
z -SignFedAvg: A Unified Sign-based Stochastic Compression for Federated Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Federated learning is a promising privacy-preserving distributed learning paradigm
2 but suffers from high communication cost when training large-scale machine
3 learning models. Sign-based methods, such as SignSGD [Bernstein et al., 2018],
4 have been proposed as a biased gradient compression technique for reducing
5 the communication cost. However, sign-based compression could diverge under
6 heterogeneous data, which motivate developments of advanced techniques, such
7 as the error-feedback method and stochastic sign-based compression, to fix this
8 issue. Nevertheless, these methods still suffer significantly slower convergence
9 rate than uncompressed algorithms. Besides, none of them allow local multiple
10 SGD updates like FedAvg [McMahan et al., 2017]. In this paper, we propose a
11 novel noisy perturbation scheme with a general symmetric noise distribution for
12 sign-based compression, which not only allows one to flexibly control the tradeoff
13 between gradient bias and convergence performance, but also provides a unified
14 viewpoint to existing sign-based methods. More importantly, we propose the very
15 first sign-based FedAvg algorithm (z -SignFedAvg). Theoretically, we show that
16 z -SignFedAvg achieves a faster convergence rate than existing sign-based methods
17 and, under the uniformly distributed noise, can even enjoy the same convergence
18 rate as its uncompressed counterpart. Extensive experiments are conducted to
19 demonstrate that our proposed z -SignFedAvg can achieve competitive empirical
20 performance on real datasets.

21 1 Introduction

22 In this paper, we consider the Federated Learning (FL) network with one parameter server and n
23 clients [McMahan et al., 2017, Li et al., 2020a], with the focus on solving the following distributed
24 learning problem

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

25 where $f_i(\cdot)$ is the local objective function for the i -th client, for $i = 1, \dots, n$. Throughout this
26 work, we assume that each f_i is smooth and possibly non-convex. The local objective functions are
27 generated from the local dataset owned by each client. When designing distributed algorithms to solve
28 (1), a crucial aspect is the communication efficiency since each client needs to transmit their local
29 gradients to the server frequently [Li et al., 2020a]. As one of the most popular FL algorithms, the
30 federated averaging (FedAvg) algorithm [McMahan et al., 2017, Konečný et al., 2016] considers local
31 multiple SGD updates with periodic communications to reduce the communication cost. Another
32 way is to compress the local gradients before sending them to the server [Li et al., 2020a, Alistarh
33 et al., 2017, Reisizadeh et al., 2020]. Among the existing compression methods, a simple yet elegant
34 technique is to take the sign of each coordinate of the local gradients, which requires only one bit for

35 transmitting each coordinate. For any $x \in \mathbb{R}$, we define the sign operator as: $\text{Sign}(x) = 1$ if $x \geq 0$
 36 and -1 otherwise.

37 Recently, optimization algorithms with sign-based compression have attracted much attention as they
 38 enjoy a great communication efficiency while still achieving comparable empirical performance as
 39 uncompressed algorithms [Bernstein et al., 2018, Karimireddy et al., 2019, Safaryan and Richtárik,
 40 2021]. However, for distributed learning, especially the scenarios with heterogeneous data, i.e.,
 41 $f_i \neq f_j$ for every $i \neq j$, a naive application of the sign-based algorithm cannot guarantee convergence
 42 [Karimireddy et al., 2019, Chen et al., 2020a, Safaryan and Richtárik, 2021]. There are mainly two
 43 approaches to fix this issue in the existing literature. The first one is the stochastic sign-based method,
 44 which introduces stochasticity into the sign operation [Jin et al., 2020, Safaryan and Richtárik, 2021,
 45 Chen et al., 2020a], and the second one is the Error-Feedback (EF) method [Karimireddy et al., 2019,
 46 Vogels et al., 2019, Tang et al., 2019]. However, these works are still unastatisfactory. Specifically, first,
 47 the theoretical convergence rates of these algorithms are still worse than uncompressed algorithms
 48 like [Ghadimi and Lan, 2013, Yu et al., 2019]. Second, none of them allows the clients to have
 49 multiple local SGD updates within one communication round like the FedAvg. This work aims at
 50 addressing these issues and closing the gaps for sign-based methods. A detailed review for related
 51 works is given in Appendix A.

52 **Main contributions.** Our contributions are summarized as follows.

- 53 (1) **A unified family of stochastic sign operators.** We show an intriguing fact: The bias brought
 54 by the sign-based compression can be flexibly controlled by injecting a proper amount of
 55 random noise before the sign operation. In particular, our analysis is based on a novel noisy
 56 perturbation scheme with a general symmetric noise distribution, and therefore provides a
 57 unified viewpoint and incorporates existing stochastic sign-based methods, including [Jin
 58 et al., 2020, Safaryan and Richtárik, 2021, Chen et al., 2020a], as special instances.
- 59 (2) **The first sign-based federated averaging algorithm.** In contrast to the existing sign-based
 60 schemes which do not allow local multiple SGD updates within one communication round,
 61 based on the proposed stochastic sign-based compression, we design a novel family of
 62 sign-based federated averaging algorithms (z -SignFedAvg) that can achieve the best of both
 63 worlds: communication efficiency and convergence performance.
- 64 (3) **New theoretical convergence rate analyses.** By a clever use of the asymptotic unbiased-
 65 ness property of the stochastic sign-based compression, we derive a series of theoretical
 66 results which exhibit better convergence rate than the existing sign-based methods. More-
 67 over, we show that by injecting a sufficiently large uniform noise, our proposed algorithm
 68 can have a matching convergence rate with the uncompressed algorithms.

69 **Organization of this paper.** In Section 2, the proposed general noisy perturbation scheme for the
 70 sign-based compression and its key propoerty about asymptotic unbiasedness are presented. Inspired
 71 by this result, the main algorithms are developed in Section 3 together with their convergence analyses
 72 under different noise distribution parameters. We evaluate our proposed algorithms on real datasets
 73 and benchmark with existing FL methods in Section 4. Finally, conclusions are drawn in Section 5.

74 **Notations.** For any $x \in \mathbb{R}^d$, we denote $x(j)$ as the j -th element of the vector x . We define
 75 the ℓ_p norm for any $p \geq 1$ as $\|x\|_p = (\sum_{j=1}^d |x(j)|^p)^{\frac{1}{p}}$. We denote that $\|\cdot\| = \|\cdot\|_2$, and
 76 $\|x\|_\infty = \max_{j \in \{1, \dots, d\}} |x(j)|$. For any function $f(x)$, we denote $f^{(k)}(x)$ as its k -th derivative, and
 77 for a vector $x = [x(1), \dots, x(d)]^\top \in \mathbb{R}^d$, we define $\text{Sign}(x) = [\text{Sign}(x(1)), \dots, \text{Sign}(x(d))]^\top$.

78 2 Sign operator with symmetric and zero-mean noise is asymptotically 79 unbiased.

80 In this section, we introduce a general noisy perturbation scheme for the sign-based compression and
 81 analyze its asymptotic unbiasedness. The result serves as the foundation of the proposed algorithm
 82 designs in subsequent sections. Specifically, let us consider the following family of noise distribution
 83 parameterized by a postive interger $z \in \mathbb{Z}_+$.

84 **Definition 1** (z -distribution). *A random variable ξ_z follows z -distribution if its p.d.f is*

$$p_z(t) = \frac{1}{2\eta_z} e^{-\frac{t^2 z}{2}}, \quad (2)$$

85 where $\eta_z = 2^{\frac{1}{2z}} \Gamma\left(1 + \frac{1}{2z}\right)$ and $\Gamma(z) = \int_0^{+\infty} t^{z-1} e^{-t} dt$ is the Gamma function.

86 It can be verified that (2) is a valid p.d.f. When $z = 1$, it corresponds to the standard Gaussian
 87 distribution. In addition, one can also show that (2) converges to the uniform distribution when
 88 $z \rightarrow +\infty$ (Lemma 2 in Appendices). This family of distribution has a nice property that can be
 89 leveraged to bound the bias caused by the sign compression, as stated in the following lemma.

90 **Lemma 1.** For any $x \in \mathbb{R}^d$ and $\sigma > 0$,

$$\|\eta_z \sigma \mathbb{E}[\text{Sign}(x + \sigma \xi_z)] - x\|^2 \leq \frac{\|x\|_{4z+2}^{4z+2}}{4(2z+1)^2 \sigma^{4z}}, \quad (3)$$

91 where $\xi_z(1), \dots, \xi_z(d)$ follow the i.i.d. z -distribution.

92 **Remark 1.** One can see that the RHS of (3) involves the term $(\|x\|_{4z+2}/\sigma)^{4z}$. Therefore, as long as
 93 $\sigma > \|x\|_\infty$, the LHS of (3) converges to zero when $z \rightarrow +\infty$. Since Lemma 2 implies that ξ_∞ is a
 94 vector whose elements follow i.i.d uniform distribution at $[-1, 1]$, we obtain $\sigma \mathbb{E}[\text{Sign}(x + \sigma \xi_\infty)] = x$
 95 as long as $\sigma > \|x\|_\infty$. It is interesting to remark that the stochastic sign operators proposed in [Jin
 96 et al., 2020, Safaryan and Richtárik, 2021] is exactly the sign operator injected with a sufficient
 97 amount of uniform noise.

98 3 Stochastic sign-based federated averaging with z -distribution.

99 In this section, based on the analysis in Section 2, we propose the following sign-based federated
 100 averaging algorithm, termed as z -SignFedAvg. The algorithm details are stated in Algorithm 1. Note
 101 that in practice, we only implement z -SignFedAvg with $z = 1$ and $z = +\infty$ which correspond to
 102 the Gaussian distribution and uniform distribution, respectively. This is because, to the best of our
 103 knowledge, there is currently no efficient way to sampling from the distribution $p_z(t)$ with other
 104 choices of z . Nevertheless, we are interested in the convergence properties of z -SignFedAvg for a
 105 general positive integer z since it provides better insights on how z impacts the convergence rate.

Algorithm 1 z -SignFedAvg (or z -SignSGD when $E = 1$)

Require: Total communication rounds T , Number of local steps E , Number of clients n , Clients
 stepsize γ , Server stepsize η , Noise coefficient σ , parameter of noise distribution z .

- 1: Initialize x_0 and for $i = 1, \dots, n$.
 - 2: **for** $t = 1$ to T **do**
 - 3: **On Clients:**
 - 4: **for** $i = 1$ to n in parallel **do**
 - 5: $x_{t-1,0}^i = x_{t-1}$
 - 6: **for** $s = 1$ to E **do**
 - 7: $g_{t-1,s}^i = g_i(x_{t-1,s-1}^i)$, where $g_i(\cdot)$ is the minibatch gradient oracle of the i -th client.
 - 8: $x_{t-1,s}^i = x_{t-1,s-1}^i - \gamma g_{t-1,s}^i$.
 - 9: **end for**
 - 10: $\Delta_{t-1}^i = \text{Sign}\left(\frac{x_{t-1} - x_{t-1,E}^i}{\gamma} + \sigma \xi_z\right)$, where $\xi_z(1), \dots, \xi_z(d) \sim p_z(t)$ i.i.d.
 - 11: Send Δ_{t-1}^i to the server.
 - 12: **end for**
 - 13: **On Server:**
 - 14: $x_t = x_{t-1} - \eta \gamma \frac{1}{n} \sum_{i=1}^n \Delta_{t-1}^i$.
 - 15: Send x_t to the clients.
 - 16: **end for**
 - 17: **return** x_T .
-

106 We first state the following standard assumptions.

107 **Assumption 1.** We assume that each $f_i(x)$ has the following properties:

108 A.1 The minibatch gradient is unbiased and has bounded variance, i.e., $\mathbb{E}[g_i(x)] =$
 109 $\nabla f_i(x)$, $\mathbb{E}[\|g_i(x) - \nabla f_i(x)\|_2^2] \leq \zeta^2$.

