator:

$$\mathbb{E} [\|\nabla f(\mathcal{C}(x))\|^2] \leq G \quad (11)$$

125 and may not hold even for quadratic functions $f(x) = x^\top \mathbf{A}x$. In addition, in the distributed
 126 setting such condition can result in vacuous bounds [23] as it does not allow to accurately capture
 127 heterogeneity.

128 3 Results in the Interpolation Case

129 To conduct a thorough theoretical analysis of methods that combine data with model parallelism,
 130 we simplify the algorithm and problem setting to isolate the unique effects of this approach. The
 131 following considerations are made:

- 132 (1) We assume that every node i computes the true gradient at the submodel $\mathbf{C}_i \nabla f_i(\mathbf{C}_i x^k)$.
- 133 (2) A notable difference from the original IST algorithm 1 is that workers perform single
 134 gradient descent step (or just gradient computation).
- 135 (3) Finally, we consider a special case of quadratic model (12) as a loss function (1).

136 Condition (1) is mainly for the sake of simplicity and clarity of exposition and can be potentially
 137 generalized to stochastic gradient computations. (2) is imposed because local steps did not bring
 138 any theoretical efficiency improvements for heterogeneous settings until very recently [32]. And
 139 even then, only with the introduction of additional control variables, which goes against resource-
 140 constrained device setting. The reason behind (3) is that despite the seeming simplicity quadratic
 141 problem has been used extensively to study properties of neural networks [46, 49]. Moreover, it is a
 142 non-trivial model which allows to understand complex optimization algorithms [4, 10, 17]. It serves
 143 as a suitable problem for observing complex phenomena and providing theoretical insights, which
 144 can also be observed in practical scenarios.

145 Having said that we consider a special case of problem (1)

$$f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x), \quad f_i(x) \equiv \frac{1}{2} x^\top \mathbf{L}_i x - \mathbf{b}_i^\top x. \quad (12)$$

146 In this case, $f(x)$ is $\bar{\mathbf{L}}$ -smooth, and $\nabla f(x) = \bar{\mathbf{L}}x - \bar{\mathbf{b}}$, where $\bar{\mathbf{L}} = \frac{1}{n} \sum_{i=1}^n \mathbf{L}_i$ and $\bar{\mathbf{b}} := \frac{1}{n} \sum_{i=1}^n \mathbf{b}_i$.

147 3.1 No linear term: problems and solutions

148 First, let us examine the case of $\mathbf{b}_i \equiv 0$, which we call interpolation for quadratics, and perform the
149 analysis for general sketches \mathbf{C}_i^k . In this case the gradient estimator (2) takes the form

$$g^k = \frac{1}{n} \sum_{i=1}^n \mathbf{C}_i^k \nabla f_i(\mathbf{C}_i^k x^k) = \frac{1}{n} \sum_{i=1}^n \mathbf{C}_i^k \mathbf{L}_i \mathbf{C}_i^k x^k = \bar{\mathbf{B}}^k x^k \quad (13)$$

150 where $\bar{\mathbf{B}}^k := \frac{1}{n} \sum_{i=1}^n \mathbf{C}_i^k \mathbf{L}_i \mathbf{C}_i^k$. We prove the following result for a method with such an estimator.

151 **Theorem 1.** Consider the method (2) with estimator (13) for a quadratic problem (12) with $\bar{\mathbf{L}} \succ 0$
152 and $\mathbf{b}_i \equiv 0$. Then if $\bar{\mathbf{W}} := \frac{1}{2} \mathbb{E} [\bar{\mathbf{L}} \bar{\mathbf{B}}^k + \bar{\mathbf{B}}^k \bar{\mathbf{L}}] \succeq 0$ and there exists constant $\theta > 0$:

$$\mathbb{E} [\bar{\mathbf{B}}^k \bar{\mathbf{L}} \bar{\mathbf{B}}^k] \preceq \theta \bar{\mathbf{W}}, \quad (14)$$

153 and the step size is chosen as $0 < \gamma \leq \frac{1}{\theta}$, the iterates satisfy

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} [\|\nabla f(x^k)\|_{\bar{\mathbf{L}}^{-1}}^2 \bar{\mathbf{W}} \bar{\mathbf{L}}^{-1}] \leq \frac{2(f(x^0) - \mathbb{E}[f(x^K)])}{\gamma K}, \quad (15)$$

154 and

$$\mathbb{E} [\|x^k - x^*\|_{\bar{\mathbf{L}}}^2] \leq \left(1 - \gamma \lambda_{\min} \left(\bar{\mathbf{L}}^{-\frac{1}{2}} \bar{\mathbf{W}} \bar{\mathbf{L}}^{-\frac{1}{2}}\right)\right)^k \|x^0 - x^*\|_{\bar{\mathbf{L}}}^2. \quad (16)$$

155 This theorem establishes an $\mathcal{O}(1/K)$ convergence rate with constant step size up to a stationary point
156 and linear convergence for the expected distance to the optimum. Note that we employ weighted
157 norms in our analysis, as the considered class of loss functions satisfies the matrix $\bar{\mathbf{L}}$ -smoothness
158 Assumption 1. The use of standard Euclidean distance may result in loose bounds that do not recover
159 correct rates for special cases like Gradient Descent.

160 It is important to highlight that inequality (14) may not hold (for any $\theta > 0$) in the general case
161 as the matrix $\bar{\mathbf{W}}$ is not guaranteed to be positive (semi-)definite in the case of general sampling.
162 The intuition behind it is that arbitrary sketches \mathbf{C}_i^k can result in gradient estimator g^k , which is
163 misaligned with the true gradient $\nabla f(x^k)$. Specifically, the inner product $\langle \nabla f(x^k), g^k \rangle$ can be
164 negative, and there is no expected descent after one step.

165 Next, we give examples of samplings for which the inequality (14) can be satisfied.

166 **1. Identity.** Consider $\mathbf{C}_i \equiv \mathbf{I}$. Then $\bar{\mathbf{B}}^k = \bar{\mathbf{L}}$, $\bar{\mathbf{B}}^k \bar{\mathbf{L}} \bar{\mathbf{B}}^k = \bar{\mathbf{L}}^3$, $\bar{\mathbf{W}} = \bar{\mathbf{L}}^2 \succ 0$ and hence (14) is
167 satisfied for $\theta = \lambda_{\max}(\bar{\mathbf{L}})$. So, (15) says that if we choose $\gamma = \frac{1}{\theta}$, then

$$\frac{1}{K} \sum_{k=0}^{K-1} \|\nabla f(x^k)\|_{\mathbf{I}}^2 \leq \frac{2\lambda_{\max}(\bar{\mathbf{L}})(f(x^0) - f(x^K))}{K},$$

168 which exactly matches the rate of Gradient Descent in the non-convex setting. As for iterates
169 convergence, the rate in (16) is $\lambda_{\max}(\bar{\mathbf{L}})/\lambda_{\min}(\bar{\mathbf{L}})$ corresponding to precise Gradient Descent result for
170 strongly convex functions.

2. Permutation. Assume $n = d^2$ and the use of Perm-1 (special case of Definition 2) sketch
 $\mathbf{C}_i^k = n e_{\pi_i^k} e_{\pi_i^k}^\top$, where $\pi^k = (\pi_1^k, \dots, \pi_n^k)$ is a random permutation of $[n]$. Then

$$\mathbb{E} [\bar{\mathbf{B}}^k] = \frac{1}{n} \sum_{i=1}^n n^2 \mathbb{E} [\mathbf{C}_i^k \mathbf{L}_i \mathbf{C}_i^k] = \frac{1}{n} \sum_{i=1}^n n \text{Diag}(\mathbf{L}_i) = \sum_{i=1}^n \mathbf{D}_i = n \bar{\mathbf{D}},$$

²This is done mainly for simplifying the presentation. Results can be generalized to the case of $n \neq d$ in the similar way as done in [40] which can be found in the Appendix.

171 where $\bar{\mathbf{D}} := \frac{1}{n} \sum_{i=1}^n \mathbf{D}_i$, $\mathbf{D}_i := \text{Diag}(\mathbf{L}_i)$. Then inequality (14) leads to

$$n \bar{\mathbf{D}} \bar{\mathbf{L}} \bar{\mathbf{D}} \preceq \frac{\theta}{2} (\bar{\mathbf{L}} \bar{\mathbf{D}} + \bar{\mathbf{D}} \bar{\mathbf{L}}), \quad (17)$$

172 which may not always hold as $\bar{\mathbf{L}} \bar{\mathbf{D}} + \bar{\mathbf{D}} \bar{\mathbf{L}}$ is not guaranteed to be positive definite even in case of
 173 $\bar{\mathbf{L}} \succ 0$. However, such kind of condition can be enforced via a slight modification of permutation
 174 sketches $\{\tilde{\mathbf{C}}_i\}_{i=1}^n$, which is done in Section 3.1.2. The limitation of such an approach is that
 175 compressors $\tilde{\mathbf{C}}_i$ become no longer unbiased.

176 **Remark 1.** Matrix $\bar{\mathbf{W}}$ in case of permutation sketches may not be positive-definite. Consider the
 177 following homogeneous ($\mathbf{L}_i \equiv \mathbf{L}$) two-dimensional problem example

$$\mathbf{L} = \begin{bmatrix} a & c \\ c & b \end{bmatrix}. \quad (18)$$

178 Then

$$\bar{\mathbf{W}} = \frac{1}{2} [\bar{\mathbf{L}} \bar{\mathbf{D}} + \bar{\mathbf{D}} \bar{\mathbf{L}}] = \begin{bmatrix} a^2 & c(a+b)/2 \\ c(a+b)/2 & b^2 \end{bmatrix}, \quad (19)$$

179 which for $c > \frac{2ab}{a+b}$ has $\det(\bar{\mathbf{W}}) < 0$, and thus $\bar{\mathbf{W}} \not\succeq 0$ according to Sylvester's criterion.

180 Next, we focus on the particular case of **Permutation** sketches, which are the most suitable for
 181 model partitioning according to Independent Subnetwork Training (IST). At the rest of the section,
 182 we discuss how the condition (14) can be enforced via a specially designed preconditioning of the
 183 problem (12) or modification of sketch mechanism (7).

184 3.1.1 Homogeneous problem preconditioning

185 To start consider a homogeneous setting $f_i(x) = \frac{1}{2} x^\top \mathbf{L} x$, so $\mathbf{L}_i \equiv \mathbf{L}$. Now define $\mathbf{D} = \text{Diag}(\mathbf{L}) -$
 186 diagonal matrix with elements equal to diagonal of \mathbf{L} . Then problem can be converted to

$$f_i(\mathbf{D}^{-\frac{1}{2}} x) = \frac{1}{2} (\mathbf{D}^{-\frac{1}{2}} x)^\top \mathbf{L} (\mathbf{D}^{-\frac{1}{2}} x) = \frac{1}{2} x^\top \underbrace{(\mathbf{D}^{-\frac{1}{2}} \mathbf{L} \mathbf{D}^{-\frac{1}{2}})}_{\tilde{\mathbf{L}}} x, \quad (20)$$