110 A.2 Each f_i is smooth, i.e., for any $x, y \in \mathbb{R}^d$, there exists some non-negative constants L_1, \dots, L_d
 111 such that $f(y) - f(x) \leq \langle \nabla f(x), y - x \rangle + \frac{\sum_{j=1}^d L_j (y^{(j)} - x^{(j)})^2}{2}$.

112 A.3 There exists some constant f^* such that $f(x) \geq f^*, \forall x \in \mathbb{R}^d$

113 A.4 There exists a constant $G \geq 0$ such that $\|\nabla f_i(x)\| \leq G, \forall i = 1, \dots, n$, and $x \in \mathbb{R}^d$.

114 For the convergence rate analysis, we consider two cases, namely, the case with $z < +\infty$ and the
115 case of $z = \infty$.

116 **3.1 Case 1:** $z < +\infty$

117 As we can see from Lemma 1, there always exists some gradient bias when $z < +\infty$. In order to
118 bound it, we further assume that a higher order moment of the minibatch gradient noise is bounded.

119 **Assumption 2.** There exists a constant $Q_z \geq 0$ such that for any $x \in \mathbb{R}^d$, we have

$$\mathbb{E}[\|g_i(x) - \nabla f_i(x)\|_{4z+2}^{4z+2}] \leq Q_z. \quad (4)$$

120 **Theorem 1.** Suppose that Assumption 1 and 2 hold. Denote $\bar{x}_{t,s} = \frac{1}{n} \sum_{i=1}^n x_{t,s}^i$ and $L_{\max} =$
121 $\max_j L_j$. Then, for $\eta = \eta_z \sigma$ and $\gamma \leq \frac{1}{L_{\max}}$, we have

$$\mathbb{E} \left[\frac{1}{TE} \sum_{t=1}^T \sum_{s=1}^E \|\nabla f(\bar{x}_{t-1,s-1})\|^2 \right] \leq \underbrace{\frac{2\mathbb{E}[f(x_0) - f^*]}{TE\gamma} + \frac{\gamma\zeta^2 L_{\max}}{n} + \frac{(E-1)(2E-1)\gamma^2 L_{\max}^2 G^2}{3}}_{\text{(a). Standard terms in FedAvg}} \quad (5a)$$

$$+ \underbrace{\frac{2^{2z} E^{2z} \sqrt{Q_z + G^{4z+2}} G}{(2z+1)\sigma^{2z}} + \frac{\gamma^{2z} E^{4z+1} (Q_z + G^{4z+2}) L_{\max}}{2(2z+1)^2 \sigma^{4z}}}_{\text{(b). Bias terms}} \quad (5b)$$

$$+ \underbrace{\frac{4\eta_z^2 \gamma \sigma^2 \sum_{j=1}^d L_j}{En}}_{\text{(c). Variance term}}. \quad (5c)$$

122 **When is the bound non-trivial?** Since we assume that the ℓ_2 -norm of gradient is bounded by G , all
123 the terms in the RHS of (5) should be no larger than G^2 . For example, one can check that to have the
124 first term in (5b) less than G^2 , one requires σ to be greater than $2E (Q_z + G^{4z})^{\frac{1}{4z}} / (2z+1)^{\frac{1}{2z}}$.

125 **Bias-variance trade-off.** An interesting observation from Theorem 1 is that there exists a trade-off
126 between the bias and variance terms. One can see that the terms in (5b) is caused by the gradient bias
127 of the sign operation (see (1)) and is an infinitesimal of σ with $\mathcal{O}(\sigma^{-2z})$, while the term in (5c) is due
128 to the injected noise and is in the order of $\mathcal{O}(\gamma\sigma^2)$. Specifically, the first term in (b) only depends on
129 the noise scale σ and mostly affects the final learning performance. In the meanwhile, the variance
130 term mainly affects the convergence speed because a smaller stepsize is required for it to diminish.

131 Theoretically, we can choose an iteration-dependent noise scale σ so as to making the algorithm
132 converge. In the following results, we denote $\tau = TE$ as the total number of gradient queries to the
133 local objective function.

134 **Corollary 1 (Informal).** Let $\gamma = \min\{n^{\frac{z}{2z+1}} \tau^{-\frac{z+1}{2z+1}}, \frac{1}{L_{\max}}\}$ and $\sigma = (n\tau)^{\frac{1}{4z+2}}$ in Theorem 1, and
135 let $E \leq n^{-\frac{3z}{4z+2}} \tau^{\frac{z+2}{4z+2}}$. We have

$$\mathbb{E} \left[\frac{1}{\tau} \sum_{t=1}^T \sum_{s=1}^E \|\nabla f(\bar{x}_{t-1,s-1})\|^2 \right] = \mathcal{O}((n\tau)^{-\frac{z}{2z+1}}). \quad (6)$$

136 **Achieving linear speedup.** From Corollary 1, we can see that the z -SignFedAvg needs $(n\tau)^{\frac{3z}{4z+2}}$
137 communication rounds to achieve a linear speedup convergence rate. Particularly, when $z = 1$, the
138 corresponding convergence speed is $\mathcal{O}((n\tau)^{-\frac{1}{3}})$ and the required communication rounds is $(n\tau)^{\frac{1}{2}}$.

139 **3.2 Case 2:** $z = +\infty$

140 In this case, the injected noise in z -SignFedAvg is uniformly distributed at $[-1, 1]$. From Remark 1
141 we have learned that the bias term in (5b) will either blow up when σ is lower than some threshold, or
142 vanish on the other hand. To quantify this threshold, we need the limit form of Assumption 2:

143 **Assumption 3.** *There exists a constant $Q_\infty \geq 0$ such that, with probability 1 we have*

$$\|g_i(x) - \nabla f_i(x)\|_\infty \leq Q_\infty, \forall x \in \mathbb{R}^d. \quad (7)$$

144 **Theorem 2 (Informal).** *Suppose that Assumption 1 and 3 hold. For $\gamma = \min\{n^{\frac{1}{2}}\tau^{-\frac{1}{2}}, \frac{1}{L_{\max}}\}$, $\eta = \sigma$,*
 146 *$E \leq n^{-\frac{3}{4}}\tau^{\frac{1}{4}}$, and $\sigma > E(G + Q_\infty)$, we have*

$$\mathbb{E} \left[\frac{1}{\tau} \sum_{t=1}^T \sum_{s=1}^E \|\nabla f(\bar{x}_{t-1, s-1})\|^2 \right] = \mathcal{O} \left((n\tau)^{-\frac{1}{2}} \right). \quad (8)$$

147 We can see that (8) implies ∞ -SignFedAvg recovers the convergence rate of uncompressed algorithms
 148 [Yu et al., 2019].

150 **Remark 2.** *It is worthwhile to point out that the condition of sufficiently large noise scale $\sigma >$*
 151 *$E(G + Q_\infty)$ is necessary and cannot be spared. By intuition, this is because when $\sigma \leq E(G + Q_\infty)$*
 152 *in Theorem 2, the injected uniform noise cannot change the sign of gradients in the worst case. For*
 153 *example, if ξ_∞ follows uniform distribution on $[-1, 1]$, and now $\sigma < A$ for some $A > 0$, we have*
 154 *$\text{Sign}(x + \sigma\xi_\infty) = \text{Sign}(x)$ for any $x \geq A$.*

155 **Remark 3.** *By comparing the required thresholds for σ in Theorem 1 and Theorem 2, we can see*
 156 *that when there is no minibatch gradient noise (i.e., $\zeta = 0$), Case 2 demands less noise injection*
 157 *and may perform better than Case 1. On the contrary, when the minibatch gradient noise has a long*
 158 *tail such as $Q_z \ll Q_\infty^{4z}$, Case 1 may be better. Despite of the distinction in theory, as we will see in*
 159 *Section 4, Case 1 and Case 2 have almost the same behavior on real datasets.*

160 More detailed theoretical results and comparison with existing methods are relegated to Appendix B.

161 4 Experiments

162 In this section, we present the experiment results on real datasets. All the figures are obtained by 10
 163 independent runs and are visualized in the form of mean \pm std. We also conduct an experiment on
 164 synthetic data where there is no minibatch gradient noise, and the results is relegated to Appendix D.

165 **Noise scale as a hyperparameter.** Although we explicitly characterize how the performance of
 166 Algorithm 1 depends on the noise scale σ in previous section, we treat σ as a tunable hyperparameter
 167 in practice. Because, on one hand, it usually impossible to compute the moment and support of
 168 the gradient noise in reality. One the other hand, since the theoretical results above only provide a
 169 worst-case guarantee, for some real problems, the optimal noise scales selected from grid search can
 170 be much smaller than the choice suggested by theory.

171 4.1 An extremely non-i.i.d setting

172 In this section, we consider an extremely non-i.i.d setting with the well-known dataset MNIST [Deng,
 173 2012] which is a hand-written digits recognition dataset. Specifically, we split the dataset into 10
 174 parts based on the labels and each client only have the data of one digit. In such a highly heterogeneous
 175 setting, there is no benefit from local computation with periodic communication. Therefore,
 176 we compare with the listed algorithms: SGDwM [Ghadimi and Lan, 2013], EF-SignSGDwM
 177 [Karimireddy et al., 2019, Vogels et al., 2019], Sto-SignSGDwM [Safaryan and Richtárik, 2021].
 178 Some baseline algorithms have an additional hyperparameter for the momentum of gradient. For all
 179 the algorithms, we select the their own optimal hyperparameters like stepsize, momentum coefficient,
 180 noise scale via grid search. For more details like hyperparameters for all the tested algorithms and the
 181 performance of 1-SignSGD and ∞ -SignSGD under different noise scales, we refer the readers to
 182 Appendix E.1.

183 **Results.** From Figure 1a, 1b, we can observe that 1-SignSGD and ∞ -SignSGD have roughly the
 184 same performance which outperform other sign-based algorithms and is slightly inferior to the
 185 uncompressed algorithm. Our theory is verified by comparing the performance of noiseless SignSGD
 186 and our proposed algorithms. If we compare the performance with respect to the total number of
 187 transmitted bits, our algorithms achieve the state-of-the-art performance on this task as we can see in
 188 Figure 1c.

189 4.2 Federated Learning on EMNIST

190 In this section, we verify the performance of our proposed Algorithm 1 on EMNIST[Cohen et al.,
 191 2017]. We mainly compare the performance of three algorithms: FedAvg without any compression
 192 [McMahan et al., 2017, Yu et al., 2019] and our proposed Algorithm 1 with $z = 1$ and $z = \infty$.

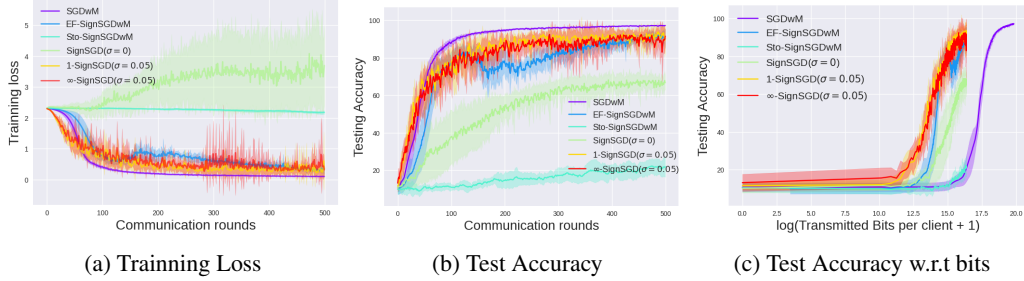


Figure 1: Performance of various algorithms on non-i.i.d MNIST

193 **Settings.** We follow a similar setting to [Reddi et al., 2020]. We also consider the scenario with
 194 partial participation. Specifically, for the EMNIST dataset, there are 3579 clients in total and we
 195 sample 100 clients uniformly to upload their local gradients at each communication round.