187 which is equivalent to the original problem after a change of variables $\tilde{x} := \mathbf{D}^{-\frac{1}{2}} x$. Note that
 188 $\mathbf{D} = \text{Diag}(\mathbf{L})$ is positive definite as $\mathbf{L} \succ 0$, and therefore $\tilde{\mathbf{L}} \succ 0$. Moreover, the preconditioned
 189 matrix $\tilde{\mathbf{L}}$ has all ones on the diagonal: $\text{Diag}(\tilde{\mathbf{L}}) = \mathbf{I}$. If we now combine it with Perm-1 sketches

$$\mathbb{E} [\bar{\mathbf{B}}^k] = \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \mathbf{C}_i \tilde{\mathbf{L}} \mathbf{C}_i \right] = n \text{Diag}(\tilde{\mathbf{L}}) = n \mathbf{I}.$$

190 Therefore, inequality (14) takes the form $\bar{\mathbf{W}} = n \tilde{\mathbf{L}} \succeq \frac{1}{\theta} n^2 \tilde{\mathbf{L}}$, which holds for $\theta \geq n$, and left hand
 191 side of (15) can be transformed the following way

$$\|\nabla f(x^k)\|_{\tilde{\mathbf{L}}^{-1} \bar{\mathbf{W}} \tilde{\mathbf{L}}^{-1}}^2 \geq n \lambda_{\min}(\tilde{\mathbf{L}}^{-1}) \|\nabla f(x^k)\|_{\mathbf{I}}^2 = n \lambda_{\max}(\tilde{\mathbf{L}}) \|\nabla f(x^k)\|_{\mathbf{I}}^2 \quad (21)$$

192 for an accurate comparison to standard methods. The resulting convergence guarantee

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \left[\|\nabla f(x^k)\|_{\mathbf{I}}^2 \right] \leq \frac{2 \lambda_{\max}(\tilde{\mathbf{L}}) (f(x^0) - \mathbb{E}[f(x^K)])}{K}, \quad (22)$$

193 which matches classical Gradient Descent.

194 3.1.2 Heterogeneous sketch preconditioning

195 In contrast to homogeneous case the heterogeneous problem $f_i(x) = \frac{1}{2} x^\top \mathbf{L}_i x$ can not be so easily
 196 preconditioned by a simple change of variables $\tilde{x} := \mathbf{D}^{-\frac{1}{2}} x$, as every client i has its own matrix
 197 \mathbf{L}_i . However, this problem can be fixed via the following modification of Perm-1, which scales the
 198 output according to the diagonal elements of local smoothness matrix \mathbf{L}_i :

$$\tilde{\mathbf{C}}_i := \sqrt{n} \left[\mathbf{L}_i^{-\frac{1}{2}} \right]_{\pi_i, \pi_i} e_{\pi_i} e_{\pi_i}^\top. \quad (23)$$

199 In this case $\mathbb{E} [\tilde{\mathbf{C}}_i \mathbf{L}_i \tilde{\mathbf{C}}_i] = \mathbf{I}$, $\mathbb{E} [\bar{\mathbf{B}}^k] = \mathbf{I}$, and $\bar{\mathbf{W}} = \bar{\mathbf{L}}$. Then inequality (14) is satisfied for $\theta \geq 1$.

200 If one plugs these results into (15), such convergence guarantee can be obtained

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} [\|\nabla f(x^k)\|_{\mathbf{I}}^2] \leq \frac{2\lambda_{\max}(\bar{\mathbf{L}})(f(x^0) - \mathbb{E}[f(x^K)])}{K}, \quad (24)$$

201 which matches the Gradient Descent result as well. Thus we can conclude that heterogeneity does not
 202 bring such a fundamental challenge in this scenario. In addition, a method with Perm-1 is significantly
 203 better in terms of computational and communication complexity as it requires calculating the local
 204 gradients with respect to much smaller submodels and transmits only sparse updates.

205 This construction also shows that for $\gamma = 1/\theta = 1$

$$\gamma \lambda_{\min} \left(\bar{\mathbf{L}}^{-\frac{1}{2}} \bar{\mathbf{W}} \bar{\mathbf{L}}^{-\frac{1}{2}} \right) = \lambda_{\min} \left(\bar{\mathbf{L}}^{-\frac{1}{2}} \bar{\mathbf{L}} \bar{\mathbf{L}}^{-\frac{1}{2}} \right) = 1, \quad (25)$$

206 which after plugging into the bound for the iterates (16) shows that the method basically converges in
 207 1 iteration. This observation that sketch preconditioning can be extremely efficient, although it uses
 208 only the diagonal elements of matrices \mathbf{L}_i .

209 Now when we understand that the method can perform very well in the special case of $\tilde{\mathbf{b}}_i \equiv 0$ we can
 210 move on to a more complicated situation.

211 4 Irreducible Bias in the General Case

212 Now we look at the most general heterogeneous case with different matrices and linear terms
 213 $f_i(x) \equiv \frac{1}{2} x^\top \mathbf{L}_i x - x^\top \mathbf{b}_i$. In this instance gradient estimator (2) takes the form

$$g^k = \frac{1}{n} \sum_{i=1}^n \mathbf{C}_i^k \nabla f_i(\mathbf{C}_i^k x^k) = \frac{1}{n} \sum_{i=1}^n \mathbf{C}_i^k (\mathbf{L}_i \mathbf{C}_i^k x^k - \mathbf{b}_i) = \bar{\mathbf{B}}^k x^k - \bar{\mathbf{C}}\mathbf{b}, \quad (26)$$

214 where $\bar{\mathbf{C}}\mathbf{b} = \frac{1}{n} \sum_{i=1}^n \mathbf{C}_i^k \mathbf{b}_i$. Herewith let us use a heterogeneous permutation sketch preconditioner
 215 (23) like in Section 3.1.2 Then $\mathbb{E} [\bar{\mathbf{B}}^k] = \mathbf{I}$ and $\mathbb{E} [\bar{\mathbf{C}}\mathbf{b}] = \frac{1}{\sqrt{n}} \bar{\mathbf{D}}\bar{\mathbf{b}}$, where $\bar{\mathbf{D}}\bar{\mathbf{b}} := \frac{1}{n} \sum_{i=1}^n \mathbf{D}_i^{-\frac{1}{2}} \mathbf{b}_i$.