196 **Results.** The hyperparameters for the algorithms are tuned via grid search and details are in Appendix
 197 E.2. Specifically, we use $\sigma = 0.01$ for both 1-SignFedAvg and ∞ -SignFedAvg on EMNIST dataset.
 198 We can see from Figure 2 that all the algorithms can benefit from multiple local steps, and more
 199 surprisingly, both 1-SignFedAvg and ∞ -SignFedAvg can outperform the uncompressed algorithm
 200 FedAvg. This is probably because the EMNIST dataset is less non-i.i.d as the dataset we use in
 201 Section 4.1. The performance of 1-SignFedAvg and ∞ -SignFedAvg under various choices of noise
 202 scale are relegated to the Figure 6 and 7 in Appendix E.2, which also matches our theoretical results.

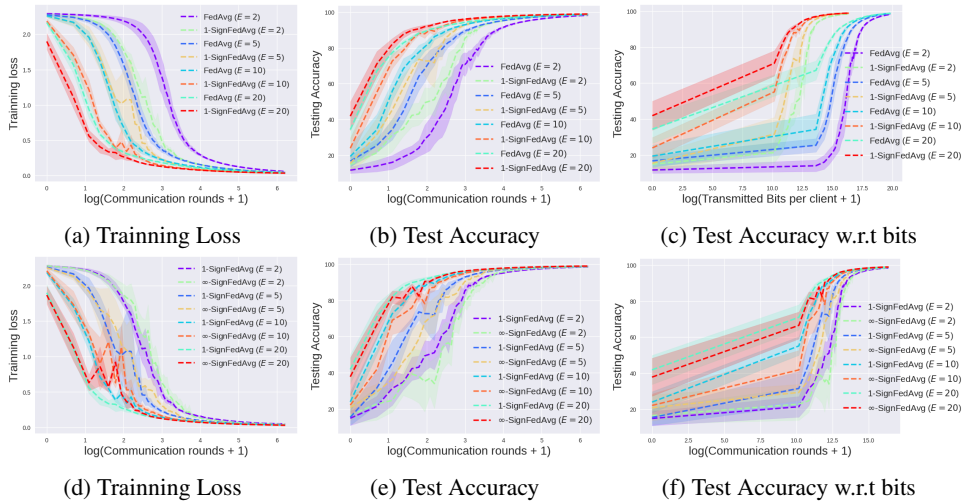


Figure 2: Performance of FedAvg, 1-SignFedAvg and ∞ -SignFedAvg on EMNIST dataset.

203 5 Conclusion

204 In this work, we have proposed the z -SignFedAvg: a FedAvg-type algorithm with a novel family of
 205 sign-based stochastic compression. Throughout extensive theoretical analysis and empirical evaluation,
 206 we have shown that z -SignFedAvg can perform comparably, sometimes even better, as the
 207 uncompressed FedAvg algorithm with a significantly reduced number of bits transmitted from the
 208 clients to the server. However, a vital issue in z -SignFedAvg is that it involves a new hyperparameter,
 209 i.e., the noise scale σ , which needs to be carefully chosen for achieving a good convergence performance.
 210 An interesting observation from the experiments is that the less heterogeneous the local data
 211 are, the smaller the optimal noise scale is, which is consistent with the theoretical insights. In the
 212 future, we will further study the relationship between the client’s heterogeneity and the optimal noise
 213 scale. As a final remark, we note that the stochastic sign-based compression proposed in this work is
 214 of independent interest and can be directly combined with other adaptive FL algorithms like those in
 215 Karimireddy et al. [2020], Reddi et al. [2020].

216 **References**

- 217 Jeremy Bernstein, Yu-Xiang Wang, Kamyar Azizzadenesheli, and Animashree Anandkumar. signsgd:
218 Compressed optimisation for non-convex problems. In *International Conference on Machine*
219 *Learning*, pages 560–569. PMLR, 2018.
- 220 Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas.
221 Communication-efficient learning of deep networks from decentralized data. In *Artificial intelli-*
222 *gence and statistics*, pages 1273–1282. PMLR, 2017.
- 223 Tian Li, Anit Kumar Sahu, Ameet Talwalkar, and Virginia Smith. Federated learning: Challenges,
224 methods, and future directions. *IEEE Signal Processing Magazine*, 37(3):50–60, 2020a.
- 225 Jakub Konečný, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh, and
226 Dave Bacon. Federated learning: Strategies for improving communication efficiency. *arXiv*
227 *preprint arXiv:1610.05492*, 2016.
- 228 Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. Qsgd: Communication-
229 efficient sgd via gradient quantization and encoding. *Advances in neural information processing*
230 *systems*, 30, 2017.
- 231 Amirhossein Reisizadeh, Aryan Mokhtari, Hamed Hassani, Ali Jadbabaie, and Ramtin Pedarsani.
232 Fedpaq: A communication-efficient federated learning method with periodic averaging and quan-
233 tization. In *International Conference on Artificial Intelligence and Statistics*, pages 2021–2031.
234 PMLR, 2020.
- 235 Sai Praneeth Karimireddy, Quentin Rebjock, Sebastian Stich, and Martin Jaggi. Error feedback
236 fixes signsgd and other gradient compression schemes. In *International Conference on Machine*
237 *Learning*, pages 3252–3261. PMLR, 2019.
- 238 Mher Safaryan and Peter Richtárik. Stochastic sign descent methods: New algorithms and better
239 theory. In *International Conference on Machine Learning*, pages 9224–9234. PMLR, 2021.
- 240 Xiangyi Chen, Tiancong Chen, Haoran Sun, Steven Z Wu, and Mingyi Hong. Distributed training
241 with heterogeneous data: Bridging median-and mean-based algorithms. *Advances in Neural*
242 *Information Processing Systems*, 33:21616–21626, 2020a.
- 243 Richeng Jin, Yufan Huang, Xiaofan He, Huaiyu Dai, and Tianfu Wu. Stochastic-sign sgd for federated
244 learning with theoretical guarantees. *arXiv preprint arXiv:2002.10940*, 2020.
- 245 Thijs Vogels, Sai Praneeth Karimireddy, and Martin Jaggi. Powersgd: Practical low-rank gradient
246 compression for distributed optimization. *Advances in Neural Information Processing Systems*, 32,
247 2019.
- 248 Hanlin Tang, Chen Yu, Xiangru Lian, Tong Zhang, and Ji Liu. Doublesqueeze: Parallel stochastic
249 gradient descent with double-pass error-compensated compression. In *International Conference*
250 *on Machine Learning*, pages 6155–6165. PMLR, 2019.
- 251 Saeed Ghadimi and Guanghui Lan. Stochastic first- and zeroth-order methods for nonconvex stochas-
252 tic programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013. doi:10.1137/120880811.
253 URL <https://doi.org/10.1137/120880811>.
- 254 Hao Yu, Sen Yang, and Shenghuo Zhu. Parallel restarted sgd with faster convergence and less
255 communication: Demystifying why model averaging works for deep learning. In *Proceedings of*
256 *the AAAI Conference on Artificial Intelligence*, volume 33, pages 5693–5700, 2019.
- 257 Li Deng. The mnist database of handwritten digit images for machine learning research [best of the
258 web]. *IEEE signal processing magazine*, 29(6):141–142, 2012.
- 259 Gregory Cohen, Saeed Afshar, Jonathan Tapson, and Andre Van Schaik. Emnist: Extending mnist to
260 handwritten letters. In *2017 international joint conference on neural networks (IJCNN)*, pages
261 2921–2926. IEEE, 2017.

- 262 Sashank Reddi, Zachary Charles, Manzil Zaheer, Zachary Garrett, Keith Rush, Jakub Konečný,
263 Sanjiv Kumar, and H Brendan McMahan. Adaptive federated optimization. *arXiv preprint*
264 *arXiv:2003.00295*, 2020.
- 265 Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and
266 Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In
267 *International Conference on Machine Learning*, pages 5132–5143. PMLR, 2020.
- 268 Jinshuo Dong, Aaron Roth, and Weijie Su. Gaussian differential privacy. *Journal of the Royal*
269 *Statistical Society*, 2021.
- 270 Yanmeng Wang, Yanqing Xu, Qingjiang Shi, and Tsung-Hui Chang. Quantized federated learning un-
271 der transmission delay and outage constraints. *IEEE Journal on Selected Areas in Communications*,
272 40(1):323–341, 2021.
- 273 Cynthia Dwork, Aaron Roth, et al. The algorithmic foundations of differential privacy. *Found. Trends*
274 *Theor. Comput. Sci.*, 9(3-4):211–407, 2014.
- 275 Robin C Geyer, Tassilo Klein, and Moin Nabi. Differentially private federated learning: A client
276 level perspective. *arXiv preprint arXiv:1712.07557*, 2017.
- 277 Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and
278 Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC*
279 *conference on computer and communications security*, pages 308–318, 2016a.
- 280 Yiwei Li, Tsung-Hui Chang, and Chong-Yung Chi. Secure federated averaging algorithm with
281 differential privacy. In *2020 IEEE 30th International Workshop on Machine Learning for Signal*
282 *Processing (MLSP)*, pages 1–6. IEEE, 2020b.
- 283 Naman Agarwal, Ananda Theertha Suresh, Felix Xinnan X Yu, Sanjiv Kumar, and Brendan McMahan.
284 cpsgd: Communication-efficient and differentially-private distributed sgd. *Advances in Neural*
285 *Information Processing Systems*, 31, 2018.
- 286 Zhiqi Bu, Jinshuo Dong, Qi Long, and Weijie J Su. Deep learning with gaussian differential privacy.
287 *Harvard data science review*, 2020(23), 2020.
- 288 Qinqing Zheng, Shuxiao Chen, Qi Long, and Weijie Su. Federated f-differential privacy. In
289 *International Conference on Artificial Intelligence and Statistics*, pages 2251–2259. PMLR, 2021.
- 290 Saba Amiri, Adam Belloum, Sander Klous, and Leon Gommans. Compressive differentially-private
291 federated learning through universal vector quantization. 2021.
- 292 Lun Wang, Ruoxi Jia, and Dawn Song. D2p-fed: Differentially private federated learning with
293 efficient communication. *arXiv preprint arXiv:2006.13039*, 2020.
- 294 Xiangru Lian, Ce Zhang, Huan Zhang, Cho-Jui Hsieh, Wei Zhang, and Ji Liu. Can decentralized
295 algorithms outperform centralized algorithms? a case study for decentralized parallel stochastic
296 gradient descent. *Advances in Neural Information Processing Systems*, 30, 2017.
- 297 Ilya Mironov. Rényi differential privacy. In *2017 IEEE 30th computer security foundations symposium*
298 *(CSF)*, pages 263–275. IEEE, 2017.
- 299 Shahab Asoodeh, Jiachun Liao, Flavio P Calmon, Oliver Kosut, and Lalitha Sankar. Three variants
300 of differential privacy: Lossless conversion and applications. *IEEE Journal on Selected Areas in*
301 *Information Theory*, 2(1):208–222, 2021.
- 302 Igal Sason and Sergio Verdú. f -divergence inequalities. *IEEE Transactions on Information Theory*,
303 62(11):5973–6006, 2016.
- 304 Laurent Condat, Kai Yi, and Peter Richtárik. Ef-bv: A unified theory of error feedback and variance
305 reduction mechanisms for biased and unbiased compression in distributed optimization. *arXiv*
306 *preprint arXiv:2205.04180*, 2022.

- 307 Shuai Zheng, Ziyue Huang, and James Kwok. Communication-efficient distributed blockwise
308 momentum sgd with error-feedback. *Advances in Neural Information Processing Systems*, 32,
309 2019.
- 310 Xiangyi Chen, Steven Z Wu, and Mingyi Hong. Understanding gradient clipping in private sgd: A
311 geometric perspective. *Advances in Neural Information Processing Systems*, 33:13773–13782,
312 2020b.
- 313 Xinwei Zhang, Xiangyi Chen, Mingyi Hong, Zhiwei Steven Wu, and Jinfeng Yi. Understanding
314 clipping for federated learning: Convergence and client-level differential privacy. *arXiv preprint*
315 *arXiv:2106.13673*, 2021.
- 316 Zhiwei Tang, Tsung-Hui Chang, Xiaojing Ye, and Hongyuan Zha. Low-rank matrix recovery with
317 unknown correspondence. *arXiv preprint arXiv:2110.07959*, 2021a.
- 318 Peter Kairouz, Ziyu Liu, and Thomas Steinke. The distributed discrete gaussian mechanism for
319 federated learning with secure aggregation. In *International Conference on Machine Learning*,
320 pages 5201–5212. PMLR, 2021.
- 321 John T Chu. On bounds for the normal integral. *Biometrika*, 42(1/2):263–265, 1955.
- 322 Hanlin Tang, Shaoduo Gan, Ammar Ahmad Awan, Samyam Rajbhandari, Conglong Li, Xiangru Lian,
323 Ji Liu, Ce Zhang, and Yuxiong He. 1-bit adam: Communication efficient large-scale training with
324 adam’s convergence speed. In *International Conference on Machine Learning*, pages 10118–10129.
325 PMLR, 2021b.
- 326 Jianyu Wang and Gauri Joshi. Cooperative sgd: A unified framework for the design and analysis of
327 local-update sgd algorithms. *Journal of Machine Learning Research*, 22, 2021.
- 328 Frank Seide, Hao Fu, Jasha Droppo, Gang Li, and Dong Yu. 1-bit stochastic gradient descent and its
329 application to data-parallel distributed training of speech dnns. In *Fifteenth annual conference of*
330 *the international speech communication association*. Citeseer, 2014.
- 331 Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu
332 Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. {TensorFlow}: a system for
333 {Large-Scale} machine learning. In *12th USENIX symposium on operating systems design and*
334 *implementation (OSDI 16)*, pages 265–283, 2016b.
- 335 Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito,
336 Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in
337 pytorch. 2017.
- 338 Alex Krizhevsky and Geoff Hinton. Convolutional deep belief networks on cifar-10. *Unpublished*
339 *manuscript*, 40(7):1–9, 2010.
- 340 Tim Van Erven and Peter Harremos. Rényi divergence and kullback-leibler divergence. *IEEE*
341 *Transactions on Information Theory*, 60(7):3797–3820, 2014.
- 342 Dimitri P Bertsekas. Nonlinear programming. *Journal of the Operational Research Society*, 48(3):
343 334–334, 1997.
- 344 Naman Agarwal, Peter Kairouz, and Ziyu Liu. The skellam mechanism for differentially private
345 federated learning. *Advances in Neural Information Processing Systems*, 34:5052–5064, 2021.
- 346 Yangsibo Huang, Samyak Gupta, Zhao Song, Kai Li, and Sanjeev Arora. Evaluating gradient
347 inversion attacks and defenses in federated learning. *Advances in Neural Information Processing*
348 *Systems*, 34:7232–7241, 2021.
- 349 Jonas Geiping, Hartmut Bauermeister, Hannah Dröge, and Michael Moeller. Inverting gradients-how
350 easy is it to break privacy in federated learning? *Advances in Neural Information Processing*
351 *Systems*, 33:16937–16947, 2020.
- 352 Ilya Mironov, Kunal Talwar, and Li Zhang. R^{ϵ} ’enyi differential privacy of the sampled gaussian
353 mechanism. *arXiv preprint arXiv:1908.10530*, 2019.

354 Checklist

355 The checklist follows the references. Please read the checklist guidelines carefully for information on
356 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or
357 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing
358 the appropriate section of your paper or providing a brief inline description. For example:

- 359 • Did you include the license to the code and datasets? **[Yes]** See Section ??.
- 360 • Did you include the license to the code and datasets? **[No]** The code and the data are
361 proprietary.
- 362 • Did you include the license to the code and datasets? **[N/A]**

363 Please do not modify the questions and only use the provided macros for your answers. Note that the
364 Checklist section does not count towards the page limit. In your paper, please delete this instructions
365 block and only keep the Checklist section heading above along with the questions/answers below.

366 1. For all authors...

- 367 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
368 contributions and scope? **[TODO]**
- 369 (b) Did you describe the limitations of your work? **[TODO]**
- 370 (c) Did you discuss any potential negative societal impacts of your work? **[TODO]**
- 371 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
372 them? **[TODO]**

373 2. If you are including theoretical results...

- 374 (a) Did you state the full set of assumptions of all theoretical results? **[TODO]**
- 375 (b) Did you include complete proofs of all theoretical results? **[TODO]**

376 3. If you ran experiments...

- 377 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
378 mental results (either in the supplemental material or as a URL)? **[TODO]**
- 379 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
380 were chosen)? **[TODO]**
- 381 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
382 ments multiple times)? **[TODO]**
- 383 (d) Did you include the total amount of compute and the type of resources used (e.g., type
384 of GPUs, internal cluster, or cloud provider)? **[TODO]**

385 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

- 386 (a) If your work uses existing assets, did you cite the creators? **[TODO]**
- 387 (b) Did you mention the license of the assets? **[TODO]**
- 388 (c) Did you include any new assets either in the supplemental material or as a URL?
389 **[TODO]**
- 390 (d) Did you discuss whether and how consent was obtained from people whose data you're
391 using/curating? **[TODO]**
- 392 (e) Did you discuss whether the data you are using/curating contains personally identifiable
393 information or offensive content? **[TODO]**

394 5. If you used crowdsourcing or conducted research with human subjects...

- 395 (a) Did you include the full text of instructions given to participants and screenshots, if
396 applicable? **[TODO]**
- 397 (b) Did you describe any potential participant risks, with links to Institutional Review
398 Board (IRB) approvals, if applicable? **[TODO]**
- 399 (c) Did you include the estimated hourly wage paid to participants and the total amount
400 spent on participant compensation? **[TODO]**

401 **Appendices**

402 **A Discussion on related works**

403 Sign-based optimization algorithms such as the SignSGD in [Bernstein et al., 2018] have gained
404 much popularity recently because of their simple compression rule and comparable performance to
405 uncompressed algorithms. In this work, we focus on the scenario with heterogeneous data, and as
406 we have discussed in Section 1, a naive application of sign-based compression in this scenario is
407 problematic. Besides, we consider using sign-based compression only for the uplink communication
408 in this work, while it is worth mentioning that [Tang et al., 2019, Jin et al., 2020, Chen et al., 2020a]
409 also compress for the downlink communication. In the following paragraphs, we review a few related
410 works on similar topics.

411 **Stochastic sign-based method.** The setting considered by [Safaryan and Richtárik, 2021] is the
412 closest to ours because [Jin et al., 2020, Chen et al., 2020a] also compresses the server-to-client
413 communication with majority vote. Aside from the difference in setting, the algorithms in them
414 achieve the same convergence rate $O(\tau^{-\frac{1}{4}})$ w.r.t different convergence metrics, where τ is the number
415 of gradient queries to the objective function. As will be discussed in Appendix A, these rates are
416 usually inferior to that of uncompressed algorithms. Our proposed algorithm also belongs to this
417 category. Compared to existing works, we require a slightly stronger assumption on the gradient
418 noise, and the convergence speed of our algorithm is either $O(\tau^{-\frac{1}{3}})$ or $O(\tau^{-\frac{1}{2}})$ with the commonly
419 used squared ℓ_2 norm of gradients as the convergence metric. Moreover, we also show that our
420 proposed sign-based algorithm can achieve a linear speedup when the number of clients increases,
421 and such a result is not known in previous works.

422 **Error Feedback method.** The error feedback (EF) method is first proposed by [Seide et al., 2014]
423 and then theoretically justified by [Karimireddy et al., 2019]. Then [Vogels et al., 2019, Tang
424 et al., 2019, 2021b] further extend this EF framework into distributed non-i.i.d setting and adaptive
425 gradient method. The key idea is to show that the sign operator multiplying with one norm is a
426 contractive compressor, and the error induced by the contractive compressor can be fixed by the
427 error-compensated gradient method. However, unlike the pure sign-based gradient method, it must
428 transmit one extra real number for the one norm. Besides, the convergence rate for the EF algorithms
429 is $O(\tau^{-\frac{1}{2}} + \tau^{-1}/\delta^2)$, where δ is the parameter of contractive compressor. In the worst case, the
430 sign operator multiplying with one norm is a contractive compressor with $\delta = 1/d$, where d is
431 the dimension of the gradient. Therefore, the convergence rate of it becomes $O(\tau^{-\frac{1}{2}} + d^2\tau^{-1})$,
432 which could become very bad especially for high-dimension optimization problem. Besides, to our
433 knowledge, no one has extended the error feedback method to the scenario with periodic aggregation.

434 It is often tricky to compare the convergence results of sign-based methods because some works
435 like [Chen et al., 2020a, Safaryan and Richtárik, 2021] do not use the standard convergence metric.
436 To better compare existing results and ours, in Appendix A, we provide a detailed discussion on
437 the existing convergence metrics and summarize the representative algorithms and their theoretical
438 results in Table 1.

439 Table 1 gives a brief summary for a few representative works related to this work. In this table, we
440 review the algorithms in these works on the convergence rate along with the used convergence metrics,
441 communication complexity, assumptions required, and also whether they can deal with periodic
442 aggregation. Particularly, [Chen et al., 2020a, Safaryan and Richtárik, 2021] adopt convergence
443 metrics other than the commonly used average squared ℓ_2 norm of gradients. Due to the additional
444 step of server-to-client compression, [Chen et al., 2020a] use a convergence metric mixed with ℓ_2
445 norm and ℓ_1 norm, while [Safaryan and Richtárik, 2021] use ℓ_2 norm instead of squared ℓ_2 norm. For
446 communication complexity, we only compare the unlink communication, and to compute the used
447 bits per communication, we assume that all the algorithms need 32 bits to represent a float number as
448 this is the most common setting in Tensorflow [Abadi et al., 2016b] and Pytorch [Paszke et al., 2017].