216 Furthermore expected gradient estimator (26) results in $\mathbb{E} [g^k] = x^k - \frac{1}{\sqrt{n}} \bar{\mathbf{D}}\bar{\mathbf{b}}$ and can be transformed
 217 the following way

$$\mathbb{E} [g^k] = \bar{\mathbf{L}}^{-1} \bar{\mathbf{L}} x^k \pm \bar{\mathbf{L}}^{-1} \bar{\mathbf{b}} - \frac{1}{\sqrt{n}} \bar{\mathbf{D}}\bar{\mathbf{b}} = \bar{\mathbf{L}}^{-1} \nabla f(x^k) + \underbrace{\bar{\mathbf{L}}^{-1} \bar{\mathbf{b}} - \frac{1}{\sqrt{n}} \bar{\mathbf{D}}\bar{\mathbf{b}}}_h, \quad (27)$$

218 which reflects the decomposition of the estimator into optimally preconditioned true gradient and a
 219 bias, depending on the linear terms \mathbf{b}_i .

220 4.1 Bias of the method

221 Estimator (27) can be directly plugged (with proper conditioning) into general SGD update (2)

$$\mathbb{E} [x^{k+1}] = x^k - \gamma \mathbb{E} [g^k] = (1 - \gamma)x^k + \frac{\gamma}{\sqrt{n}} \bar{\mathbf{D}}\bar{\mathbf{b}} = (1 - \gamma)^{k+1} x^0 + \frac{\gamma}{\sqrt{n}} \bar{\mathbf{D}}\bar{\mathbf{b}} \sum_{j=0}^k (1 - \gamma)^j. \quad (28)$$

222 The resulting recursion (28) is exact, and its asymptotic limit can be analyzed. Thus for constant
 223 $\gamma < 1$ by using the formula for the sum of the first k terms of a geometric series, one gets

$$\mathbb{E} [x^k] = (1 - \gamma)^k x^0 + \frac{1 - (1 - \gamma)^k}{\sqrt{n}} \bar{\mathbf{D}}\bar{\mathbf{b}} \xrightarrow[k \rightarrow \infty]{} \frac{1}{\sqrt{n}} \bar{\mathbf{D}}\bar{\mathbf{b}},$$

which shows that in the limit, the first initialization term (with x^0) vanishes while the second converges
 to $\frac{1}{\sqrt{n}} \bar{\mathbf{D}}\bar{\mathbf{b}}$. This reasoning shows that the method does not converge to the exact solution

$$x^k \rightarrow x^\infty \neq x^* \in \arg \min_{x \in \mathbb{R}^d} \left\{ \frac{1}{2} x^\top \bar{\mathbf{L}} x - x^\top \bar{\mathbf{b}} \right\},$$

224 which for the positive-definite $\bar{\mathbf{L}}$ can be defined as $x^* = \bar{\mathbf{L}}^{-1} \bar{\mathbf{b}}$, while $x^\infty = \frac{1}{n\sqrt{n}} \sum_{i=1}^n \mathbf{D}_i^{-\frac{1}{2}} \mathbf{b}_i$. So,
 225 in general, there is an unavoidable bias. However, in the limit case: $n = d \rightarrow \infty$, the bias diminishes.

226 **4.2 Generic convergence analysis**

227 While the analysis in Section 4.1 is precise, it does not allow us to compare the convergence of IST
 228 to standard optimization methods. Due to this, we also analyze the non-asymptotic behavior of the
 229 method to understand the convergence speed. Our result is formalized in the following theorem.

230 **Theorem 2.** Consider the method (2) with estimator (26) for a quadratic problem (12) with the
 231 positive definite matrix $\bar{\mathbf{L}} \succ 0$. Assume that for every $\mathbf{D}_i := \text{Diag}(\mathbf{L}_i)$ matrices $\mathbf{D}_i^{-\frac{1}{2}}$ exist, scaled
 232 permutation sketches (23) are used and heterogeneity is bounded as $\mathbb{E} \left[\|g^k - \mathbb{E}[g^k]\|_{\bar{\mathbf{L}}}^2 \right] \leq \sigma^2$.
 233 Then for step size is chosen as

$$0 < \gamma \leq \gamma_{c,\beta} := \frac{1/2-\beta}{\beta+1/2}, \quad (29)$$

234 where $\gamma_{c,\beta} \in (0, 1]$ for $\beta \in (0, 1/2)$, the iterates satisfy

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \left[\|\nabla f(x^k)\|_{\bar{\mathbf{L}}^{-1}}^2 \right] \leq \frac{2(f(x^0) - \mathbb{E}[f(x^K)])}{\gamma K} + (2\beta^{-1}(1-\gamma) + \gamma) \|h\|_{\bar{\mathbf{L}}}^2 + \gamma \sigma^2, \quad (30)$$

235 where $\bar{\mathbf{L}} = \frac{1}{n} \sum_{i=1}^n \mathbf{L}_i$, $h = \bar{\mathbf{L}}^{-1} \bar{\mathbf{b}} - \frac{1}{\sqrt{n}} \frac{1}{n} \sum_{i=1}^n \mathbf{D}_i^{-\frac{1}{2}} \mathbf{b}_i$ and $\bar{\mathbf{b}} = \frac{1}{n} \sum_{i=1}^n \mathbf{b}_i$.

236 Note that the derived convergence upper bound has a neighborhood proportional to the bias of
 237 the gradient estimator h and level of heterogeneity σ^2 . Some of these terms with factor γ can be
 238 eliminated via decreasing learning rate schedule (e.g., $\sim 1/\sqrt{k}$). However, such a strategy does not
 239 diminish the term with a multiplier $2\beta^{-1}(1-\gamma)$, making the neighborhood irreducible. Moreover,
 240 this term can be eliminated for $\gamma = 1$, which also minimizes the first term that decreases as $1/K$.
 241 Though, such step size choice maximizes the terms with factor γ . Furthermore, there exists an
 242 inherent trade-off between convergence speed and the size of the neighborhood.

243 In addition, convergence to the stationary point is measured in the weighted by $\bar{\mathbf{L}}^{-1}$ squared norm of
 244 the gradient. At the same time, the neighborhood term depends on the weighted by $\bar{\mathbf{L}}$ norm of h . This
 245 fine-grained decoupling is achieved by carefully applying Fenchel-Young inequality and provides a
 246 tighter characterization of the convergence compared to using standard Euclidean distances.

247 **Homogeneous case.** In this scenario, every worker has access to the all data $f_i(x) \equiv \frac{1}{2}x^\top \mathbf{L}x - x^\top \mathbf{b}$.
 248 Then diagonal preconditioning of the problem can be used as in the previous Section 3.1.1. This
 249 results in a gradient $\nabla f(x) = \tilde{\mathbf{L}}x - \tilde{\mathbf{b}}$ for $\tilde{\mathbf{L}} = \mathbf{D}^{-\frac{1}{2}} \mathbf{L} \mathbf{D}^{-\frac{1}{2}}$ and $\tilde{\mathbf{b}} = \mathbf{D}^{-\frac{1}{2}} \mathbf{b}$. If it is further combined
 250 with a scaled by $1/\sqrt{n}$ Permutation sketch $\mathbf{C}_i := \sqrt{n} e_{\pi_i} e_{\pi_i}^\top$, the resulting gradient estimator is

$$g^k = x^k - \frac{1}{\sqrt{n}} \tilde{\mathbf{b}} = \tilde{\mathbf{L}}^{-1} \nabla f(x^k) + \tilde{h}, \quad (31)$$

251 for $\tilde{h} = \tilde{\mathbf{L}}^{-1} \tilde{\mathbf{b}} - \frac{1}{\sqrt{n}} \tilde{\mathbf{b}}$. In this case heterogeneity term σ^2 from upper bound (30) disappears
 252 as $\mathbb{E} \left[\|g^k - \mathbb{E}[g^k]\|_{\bar{\mathbf{L}}}^2 \right] = 0$, thus the neighborhood size can significantly decrease. However,
 253 the bias term depending on \tilde{h} still remains as the method does not converge to the exact solution
 254 $x^k \rightarrow x^\infty \neq x^* = \tilde{\mathbf{L}}^{-1} \tilde{\mathbf{b}}$ for positive-definite $\tilde{\mathbf{L}}$. Nevertheless the method's fixed point $x^\infty = \tilde{\mathbf{b}}/\sqrt{n}$
 255 and solution x^* can coincide when $\tilde{\mathbf{L}}^{-1} \tilde{\mathbf{b}} = \frac{1}{\sqrt{n}} \tilde{\mathbf{b}}$, which means that $\tilde{\mathbf{b}}$ is the right eigenvector of
 256 matrix $\tilde{\mathbf{L}}^{-1}$ with eigenvalue $\frac{1}{\sqrt{n}}$.

257 Let us contrast obtained result (30) with non-convex rate of SGD [25] with constant step size γ for
 258 L -smooth and lower-bounded f

$$\min_{k \in \{0, \dots, K-1\}} \|\nabla f(x^k)\|^2 \leq \frac{6(f(x^0) - \inf f)}{\gamma K} + \gamma LC, \quad (32)$$

259 where constant C depends, for example, on the variance of stochastic gradient estimates. Observe
 260 that the first term in the compared upper bounds (32) and (30) is almost identical and decreases with
 261 speed $1/K$. But unlike (30) the neighborhood for SGD can be completely eliminated by reducing the
 262 step size γ . This highlights a fundamental difference of our results to unbiased methods.

263 The intuition behind this issue is that for SGD-type methods like Compressed Gradient Descent

$$x^{k+1} = x^k - \mathcal{C}(\nabla f(x^k)) \quad (33)$$

264 the gradient estimate is unbiased and enjoys the property that variance

$$\mathbb{E} [\|\mathcal{C}(\nabla f(x^k)) - \nabla f(x^k)\|^2] \leq \omega \|\nabla f(x^k)\|^2 \quad (34)$$

265 goes down to zero as the method progresses because $\nabla f(x^k) \rightarrow \nabla f(x^*) = 0$ in the unconstrained
266 case. In addition, any stationary point x^* ceases to be a fixed point of the iterative procedure as

$$x^* \neq x^* - \nabla f(\mathcal{C}(x^*)), \quad (35)$$

267 in the general case, unlike for Compressed Gradient Descent with both biased and unbiased compres-
268 sors \mathcal{C} . So, even if the method (computing gradient at sparse model) is initialized from the *solution*
269 after one gradient step, it may get away from there.

270 4.3 Comparison to previous works

271 **Independent Subnetwork Training [45].** There are several improvements over the previous works
272 that tried to theoretically analyze the convergence of Distributed IST.