449 Among the works in Table 1, the setting considered by [Safaryan and Richtárik, 2021] is the closest to
450 ours. [Safaryan and Richtárik, 2021] propose an algorithm that can achieve convergence rate $O(\tau^{-\frac{1}{4}})$
451 with average ℓ_2 norm of gradients as the metric. We remark that this is inferior to the convergence
452 rate $O(\tau^{-\frac{1}{2}})$ with squared ℓ_2 norm as the metric. To illustrate this point, we denote a serie of vector

Algorithm	Convergence metric / rate	Used bits per communication	Extra Assumptions?	Can achieve linear speedup?	Can allow multiple local steps?
[Ghadimi and Lan, 2013]	$\mathcal{O}(\tau^{-\frac{1}{2}})$ squared ℓ_2	$32d$	No	✓	✗
[Yu et al., 2019]	$\mathcal{O}(\tau^{-\frac{1}{2}})$ squared ℓ_2	$32d$	• Bounded gradient	✓	✓
[Karimireddy et al., 2019]	$\mathcal{O}(\tau^{-\frac{1}{2}} + d^2\tau^{-1})$ squared ℓ_2	$d + 32$	• Bounded gradient	✗	✗
[Safaryan and Richtárik, 2021]	$\mathcal{O}(\tau^{-\frac{1}{4}})$ ℓ_2	d	No	✗	✗
[Jin et al., 2020]	$\mathcal{O}(\tau^{-\frac{1}{4}})$ squared ℓ_2	d	• Bounded gradient • n is an odd number	✗	✗
[Chen et al., 2020a]	$\mathcal{O}(\tau^{-\frac{1}{4}})$ mixed	d	• Bounded gradient • n is an odd number	✗	✗
1-SignFedAvg (ALG. 1) This work	$\mathcal{O}(\tau^{-\frac{1}{3}})$ squared ℓ_2	d	• Bounded gradient • Bounded 6th moment of gradient noise	✓	✓
∞ -SignFedAvg (ALG. 1) This work	$\mathcal{O}(\tau^{-\frac{1}{2}})$ squared ℓ_2	d	• Bounded gradient • Bounded support of gradient noise	✓	✓

Table 1: Summary for related works.

453 $\{\alpha_1, \dots, \alpha_\tau, \dots\}$ with $\alpha_i \in \mathbb{R}^d$. If now

$$\frac{1}{\tau} \sum_{i=1}^{\tau} \|\alpha_i\| = \mathcal{O}(\tau^{-\frac{1}{4}}), \quad (9)$$

454 in the worst case, we can only guarantee that

$$\frac{1}{\tau} \sum_{i=1}^{\tau} \|\alpha_i\|^2 \leq \tau \left(\frac{1}{\tau} \sum_{i=1}^{\tau} \|\alpha_i\| \right)^2 = \mathcal{O}(\tau^{\frac{1}{2}}) \quad (10)$$

455 for squared ℓ_2 norm. As a simple example, the equality in (10) holds when there is only one non-zero
456 term in $\{\alpha_1, \dots, \alpha_\tau\}$.

457 On the contrary, if

$$\frac{1}{\tau} \sum_{i=1}^{\tau} \|\alpha_i\|^2 = \mathcal{O}(\tau^{-\frac{1}{2}}), \quad (11)$$

458 we have

$$\frac{1}{\tau} \sum_{i=1}^{\tau} \|\alpha_i\| \leq \sqrt{\frac{1}{\tau} \sum_{i=1}^{\tau} \|\alpha_i\|^2} = \mathcal{O}(\tau^{-\frac{1}{4}}). \quad (12)$$

459 Consider the scenario $E = 1$, the algorithm in [Safaryan and Richtárik, 2021] is equivalent to our
460 Algorithm 1 with σ chosen to be $\|g_{t-1,s}^i\|$. On one hand, this choice of noise scale σ make it unable
461 to be extended to the federated averaging algorithm, because each client use a different noise scale.
462 On the other hand, this choice is linearly increasing w.r.t problem dimension and hence is too
463 conservative. From Figure 3 and 1 we can see that this input-dependent noise scale result in an
464 extremely slow convergence for high-dimension problems.

465 B Theoretical results

466 In this section, we state the formal version of Corollary 1 and Theorem 2.

467 **Corollary 2** (Formal version of Corollary 1). *If we choose $\gamma = \min\{n^{\frac{z}{2z+1}}\tau^{-\frac{z+1}{2z+1}}, \frac{1}{L_{\max}}\}$ and*
 468 *$\sigma = (n\tau)^{\frac{1}{4z+2}}$ in Theorem 1, we have*

$$\mathbb{E}\left[\frac{1}{\tau} \sum_{t=1}^T \sum_{s=1}^E \|\nabla f(\bar{x}_{t-1,s-1})\|^2\right] \leq \frac{2\mathbb{E}[f(x_0) - f^*]}{(n\tau)^{\frac{z}{2z+1}}} + \frac{\zeta^2 L_{\max}}{(n\tau)^{\frac{z+1}{2z+1}}} + \frac{(E-1)(2E-1)n^{\frac{2z}{2z+1}}L_{\max}^2 G^2}{3\tau^{\frac{2z+2}{2z+1}}} \quad (13a)$$

$$+ \frac{2^{2z}E^{2z}\sqrt{Q_z + G^{4z+2}}G}{(2z+1)(n\tau)^{\frac{z}{2z+1}}} + \frac{2^{4z}E^{4z+1}(Q_z + G^{4z+2})L_{\max}}{2(2z+1)^2 n^{\frac{z}{2z+1}} \tau^{\frac{3z+1}{2z+1}}} \quad (13b)$$

$$+ \frac{4\eta_z^2 \sum_{j=1}^d L_j}{E(n\tau)^{\frac{z}{2z+1}}}. \quad (13c)$$

469 *Furthermore, if $E \leq n^{-\frac{3z}{4z+2}}\tau^{\frac{z+2}{4z+2}}$, the upper bound above will converge as $\mathcal{O}\left((n\tau)^{-\frac{z}{2z+1}}\right)$.*

470 **Relationship to [Chen et al., 2020a].** [Chen et al., 2020a] also studies the sign-based optimization
 471 algorithm with symmetric and zero-mean noise and prove that the convergence rate is $\mathcal{O}(\tau^{-\frac{1}{4}})$ using
 472 a similar iteration-dependent noise scale like us. However, there are two difference between their
 473 result and our Theorem 1. First, since they also consider the downlink compression, the convergence
 474 metric they used is no longer ℓ_2 norm and hard to interpret. Second, unlike [Chen et al., 2020a]
 475 whose result is rooted in median-based algorithm, our analysis directly exploits the property of sign
 476 operation and hence can provide a better and more interpretable result.

477 **Theorem 3** (Formal version of Theorem 2). *Given that Assumption 1 and 3 hold, and we choose*
 478 *$\eta = \sigma$, if $\gamma \leq \frac{1}{L_{\max}}$, if $\sigma > E(G + Q_{\infty})$, we have*

$$\mathbb{E}\left[\frac{1}{TE} \sum_{t=1}^T \sum_{s=1}^E \|\nabla f(\bar{x}_{t-1,s-1})\|^2\right] \leq \underbrace{\frac{2\mathbb{E}[f(x_0) - f^*]}{TE\gamma} + \frac{\gamma\zeta^2 L_{\max}}{n} + \frac{(E-1)(2E-1)\gamma^2 L_{\max}^2 G^2}{3}}_{\text{standard terms in federated averaging}} \quad (14a)$$

$$+ \underbrace{\frac{4\gamma\sigma^2 \sum_{j=1}^d L_j}{En}}_{\text{variance term}}. \quad (14b)$$

479 *otherwise, if $\sigma \leq E(G + Q_{\infty})$, there exists a problem where the algorithm cannot converge.*

480 *If we further choose $\gamma = \min\{n^{\frac{1}{2}}\tau^{-\frac{1}{2}}, \frac{1}{L_{\max}}\}$, we have*

$$\mathbb{E}\left[\frac{1}{\tau} \sum_{t=1}^T \sum_{s=1}^E \|\nabla f(\bar{x}_{t-1,s-1})\|^2\right] \leq \frac{2\mathbb{E}[f(x_0) - f^*]}{(n\tau)^{\frac{1}{2}}} + \frac{\zeta^2 L_{\max}}{(n\tau)^{\frac{1}{2}}} + \frac{(E-1)(2E-1)nL_{\max}^2 G^2}{3\tau} \quad (15)$$

$$+ \frac{4\sigma^2 \sum_{j=1}^d L_j}{E(n\tau)^{\frac{1}{2}}}. \quad (16)$$

481 *Furthermore, if $E \leq n^{-\frac{3}{4}}\tau^{\frac{1}{4}}$, the upper bound above will converge as $\mathcal{O}\left((n\tau)^{-\frac{1}{2}}\right)$, which recovers*
 482 *the convergence result of uncompressed algorithm [Yu et al., 2019].*

483 **Remark 4.** *When $\sigma \leq E(G + Q_{\infty})$ in Theorem 3, the injected uniform noise cannot change the sign*
 484 *of gradients in the worst case. For example, if ξ_{∞} follows uniform distribution on $[-1, 1]$, and now*
 485 *$\sigma < A$ for some $A > 0$, we have $\text{Sign}(x + \sigma\xi_{\infty}) = \text{Sign}(x)$ for any $x \geq A$.*

486 **Relationship to [Jin et al., 2020, Safaryan and Richtárik, 2021].** We remark that both the
 487 stochastic sign operators proposed in [Jin et al., 2020, Safaryan and Richtárik, 2021] are equivalent
 488 to the sign operator with uniform noise considered in Case 2. In particular, [Jin et al., 2020] also

489 consider downlink compression and hence its convergence results are not directly comparable to the
490 Case 2. [Safaryan and Richtárik, 2021] adopts an input-dependent noise scale and proves $\mathcal{O}(\tau^{-\frac{1}{4}})$
491 convergence rate with ℓ_2 norm of gradient as the metric. We remark that this rate is usually worse
492 than the rate $\mathcal{O}(\tau^{-\frac{1}{2}})$ with squared ℓ_2 norm as the metric. Such input-dependent noise scale make
493 it possible to prove convergence without the bounded support of gradient noise assumed in this
494 work. But there are two disadvantages for their choice of noise scale. First, it can not be extended
495 to federated averaging algorithm. Second, it often leads to slow convergence in practice when the
496 problem dimension is very high. More discussions on Safaryan and Richtárik [2021] are provided in
497 Appendix A.

498 C Missing proofs

499 **Lemma 2.** z -distribution weakly converges to uniform distribution at $[-1, 1]$ when $z \rightarrow +\infty$.

500 *Proof of Lemma 2.* Now we denote the p.d.f of uniform distribution as

$$p_\infty(x) = \begin{cases} \frac{1}{2} & |x| \leq 1, \\ 0 & |x| > 1. \end{cases} \quad (17)$$

501 Without loss of generality, for any $x > 1$ and $z \in \mathbb{Z}_+$, we have

$$\left| \int_{-\infty}^x \frac{1}{2\eta_z} e^{-\frac{t^2 z}{2}} dt - \int_{-\infty}^x p_\infty(t) dt \right| = \left| \int_0^x \left(\frac{1}{2\eta_z} e^{-\frac{t^2 z}{2}} - p_\infty(t) \right) dt \right| \quad (18a)$$

$$\leq \int_0^1 \left| \frac{1}{2\eta_z} e^{-\frac{t^2 z}{2}} - \frac{1}{2} \right| dt + \int_1^x \frac{1}{2\eta_z} e^{-\frac{t^2 z}{2}} dt. \quad (18b)$$

502 For any $0 < \epsilon < \min\{1, x-1\}$, we have

$$\int_0^1 \left| \frac{1}{2\eta_z} e^{-\frac{t^2 z}{2}} - \frac{1}{2} \right| dt = \int_0^{1-\epsilon} \left| \frac{1}{2\eta_z} e^{-\frac{t^2 z}{2}} - \frac{1}{2} \right| dt + \int_{1-\epsilon}^1 \left| \frac{1}{2\eta_z} e^{-\frac{t^2 z}{2}} - \frac{1}{2} \right| dt \quad (19a)$$

$$\leq \left| \frac{1}{2\eta_z} e^{-\frac{(1-\epsilon)^2 z}{2}} - \frac{1}{2} \right| + \epsilon. \quad (19b)$$