273 The first difference is that our results allow for an almost arbitrary level of model sparsification,
274 i.e., work for any $\omega \geq 0$ as permutation sketches can be viewed as a special case of compression
275 operators (1). This improves significantly over the work of [45], which demands³ $\omega \lesssim \mu^2/L^2$. Such a
276 requirement is very restrictive as the condition number L/μ of the loss function f is typically very
277 large for any non-trivial optimization problem. Thus, the sparsifier’s (4) variance $\omega = d/q - 1$ has to
278 be very close to 0 and $q \approx d$. So, the previous theory allows almost no compression (sparsification)
279 because it is based on the analysis of Gradient Descent with Compressed Iterates [24].

280 The second distinction is that the original IST work [45] considered a single node setting and thus
281 their convergence bounds did not capture the effect of heterogeneity, which we believe is of crucial
282 importance for distributed setting [9, 39]. Besides, they consider Lipschitz continuity of the loss
283 function f , which is not satisfied for a simple quadratic model. A more detailed comparison including
284 additional assumptions on the gradient estimator made in [45] is presented in the Appendix.

285 **FL with Model Pruning.** In a recent work [48] made an attempt to analyze a variant of the FedAvg
286 algorithm with sparse local initialization and compressed gradient training (pruned local models).
287 They considered a case of L -smooth loss and sparsification operator satisfying a similar condition to
288 (1). However, they also assumed that the squared norm of stochastic gradient is uniformly bounded
289 (11), which is “pathological” [23] especially in the case of local methods as it does not allow to
290 capture the very important effect of heterogeneity and can result in vacuous bounds.

291 In the Appendix we show some limitations of other relevant previous approaches to training with
292 compressed models: too restrictive assumptions on the algorithm [33] or not applicability in our
293 problem setting [7].

294 5 Conclusions and Future Work

295 In this study, we introduced a novel approach to understanding training with combined model and
296 data parallelism for a quadratic model. This framework allowed to shed light on distributed submodel
297 optimization which revealed the advantages and limitations Independent Subnetwork Training (IST).
298 Moreover, we accurately characterized the behavior of the considered method in both homogeneous
299 and heterogeneous scenarios without imposing restrictive assumptions on gradient estimators.

300 In future research, it would be valuable to explore extensions of our findings to settings that are closer
301 to practical scenarios, such as cross-device federated learning. This could involve investigating partial
302 participation support, leveraging local training benefits, and ensuring robustness against stragglers.
303 Additionally, it would be interesting to generalize our results to non-quadratic scenarios without
304 relying on pathological assumptions.

³ μ refers to constant from Polyak-Łojasiewicz (or strong convexity) condition. In case of a quadratic problem with positive-definite matrix \mathbf{A} : $\mu = \lambda_{\min}(\mathbf{A})$

References

- 305
- 306 [1] Ahmad Ajalloeian and Sebastian U Stich. On the convergence of SGD with biased gradients.
307 *arXiv preprint arXiv:2008.00051*, 2020.
- 308 [2] Samiul Alam, Luyang Liu, Ming Yan, and Mi Zhang. FedRolex: Model-heterogeneous federated
309 learning with rolling sub-model extraction. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave,
310 and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022.
- 311 [3] Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. QSGD:
312 Communication-efficient SGD via gradient quantization and encoding. In *Advances in Neural
313 Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- 314 [4] Yossi Arjevani, Ohad Shamir, and Nathan Srebro. A tight convergence analysis for stochastic
315 gradient descent with delayed updates. In *Algorithmic Learning Theory*, pages 111–132. PMLR,
316 2020.
- 317 [5] Sebastian Caldas, Jakub Konečný, H Brendan McMahan, and Ameet Talwalkar. Expanding
318 the reach of federated learning by reducing client resource requirements. *arXiv preprint
319 arXiv:1812.07210*, 2018.
- 320 [6] Zachary Charles, Kallista Bonawitz, Stanislav Chiknavaryan, Brendan McMahan, et al. Fed-
321 erated select: A primitive for communication-and memory-efficient federated learning. *arXiv
322 preprint arXiv:2208.09432*, 2022.
- 323 [7] El Mahdi Chayti and Sai Praneeth Karimireddy. Optimization with access to auxiliary informa-
324 tion. *arXiv preprint arXiv:2206.00395*, 2022.
- 325 [8] Yuanyuan Chen, Zichen Chen, Pengcheng Wu, and Han Yu. Fedobd: Opportunistic block
326 dropout for efficiently training large-scale neural networks through federated learning. *arXiv
327 preprint arXiv:2208.05174*, 2022.
- 328 [9] Sélím Chraïbi, Ahmed Khaled, Dmitry Kovalev, Peter Richtárik, Adil Salim, and Martin Takáč.
329 Distributed fixed point methods with compressed iterates. *arXiv preprint arXiv:2102.07245*,
330 2019.
- 331 [10] Leonardo Cunha, Gauthier Gidel, Fabian Pedregosa, Damien Scieur, and Courtney Paque-
332 tte. Only tails matter: Average-case universality and robustness in the convex regime. In
333 *International Conference on Machine Learning*, pages 4474–4491. PMLR, 2022.
- 334 [11] Jeffrey Dean, Greg Corrado, Rajat Monga, Kai Chen, Matthieu Devin, Mark Mao, Marc’ aurelio
335 Ranzato, Andrew Senior, Paul Tucker, Ke Yang, et al. Large scale distributed deep networks.
336 *Advances in neural information processing systems*, 25, 2012.
- 337 [12] Enmao Diao, Jie Ding, and Vahid Tarokh. Heterofl: Computation and communication efficient
338 federated learning for heterogeneous clients. *arXiv preprint arXiv:2010.01264*, 2020.
- 339 [13] Chen Dun, Mirian Hipolito, Chris Jermaine, Dimitrios Dimitriadis, and Anastasios Kyrillidis.
340 Efficient and light-weight federated learning via asynchronous distributed dropout. *arXiv
341 preprint arXiv:2210.16105*, 2022.
- 342 [14] Chen Dun, Cameron R Wolfe, Christopher M Jermaine, and Anastasios Kyrillidis. ResIST:
343 Layer-wise decomposition of resnets for distributed training. In *Uncertainty in Artificial
344 Intelligence*, pages 610–620. PMLR, 2022.
- 345 [15] Philipp Farber and Krste Asanovic. Parallel neural network training on multi-spert. In *Proceed-
346 ings of 3rd International Conference on Algorithms and Architectures for Parallel Processing*,
347 pages 659–666. IEEE, 1997.
- 348 [16] Eduard Gorbunov, Filip Hanzely, and Peter Richtárik. A unified theory of SGD: Variance
349 reduction, sampling, quantization and coordinate descent. In *International Conference on
350 Artificial Intelligence and Statistics*, pages 680–690. PMLR, 2020.