Since $\lim_{z \rightarrow \infty} \frac{1}{2\eta_z} = \lim_{z \rightarrow \infty} \frac{z}{2^{\frac{1}{2z}} \Gamma(\frac{1}{2z})} = \frac{1}{2}$ and $\lim_{z \rightarrow \infty} e^{-\frac{(1-\epsilon)^2 z}{2}} = 1$, there exists an interger $Z_1 > 0$ such that if $z > Z_1$, we have

$$\left| \frac{1}{2\eta_z} e^{-\frac{(1-\epsilon)^2 z}{2}} - \frac{1}{2} \right| \leq \epsilon.$$

503 Similarly, we have

$$\int_1^x \frac{1}{2\eta_z} e^{-\frac{t^2 z}{2}} dt = \int_1^{1+\epsilon} \frac{1}{2\eta_z} e^{-\frac{t^2 z}{2}} dt + \int_{1+\epsilon}^x \frac{1}{2\eta_z} e^{-\frac{t^2 z}{2}} dt \quad (20a)$$

$$\leq \epsilon + \frac{1}{2\eta_z} e^{-\frac{(1+\epsilon)^2 z}{2}} (x-1-\epsilon). \quad (20b)$$

504 Since $\lim_{z \rightarrow \infty} e^{-\frac{(1+\epsilon)^2 z}{2}} = 0$, there exists an interger $Z_2 > 0$ such that if $z > Z_2$, we have

$$\int_1^x \frac{1}{2\eta_z} e^{-\frac{t^2 z}{2}} dt \leq \epsilon. \quad (21)$$

505 In all, for any $0 < \epsilon < 1$, if z is sufficiently large, we have

$$\left| \int_{-\infty}^x \frac{1}{2\eta_z} e^{-\frac{t^2 z}{2}} dt - \int_{-\infty}^x p_\infty(t) dt \right| \leq 4\epsilon. \quad (22)$$

506 Take $\epsilon \rightarrow 0$ and $z \rightarrow \infty$, we have

$$\lim_{z \rightarrow \infty} \left| \int_{-\infty}^x \frac{1}{2\eta_z} e^{-\frac{t^{2z}}{2}} dt - \int_{-\infty}^x p_\infty(t) dt \right| = 0. \quad (23)$$

507

□

508 *Proof of Lemma 1.* We first state a useful inequality on the c.d.f of z distribution:

509 **Lemma 3.** For any $x \in \mathbb{R}$

$$|x| - \frac{|x|^{2z+1}}{2(2z+1)} \leq |\Psi_z(x)| \leq |x|, \text{ where } \Psi_z(x) \stackrel{\text{def.}}{=} \int_0^x e^{-\frac{t^{2z}}{2}} dt. \quad (24)$$

510 Then, we have

$$\|\eta_z \sigma \mathbb{E}[\text{Sign}(x + \sigma \xi_z)] - x\|^2 = \left\| x - \sigma \Psi_z\left(\frac{x}{\sigma}\right) \right\|^2 = \sum_{j=1}^d \left(x(j) - \sigma \Psi_z\left(\frac{x(j)}{\sigma}\right) \right)^2 \quad (25a)$$

$$\leq \sum_{j=1}^d \frac{x(j)^{4z+2}}{4(2z+1)^2 \sigma^{4z}} = \frac{\|x\|_{4z+2}^{4z+2}}{4(2z+1)^2 \sigma^{4z}}. \quad (25b)$$

511

□

512 *Proof of Lemma 3.* Without loss of generality, we prove it for $x \geq 0$.

513 First,

$$\int_0^x e^{-\frac{t^{2z}}{2}} dt \leq \int_0^x 1 dt \leq x. \quad (26)$$

514 Now we define $F(x) \stackrel{\text{def.}}{=} \int_0^x e^{-\frac{t^{2z}}{2}} dt - x + \frac{x^{2z+1}}{2(2z+1)}$. Note that $F(0) = 0$.

515 Then, we can prove that $F(x) \geq 0$ by

$$F'(x) = e^{-\frac{x^{2z}}{2}} - x + \frac{x^{2z}}{2} \geq 0. \quad (27)$$

516 (27) is due to the inequality $e^{-x} - 1 + x \geq 0$ for any $x \geq 0$.

□

517 *Proof of Theorem 1.* Here we define the virtual aggregated update:

$$\bar{x}_{t,s} = \frac{1}{n} \sum_{i=1}^n x_{t,s}^i, \quad (28)$$

$$\bar{x}_t = \bar{x}_{t-1,E}. \quad (29)$$

518 We now state the two useful lemmas:

Lemma 4.

$$\mathbb{E}[f(x_t) - f(\bar{x}_t)] \leq \frac{\gamma 2^{2z} E^{2z+1} \sqrt{Q_z + G^{4z+2}} G}{2(2z+1)\sigma^{2z}} + \frac{\gamma^2 2^{4z} E^{4z+2} (Q_z + G^{4z+2}) L_{\max}}{4(2z+1)^2 \sigma^{4z}} \quad (30a)$$

$$+ \frac{2\eta_z^2 \gamma^2 \sigma^2 \sum_{j=1}^d L_j}{n}. \quad (30b)$$

Lemma 5.

$$\mathbb{E}[f(\bar{x}_t) - f(x_{t-1})] \leq -\frac{\gamma}{2} \sum_{s=1}^E \|\nabla f(\bar{x}_{t-1,s-1})\|^2 + \frac{E\gamma^2 \zeta^2 L_{\max}}{2n} + \frac{E(E-1)(2E-1)\gamma^3 L_{\max}^2 G^2}{6}. \quad (31)$$

519 With this two lemma, we have

$$\mathbb{E}[f(x_t) - f(x_{t-1})] = \mathbb{E}[f(x_t) - f(\bar{x}_t)] + E[f(\bar{x}_t) - f(x_{t-1})] \quad (32a)$$

$$\leq -\frac{\gamma}{2} \sum_{s=1}^E \|\nabla f(\bar{x}_{t-1,s-1})\|^2 + \frac{E\gamma^2\zeta^2 L_{\max}}{2n} + \frac{E(E-1)(2E-1)\gamma^3 L_{\max}^2 G^2}{6} \quad (32b)$$

$$+ \frac{\gamma 2^{2z} E^{2z+1} \sqrt{Q_z + G^{4z+2}} G}{2(2z+1)\sigma^{2z}} + \frac{\gamma^2 2^{4z} E^{4z+2} (Q_z + G^{4z+2}) L_{\max}}{4(2z+1)^2 \sigma^{4z}} \quad (32c)$$

$$+ \frac{2\eta_z^2 \gamma^2 \sigma^2 \sum_{j=1}^d L_j}{n}. \quad (32d)$$

520 Rearranging the terms, we have

$$\frac{1}{E} \sum_{s=1}^E \|\nabla f(\bar{x}_{t-1,s-1})\|^2 \leq \frac{2\mathbb{E}[f(x_{t-1}) - f(x_t)]}{E\gamma} + \frac{\gamma\zeta^2 L_{\max}}{n} + \frac{(E-1)(2E-1)\gamma^2 L_{\max}^2 G^2}{3} \quad (33a)$$

$$+ \frac{2^{2z} E^{2z} \sqrt{Q_z + G^{4z+2}} G}{(2z+1)\sigma^{2z}} + \frac{\gamma 2^{4z} E^{4z+1} (Q_z + G^{4z+2}) L_{\max}}{2(2z+1)^2 \sigma^{4z}} \quad (33b)$$

$$+ \frac{4\eta_z^2 \gamma \sigma^2 \sum_{j=1}^d L_j}{En}. \quad (33c)$$

521 Form the telescopic sum

$$\mathbb{E}\left[\frac{1}{TE} \sum_{t=1}^T \sum_{s=1}^E \|\nabla f(\bar{x}_{t-1,s-1})\|^2\right] \leq \frac{2\mathbb{E}[f(x_0) - f^*]}{TE\gamma} + \frac{\gamma\zeta^2 L_{\max}}{n} + \frac{(E-1)(2E-1)\gamma^2 L_{\max}^2 G^2}{3} \quad (34a)$$

$$+ \frac{2^{2z} E^{2z} \sqrt{Q_z + G^{4z+2}} G}{(2z+1)\sigma^{2z}} + \frac{\gamma 2^{4z} E^{4z+1} (Q_z + G^{4z+2}) L_{\max}}{2(2z+1)^2 \sigma^{4z}} \quad (34b)$$

$$+ \frac{4\eta_z^2 \gamma \sigma^2 \sum_{j=1}^d L_j}{En}. \quad (34c)$$

522

□

523 *Proof of Lemma 4.* Therefore, from smoothness we have,

$$f(x_t) - f(\bar{x}_t) \leq \langle \nabla f(\bar{x}_t), x_t - \bar{x}_t \rangle + \frac{\sum_{j=1}^d L_j (x_t(j) - \bar{x}_t(j))^2}{2}. \quad (35)$$

524 The following equation and inequality can be checked, where the expectation is taken over the noise
525 vector ξ_z ,

$$x_t - \bar{x}_t = \frac{\gamma}{n} \sum_{i=1}^n \left(\eta_z \sigma \text{Sign} \left(\sum_{s=1}^E g_{t,s}^i + \sigma \xi_z \right) - \sum_{s=1}^E g_{t,s}^i \right), \quad (36)$$

$$\mathbb{E}[x_t - \bar{x}_t] = \frac{\gamma}{n} \sum_{i=1}^n \left(\sigma \Psi_z \left(\frac{1}{\sigma} \sum_{s=1}^E g_{t,s}^i \right) - \sum_{s=1}^E g_{t,s}^i \right). \quad (37)$$

526 For any $j = 1, \dots, d$, we have

$$\mathbb{E}[(x_t(j) - \bar{x}_t(j))^2] \leq \frac{\gamma^2}{n^2} \left(\sum_{i=1}^n \left(\sigma \Psi_z \left(\frac{1}{\sigma} \sum_{s=1}^E g_{t,s}^i(j) \right) - \sum_{s=1}^E g_{t,s}^i(j) \right) \right)^2 \quad (38a)$$

$$+ \frac{\gamma^2}{n^2} \mathbb{E} \left[\left(\sum_{i=1}^n \left(\eta_z \sigma \text{Sign} \left(\sum_{s=1}^E g_{t,s}^i(j) + \sigma \xi_z \right) - \sigma \Psi_z \left(\frac{1}{\sigma} \sum_{s=1}^E g_{t,s}^i(j) \right) \right) \right)^2 \right] \quad (38b)$$

$$\leq \frac{\gamma^2}{n} \sum_{i=1}^n \left(\sigma \Psi_z \left(\frac{1}{\sigma} \sum_{s=1}^E g_{t,s}^i(j) \right) - \sum_{s=1}^E g_{t,s}^i(j) \right)^2 \quad (38c)$$

$$+ \frac{\gamma^2}{n^2} \sum_{i=1}^n \mathbb{E} \left[\left(\eta_z \sigma \text{Sign} \left(\sum_{s=1}^E g_{t,s}^i(j) + \sigma \xi_z \right) - \sigma \Psi_z \left(\frac{1}{\sigma} \sum_{s=1}^E g_{t,s}^i(j) \right) \right)^2 \right] \quad (38d)$$

$$\leq \frac{\gamma^2}{n} \sum_{i=1}^n \left(\sigma \Psi_z \left(\frac{1}{\sigma} \sum_{s=1}^E g_{t,s}^i(j) \right) - \sum_{s=1}^E g_{t,s}^i(j) \right)^2 \quad (38e)$$

$$+ \frac{2\gamma^2}{n^2} \sum_{i=1}^n \mathbb{E} \left[\left(\eta_z \sigma \text{Sign} \left(\sum_{s=1}^E g_{t,s}^i(j) + \sigma \xi_z \right) \right)^2 \right] \quad (38f)$$