- 351 [17] Baptiste Goujaud, Damien Scieur, Aymeric Dieuleveut, Adrien B Taylor, and Fabian Pedregosa.
352 Super-acceleration with cyclical step-sizes. In *International Conference on Artificial Intelligence*
353 *and Statistics*, pages 3028–3065. PMLR, 2022.
- 354 [18] Robert Mansel Gower, Nicolas Loizou, Xun Qian, Alibek Sailanbayev, Egor Shulgin, and Peter
355 Richtárik. SGD: General analysis and improved rates. *Proceedings of the 36th International*
356 *Conference on Machine Learning, Long Beach, California*, 2019.
- 357 [19] Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola,
358 Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch SGD: Training
359 imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2018.
- 360 [20] Samuel Horvath, Stefanos Laskaridis, Mario Almeida, Ilias Leontiadis, Stylianos Venieris, and
361 Nicholas Lane. FjORD: Fair and accurate federated learning under heterogeneous targets with
362 ordered dropout. *Advances in Neural Information Processing Systems*, 34:12876–12889, 2021.
- 363 [21] Yuang Jiang, Shiqiang Wang, Victor Valls, Bong Jun Ko, Wei-Han Lee, Kin K Leung, and
364 Leandros Tassiulas. Model pruning enables efficient federated learning on edge devices. *IEEE*
365 *Transactions on Neural Networks and Learning Systems*, 2022.
- 366 [22] Peter Kairouz, H. Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Ar-
367 jun Nitin Bhagoji, Kallista A. Bonawitz, Zachary Charles, Graham Cormode, Rachel Cum-
368 mings, Rafael G. L. D’Oliveira, Hubert Eichner, Salim El Rouayheb, David Evans, Josh Gardner,
369 Zachary Garrett, Adrià Gascón, Badih Ghazi, Phillip B. Gibbons, Marco Gruteser, Zaïd Har-
370 chaoui, Chaoyang He, Lie He, Zhouyuan Huo, Ben Hutchinson, Justin Hsu, Martin Jaggi, Tara
371 Javidi, Gauri Joshi, Mikhail Khodak, Jakub Konečný, Aleksandra Korolova, Farinaz Koushanfar,
372 Sanmi Koyejo, Tancrède Lepoint, Yang Liu, Prateek Mittal, Mehryar Mohri, Richard Nock,
373 Ayfer Özgür, Rasmus Pagh, Hang Qi, Daniel Ramage, Ramesh Raskar, Mariana Raykova, Dawn
374 Song, Weikang Song, Sebastian U. Stich, Ziteng Sun, Ananda Theertha Suresh, Florian Tramèr,
375 Praneeth Vepakomma, Jianyu Wang, Li Xiong, Zheng Xu, Qiang Yang, Felix X. Yu, Han Yu,
376 and Sen Zhao. Advances and open problems in federated learning. *Found. Trends Mach. Learn.*,
377 14(1-2):1–210, 2021.
- 378 [23] Ahmed Khaled, Konstantin Mishchenko, and Peter Richtárik. Tighter theory for local SGD on
379 identical and heterogeneous data. In *International Conference on Artificial Intelligence and*
380 *Statistics*, pages 4519–4529. PMLR, 2020.
- 381 [24] Ahmed Khaled and Peter Richtárik. Gradient descent with compressed iterates. *arXiv preprint*
382 *arXiv:1909.04716*, 2019.
- 383 [25] Ahmed Khaled and Peter Richtárik. Better theory for SGD in the nonconvex world. *Transactions*
384 *on Machine Learning Research*, 2023. Survey Certification.
- 385 [26] Sarit Khirirat, Hamid Reza Feyzmahdavian, and Mikael Johansson. Distributed learning with
386 compressed gradients. *arXiv preprint arXiv:1806.06573*, 2018.
- 387 [27] Jakub Konečný, H. Brendan McMahan, Felix X. Yu, Peter Richtárik, Ananda Theertha Suresh,
388 and Dave Bacon. Federated learning: Strategies for improving communication efficiency. *NIPS*
389 *Private Multi-Party Machine Learning Workshop*, 2016.
- 390 [28] Fangshuo Liao and Anastasios Kyriillidis. On the convergence of shallow neural network training
391 with randomly masked neurons. *Transactions on Machine Learning Research*, 2022.
- 392 [29] Rongmei Lin, Yonghui Xiao, Tien-Ju Yang, Ding Zhao, Li Xiong, Giovanni Motta, and
393 Françoise Beaufays. Federated pruning: Improving neural network efficiency with federated
394 learning. *arXiv preprint arXiv:2209.06359*, 2022.
- 395 [30] Tao Lin, Sebastian U Stich, Luis Barba, Daniil Dmitriev, and Martin Jaggi. Dynamic model
396 pruning with feedback. In *International Conference on Learning Representations*, 2019.
- 397 [31] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas.
398 Communication-Efficient Learning of Deep Networks from Decentralized Data. In *Proceedings*
399 *of the 20th International Conference on Artificial Intelligence and Statistics*, volume 54 of
400 *Proceedings of Machine Learning Research*, pages 1273–1282, 20–22 Apr 2017.