$$+ \frac{2\gamma^2}{n^2} \sum_{i=1}^n \left(\sigma \Psi_z \left(\frac{1}{\sigma} \sum_{s=1}^E g_{t,s}^i(j) \right) \right)^2 \quad (38g)$$

$$\leq \frac{\gamma^2}{n} \sum_{i=1}^n \left(\sigma \Psi_z \left(\frac{1}{\sigma} \sum_{s=1}^E g_{t,s}^i(j) \right) - \sum_{s=1}^E g_{t,s}^i(j) \right)^2 + \frac{4\eta_z^2 \gamma^2 \sigma^2}{n}. \quad (38h)$$

$$(38i)$$

527 Therefore, from Lemma 1, we have

$$\mathbb{E} \left[\sum_{j=1}^d L_j (x_t(j) - \bar{x}_t(j))^2 \right] \leq \frac{\gamma^2}{n} \sum_{i=1}^n \sum_{j=1}^d L_j \left(\sigma \Psi_z \left(\frac{1}{\sigma} \sum_{s=1}^E g_{t,s}^i(j) \right) - \sum_{s=1}^E g_{t,s}^i(j) \right)^2 \quad (39a)$$

$$+ \frac{4\eta_z^2 \gamma^2 \sigma^2 \sum_{j=1}^d L_j}{n} \quad (39b)$$

$$\leq \frac{\gamma^2 L_{\max}}{n} \sum_{i=1}^n \left\| \sigma \Psi_z \left(\frac{1}{\sigma} \sum_{s=1}^E g_{t,s}^i \right) - \sum_{s=1}^E g_{t,s}^i \right\|^2 + \frac{4\eta_z^2 \gamma^2 \sigma^2 \sum_{j=1}^d L_j}{n} \quad (39c)$$

$$\leq \frac{\gamma^2 L_{\max}}{4(2z+1)^2 \sigma^{4z} n} \sum_{i=1}^n \left\| \sum_{s=1}^E g_{t,s}^i \right\|_{4z+2}^{4z+2} + \frac{4\eta_z^2 \gamma^2 \sigma^2 \sum_{j=1}^d L_j}{n}. \quad (39d)$$

528 Now we need to bound

$$\mathbb{E} \left[\left\| \sum_{s=1}^E g_{t,s}^i \right\|_{4z+2}^{4z+2} \right], \quad (40)$$

529 where the expectation is taken over gradient noise.

$$\mathbb{E} \left[\left\| \sum_{s=1}^E g_{t,s}^i \right\|_{4z+2}^{4z+2} \right] \leq \mathbb{E} \left[E^{4z+1} \sum_{s=1}^E \|g_{t,s}^i\|_{4z+2}^{4z+2} \right] \quad (41a)$$

$$= \mathbb{E} \left[E^{4z+1} \sum_{s=1}^E \|g_{t,s}^i - \nabla f_i(x_{t,s-1}^i) + \nabla f_i(x_{t,s-1}^i)\|_{4z+2}^{4z+2} \right] \quad (41b)$$

$$\leq \mathbb{E} \left[(2E)^{4z+1} \sum_{s=1}^E \|g_{t,s}^i - \nabla f_i(x_{t,s-1}^i)\|_{4z+2}^{4z+2} + (2E)^{4z+1} \sum_{s=1}^E \|\nabla f_i(x_{t,s-1}^i)\|_{4z+2}^{4z+2} \right] \quad (41c)$$

$$\leq (2E)^{4z+1} E Q_z + (2E)^{4z+1} \sum_{s=1}^E \|\nabla f_i(x_{t,s-1}^i)\|_2^{4z+2} \leq 2^{4z+1} E^{4z+2} (Q_z + G^{4z+2}). \quad (41d)$$

530 In the derivation above, we use a classical result on the monotonicity of ℓ_p norm: For any $x \in \mathbb{R}^d$
531 and $1 < r < p$, we have

$$\|x\|_p \leq \|x\|_r \leq d^{\frac{1}{r} - \frac{1}{p}} \|x\|_p. \quad (42)$$

532 Therefore, by taking expectation over both ξ_z and Gradient noise, we have

$$\mathbb{E} \left[\sum_{j=1}^d L_j (x_t(j) - \bar{x}_t(j))^2 \right] \leq \frac{\gamma^2 2^{4z+1} E^{4z+2} (Q_z + G^{4z+2}) L_{\max}}{4(2z+1)^2 \sigma^{4z}} + \frac{4\eta_z^2 \gamma^2 \sigma^2 \sum_{j=1}^d L_j}{n}. \quad (43)$$

533 Hence, we have

$$\mathbb{E}[f(x_t) - f(\bar{x}_t)] \leq \left\langle \nabla f(\bar{x}_t), \frac{\gamma}{n} \sum_{i=1}^n \left(\sigma \Psi \left(\frac{1}{\sigma} \sum_{s=1}^E g_{t,s}^i \right) - \sum_{s=1}^E g_{t,s}^i \right) \right\rangle + \frac{\sum_{j=1}^d L_j (x_t(j) - \bar{x}_t(j))^2}{2} \quad (44a)$$

$$\leq \|\nabla f(\bar{x}_t)\| \left\| \frac{\gamma}{n} \sum_{i=1}^n \left(\sigma \Psi \left(\frac{1}{\sigma} \sum_{s=1}^E g_{t,s}^i \right) - \sum_{s=1}^E g_{t,s}^i \right) \right\| + \frac{\sum_{j=1}^d L_j (x_t(j) - \bar{x}_t(j))^2}{2} \quad (44b)$$

$$\leq \frac{\gamma^{2z} E^{2z+1} \sqrt{Q_z + G^{4z+2}} G}{2(2z+1) \sigma^{2z}} + \frac{\gamma^2 2^{4z} E^{4z+2} (Q_z + G^{4z+2}) L_{\max}}{4(2z+1)^2 \sigma^{4z}} + \frac{2\eta_z^2 \gamma^2 \sigma^2 \sum_{j=1}^d L_j}{n}. \quad (44c)$$

534

□

Proof of Lemma 5.

$$f(\bar{x}_t) - f(x_{t-1}) = f(\bar{x}_{t-1,E}) - f(\bar{x}_{t-1,0}) = \sum_{s=1}^E f(\bar{x}_{t-1,s}) - f(\bar{x}_{t-1,s-1}) \quad (45a)$$

$$\leq \sum_{s=1}^E \left(-\langle \nabla f(\bar{x}_{t-1,s-1}), \bar{x}_{t-1,s-1} - \bar{x}_{t-1,s} \rangle + \frac{L_{\max}}{2} \|\bar{x}_{t-1,s} - \bar{x}_{t-1,s-1}\|^2 \right) \quad (45b)$$

$$= \sum_{s=1}^E \left(-\gamma \langle \nabla f(\bar{x}_{t-1,s-1}), \frac{1}{n} \sum_{i=1}^n g_{t-1,s}^i \rangle + \frac{\gamma^2 L_{\max}}{2} \left\| \frac{1}{n} \sum_{i=1}^n g_{t-1,s}^i \right\|^2 \right). \quad (45c)$$

535 Taking expectation over gradient noise, we have

$$\mathbb{E}[\|\frac{1}{n} \sum_{i=1}^n g_{t-1,s}^i\|^2] \leq \|\frac{1}{n} \sum_{i=1}^n \nabla f_i(x_{t-1,s-1}^i)\|^2 + \frac{\zeta^2}{n}, \quad (46a)$$

$$\mathbb{E}[-\langle \nabla f(\bar{x}_{t-1,s-1}), \frac{1}{n} \sum_{i=1}^n g_{t-1,s}^i \rangle] = -\langle \nabla f(\bar{x}_{t-1,s-1}), \frac{1}{n} \sum_{i=1}^n \nabla f_i(x_{t-1,s-1}^i) \rangle \quad (46b)$$

$$= -\frac{1}{2} \|\nabla f(\bar{x}_{t-1,s-1})\|^2 - \frac{1}{2} \|\frac{1}{n} \sum_{i=1}^n \nabla f_i(x_{t-1,s-1}^i)\|^2 \quad (46c)$$

$$+ \frac{1}{2} \|\nabla f(\bar{x}_{t-1,s-1}) - \frac{1}{n} \sum_{i=1}^n \nabla f_i(x_{t-1,s-1}^i)\|^2. \quad (46d)$$

536 Notice that from smoothness, we have for arbitrary $x, y \in \mathbb{R}^d$,

$$f(y) \leq \langle \nabla f(x), y - x \rangle + \frac{L_{\max}}{2} \|y - x\|^2, \quad (47)$$

537 which is equivalent to

$$\|\nabla f(x) - \nabla f(y)\| \leq L_{\max} \|y - x\|. \quad (48)$$

538 Now for every s , we have

$$\|\nabla f(\bar{x}_{t-1,s-1}) - \frac{1}{n} \sum_{i=1}^n \nabla f_i(x_{t-1,s-1}^i)\|^2 \quad (49a)$$

$$= \|\frac{1}{n} \sum_{i=1}^n \nabla f_i(\bar{x}_{t-1,s-1}) - \frac{1}{n} \sum_{i=1}^n \nabla f_i(x_{t-1,s-1}^i)\|^2 \quad (49b)$$

$$\leq \frac{L^2}{n} \sum_{i=1}^n \|\bar{x}_{t-1,s-1} - x_{t-1,s-1}^i\|^2 \quad (49c)$$

$$= \frac{\gamma^2 L_{\max}^2}{n} \sum_{i=1}^n \left\| \sum_{q=1}^{s-1} \left(\frac{1}{n} \sum_{j=1}^n g_{t-1,q}^j - g_{t-1,q}^i \right) \right\|^2 \quad (49d)$$

$$\leq \frac{(s-1)\gamma^2 L_{\max}^2}{n} \sum_{i=1}^n \sum_{q=1}^{s-1} \left\| \frac{1}{n} \sum_{j=1}^n g_{t-1,q}^j - g_{t-1,q}^i \right\|^2 \quad (49e)$$

$$\leq 2(s-1)^2 \gamma^2 L_{\max}^2 G^2. \quad (49f)$$

539 In all, we have

$$\mathbb{E}[f(\bar{x}_t) - f(x_{t-1})] \leq \sum_{s=1}^E \left(-\frac{\gamma}{2} \|\nabla f(\bar{x}_{t-1,s-1})\|^2 - \frac{\gamma}{2} \|\frac{1}{n} \sum_{i=1}^n \nabla f_i(x_{t-1,s-1}^i)\|^2 + \frac{\gamma^2 \zeta^2 L_{\max}}{2n} \right) \quad (50a)$$

$$+ \frac{\gamma}{2} \|\nabla f(\bar{x}_{t-1,s-1}) - \frac{1}{n} \sum_{i=1}^n \nabla f_i(x_{t-1,s-1}^i)\|^2 + \frac{\gamma^2 L_{\max}}{2} \|\frac{1}{n} \sum_{i=1}^n \nabla f_i(x_{t-1,s-1}^i)\|^2 \quad (50b)$$

$$\leq \sum_{s=1}^E \left(-\frac{\gamma}{2} \|\nabla f(\bar{x}_{t-1,s-1})\|^2 + \frac{\gamma^2 \zeta^2 L_{\max}}{2n} + (s-1)^2 \gamma^3 L_{\max}^2 G^2 \right) \quad (50c)$$

$$= -\frac{\gamma}{2} \sum_{s=1}^E \|\nabla f(\bar{x}_{t-1,s-1})\|^2 + \frac{E\gamma^2 \zeta^2 L_{\max}}{2n} + \frac{E(E-1)(2E-1)\gamma^3 L_{\max}^2 G^2}{6}. \quad (50d)$$

540

□

541 *Proof of Theorem 3.* We need a similar lemma like Lemma 4.