- 401 [32] Konstantin Mishchenko, Grigory Malinovsky, Sebastian Stich, and Peter Richtárik. ProxSkip:
402 Yes! local gradient steps provably lead to communication acceleration! finally! In *International*
403 *Conference on Machine Learning*, pages 15750–15769. PMLR, 2022.
- 404 [33] Amirkeivan Mohtashami, Martin Jaggi, and Sebastian Stich. Masked training of neural networks
405 with partial gradients. In *International Conference on Artificial Intelligence and Statistics*,
406 pages 5876–5890. PMLR, 2022.
- 407 [34] Yu Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems.
408 *SIAM Journal on Optimization*, 22(2):341–362, 2012.
- 409 [35] Xinchi Qiu, Javier Fernandez-Marques, Pedro PB Gusmao, Yan Gao, Titouan Parcollet, and
410 Nicholas Donald Lane. ZeroFL: Efficient on-device training for federated learning with local
411 sparsity. In *International Conference on Learning Representations*, 2022.
- 412 [36] Peter Richtárik and Martin Takáč. Iteration complexity of randomized block-coordinate descent
413 methods for minimizing a composite function. *Mathematical Programming*, 144(1-2):1–38,
414 2014.
- 415 [37] Peter Richtárik and Martin Takáč. Distributed coordinate descent method for learning with big
416 data. *Journal of Machine Learning Research*, 17(75):1–25, 2016.
- 417 [38] Frank Seide, Hao Fu, Jasha Droppo, Gang Li, and Dong Yu. 1-bit stochastic gradient descent
418 and its application to data-parallel distributed training of speech dnns. In *Fifteenth Annual*
419 *Conference of the International Speech Communication Association*, 2014.
- 420 [39] Egor Shulgin and Peter Richtárik. Shifted compression framework: Generalizations and
421 improvements. In *The 38th Conference on Uncertainty in Artificial Intelligence*, 2022.
- 422 [40] Rafał Szlendak, Alexander Tyurin, and Peter Richtárik. Permutation compressors for prov-
423 ably faster distributed nonconvex optimization. In *International Conference on Learning*
424 *Representations*, 2022.
- 425 [41] Jianyu Wang, Zachary Charles, Zheng Xu, Gauri Joshi, H Brendan McMahan, Maruan Al-
426 Shedivat, Galen Andrew, Salman Avestimehr, Katharine Daly, Deepesh Data, et al. A field
427 guide to federated optimization. *arXiv preprint arXiv:2107.06917*, 2021.
- 428 [42] Dingzhu Wen, Ki-Jun Jeon, and Kaibin Huang. Federated dropout—a simple approach for
429 enabling federated learning on resource constrained devices. *IEEE Wireless Communications*
430 *Letters*, 11(5):923–927, 2022.
- 431 [43] Cameron R Wolfe, Jingkang Yang, Arindam Chowdhury, Chen Dun, Artun Bayer, Santiago
432 Segarra, and Anastasios Kyrillidis. Gist: Distributed training for large-scale graph convolutional
433 networks. *arXiv preprint arXiv:2102.10424*, 2021.
- 434 [44] Tien-Ju Yang, Dhruv Guliani, Françoise Beaufays, and Giovanni Motta. Partial variable training
435 for efficient on-device federated learning. In *ICASSP 2022-2022 IEEE International Conference*
436 *on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4348–4352. IEEE, 2022.
- 437 [45] Binhang Yuan, Cameron R Wolfe, Chen Dun, Yuxin Tang, Anastasios Kyrillidis, and Chris
438 Jermaine. Distributed learning of fully connected neural networks using independent subnet
439 training. *Proceedings of the VLDB Endowment*, 15(8):1581–1590, 2022.
- 440 [46] Guodong Zhang, Lala Li, Zachary Nado, James Martens, Sushant Sachdeva, George Dahl, Chris
441 Shallue, and Roger B Grosse. Which algorithmic choices matter at which batch sizes? insights
442 from a noisy quadratic model. *Advances in neural information processing systems*, 32, 2019.
- 443 [47] Xiru Zhang, Michael Mckenna, Jill Mesirov, and David Waltz. An efficient implementation
444 of the back-propagation algorithm on the connection machine cm-2. *Advances in neural*
445 *information processing systems*, 2, 1989.
- 446 [48] Hanhan Zhou, Tian Lan, Guru Venkataramani, and Wenbo Ding. On the convergence of
447 heterogeneous federated learning with arbitrary adaptive online model pruning. *arXiv preprint*
448 *arXiv:2201.11803*, 2022.

- 449 [49] Libin Zhu, Chaoyue Liu, Adityanarayanan Radhakrishnan, and Mikhail Belkin. Quadratic
450 models for understanding neural network dynamics. *arXiv preprint arXiv:2205.11787*, 2022.
- 451 [50] Martin Zinkevich, Markus Weimer, Lihong Li, and Alex Smola. Parallelized stochastic gradient
452 descent. *Advances in neural information processing systems*, 23, 2010.