542 **Lemma 6.** If $\sigma > E(G + Q_\infty)$, then

$$\mathbb{E}[f(x_t) - f(\bar{x}_t)] \leq \frac{2\gamma^2\sigma^2 \sum_{j=1}^d L_j}{n}. \quad (51)$$

543 Following similar idea in the proof of Theorem 1, we have

$$\mathbb{E}[f(x_t) - f(x_{t-1})] = \mathbb{E}[f(x_t) - f(\bar{x}_t)] + \mathbb{E}[f(\bar{x}_t) - f(x_{t-1})] \quad (52a)$$

$$\leq -\frac{\gamma}{2} \sum_{s=1}^E \|\nabla f(\bar{x}_{t-1, s-1})\|^2 + \frac{E\gamma^2\zeta^2 L_{\max}}{2n} + \frac{E(E-1)(2E-1)\gamma^3 L_{\max}^2 G^2}{6} + \frac{2\gamma^2\sigma^2 \sum_{j=1}^d L_j}{n}. \quad (52b)$$

544 Rearranging the terms, we have

$$\frac{1}{E} \sum_{s=1}^E \|\nabla f(\bar{x}_{t-1, s-1})\|^2 \leq \frac{2\mathbb{E}[f(x_{t-1}) - f(x_t)]}{E\gamma} + \frac{\gamma\zeta^2 L_{\max}}{n} + \frac{(E-1)(2E-1)\gamma^2 L_{\max}^2 G^2}{3} \quad (53a)$$

$$+ \frac{4\gamma\sigma^2 \sum_{j=1}^d L_j}{En}. \quad (53b)$$

545 Form the telescopic sum

$$\mathbb{E}\left[\frac{1}{TE} \sum_{t=1}^T \sum_{s=1}^E \|\nabla f(\bar{x}_{t-1, s-1})\|^2\right] \leq \frac{2\mathbb{E}[f(x_0) - f^*]}{TE\gamma} + \frac{\gamma\zeta^2 L_{\max}}{n} + \frac{(E-1)(2E-1)\gamma^2 L_{\max}^2 G^2}{3} \quad (54a)$$

$$+ \frac{4\gamma\sigma^2 \sum_{j=1}^d L_j}{En}. \quad (54b)$$

546 Here we provide a simple example where $\sigma < E(G + Q_\infty)$ and the algorithm cannot converge.

Consider $E = 1$, $Q_\infty = 0$ and the problem

$$\min_{x \in \mathbb{R}} (x - A)^2 + (x + A)^2,$$

547 where $A > 0$ is some positive number. If we choose the initial to be $x_0 = \frac{A}{2}$. As we can, the gradient
548 at x_0 for the two parts of the objective function are $-A$ and $3A$ respectively. We denote that ξ_∞ as
549 the random noise following uniform distribution at $[-1, 1]$. If now $\sigma < A$, we have

$$\text{Sign}(-A + \sigma\xi_\infty) + \text{Sign}(3A + \sigma\xi_\infty) = 0, \quad (55)$$

550 i.e., this algorithm never update the variable. \square

551 *Proof of Lemma 6.* We first note that, when $z = +\infty$, we have

$$\Psi_\infty(x) = \begin{cases} x & x \in [-1, 1], \\ 1 & x < -1, \\ 1 & x > 1. \end{cases} \quad (56)$$

552 Now, from L -smoothness we have,

$$f(x_t) - f(\bar{x}_t) \leq \langle \nabla f(\bar{x}_t), x_t - \bar{x}_t \rangle + \frac{\sum_{j=1}^d L_j (x_t(j) - \bar{x}_t(j))^2}{2}. \quad (57a)$$

553 The following equation and inequality can be checked, where the expectation is taken over ξ_∞ ,

$$\mathbb{E}[x_t - \bar{x}_t] = \mathbb{E}\left[\frac{\gamma}{n} \sum_{i=1}^n \left(\sigma \text{Sign}\left(\sum_{s=1}^E g_{t,s}^i + \sigma\xi_\infty\right) - \sum_{s=1}^E g_{t,s}^i \right)\right] = 0, \quad (58)$$

554 because from the condition of σ we can see that $\sigma > \|\sum_{s=1}^E g_{t,s}^i\|_\infty$ almost surely. \square

555 For any $j = 1, \dots, d$, we have

$$\mathbb{E}[(x_t(j) - \bar{x}_t(j))^2] \leq \frac{\gamma^2}{n^2} \mathbb{E} \left[\left(\sum_{i=1}^n \left(\sigma \text{Sign} \left(\sum_{s=1}^E g_{t,s}^i(j) + \sigma \xi_\infty(j) \right) - \sigma \Psi_\infty \left(\frac{1}{\sigma} \sum_{s=1}^E g_{t,s}^i(j) \right) \right) \right)^2 \right] \quad (59)$$

$$\leq \frac{4\gamma^2\sigma^2}{n}. \quad (60)$$

556 Hence, we have

$$\mathbb{E}[f(x_t) - f(\bar{x}_t)] \leq + \frac{\sum_{j=1}^d L_j (x_t(j) - \bar{x}_t(j))^2}{2} \quad (61)$$

$$\leq \frac{2\gamma^2\sigma^2 \sum_{j=1}^d L_j}{n}. \quad (62)$$

557 D A simple simulated experiment.

558 In this section, we verify our theoretical results in Section 3 on a simple simulated experiment without
 559 any gradient noise. Specifically, we consider the following distributed optimization problem with 10
 560 clients,

$$\min_{x \in \mathbb{R}^d} \frac{1}{2} \sum_{i=1}^{10} \|x - y_i\|^2. \quad (63)$$

561 Here we generate $y_1, \dots, y_{10} \in \mathbb{R}^d$ using i.i.d standard Gaussian distribution, where d is the problem
 562 dimension. We compare the performance of the following algorithms. For all the algorithms, we use
 563 the same stepsize 0.01 and all-zero initialization. We denote the tested algorithms as:

- 564 • GD: Distributed gradient descent without any compression.
- 565 • Sto-SignSGD: The algorithm proposed by [Safaryan and Richtárik, 2021].
- 566 • SignSGD: (Algorithm 1 with $z = 1$, $E = 1$ and $\sigma = 0$.)
- 567 • 1-SignSGD (Algorithm 1 with $z = 1$ and $E = 1$.)
- 568 • ∞ -SignSGD (Algorithm 1 with $z = +\infty$ and $E = 1$.)

569 **Results.** As we can see from Figure 3, all the stochastic sign-based algorithms can converge to
 570 the optimal solution, while the SignSGD without any noise fail to converge to the optimal solution.
 571 Besides, 1-SignSGD and ∞ -SignSGD have roughly the same convergence speed which is slightly
 572 slower than the uncompressed gradient descent. It is also verified that the input-dependent noise
 573 scale adopted by [Safaryan and Richtárik, 2021] could lead to slow convergence when the problem
 574 dimension is high, as we have discussed in Section 3.2. The optimal noise scales of 1-SignSGD and
 575 ∞ -SignSGD are selected based on Figure 4. We can see that there is a clear bias-variance trade-off
 576 in 4 which corroborates our prediction in Section 3. Moreover, it worth to mention that in this
 577 experiment, the optimal σ for ∞ -SignSGD is much smaller than the conservative choice suggested
 578 by theory.

579 E Experiment details

580 E.1 Details for the experiment in Section 4.1

581 In Table 2, we provide the tuned hyperparameters for all the tested algorithms on non-i.i.d
 582 MNIST. Generally, we tune the hyperparameters via grid search: $[0.1, 0.05, 0.01, 0.005]$ for stepsize,
 583 $[0, 0.3, 0.5, 0.7, 0.9]$ for momentum coefficient, $[0, 0.02, 0.05, 0.01, 0.03, 0.05, 0.1, 0.3, 0.5]$ for noise
 584 scale.

585 In Figure 5, we visualize the performance of 1-SignSGD and ∞ -SignSGD under different noise
 586 scales. As we can see, the results for 1-SignSGD and ∞ -SignSGD are almost the same, except that
 587 the ∞ -SignSGD is slightly better than 1-SignSGD when the noise scale is large.

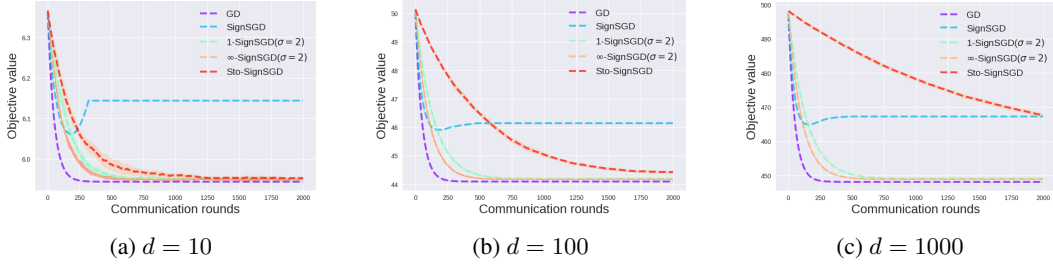


Figure 3: Performance of algorithms under different problem dimension.

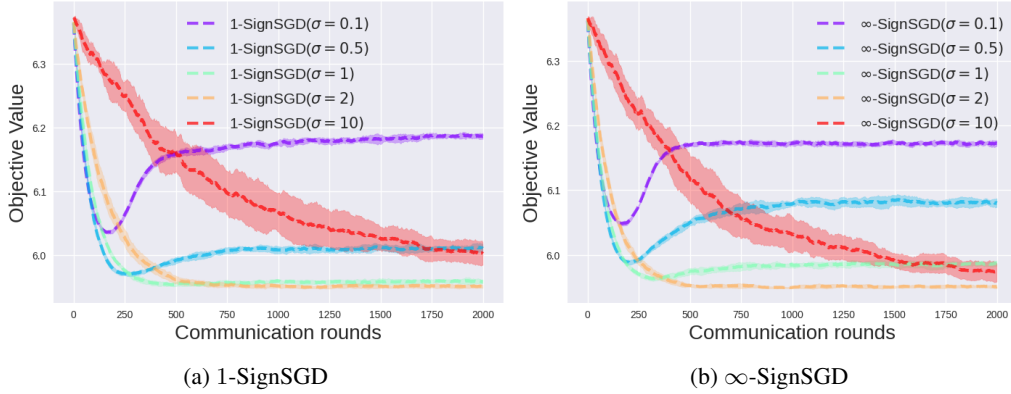


Figure 4: Algorithm 1 under different noise scales

588 **E.2 Details for the experiment in Section 4.2**

589 For the experiment on EMNIST, we fix the client stepsize as 0.05. Then we tune
 590 the server stepsize, noise scales via grid search: $[1, 0.5, 0.1, 0.05, 0.01, 0.005]$ for stepsize,
 591 $[0, 0.005, 0.02, 0.05, 0.01, 0.03, 0.05, 0.1, 0.2]$ for noise scale. The used hyperparameter in the Figure
 592 2 are summarized in Table 3. We also visualize the performance of 1-SignFedAvg and ∞ -SignFedAvg
 593 under various noise scales and local steps in Figure 6, 7, where we use SignFedAvg to represent
 594 Algorithm 1 with $\sigma = 0$.

Algorithm	Stepsize	Momentum coefficient	Noise scale
SGDwM	0.05	0.9	
EF-SignSGDwM	0.05	0.9	
Sto-SignSGDwM	0.01	0.9	
SignSGD	0.01	0	0
1-SignSGD	0.01	0	0.05
∞ -SignSGD	0.01	0	0.05

Table 2: Hyperparameters for tested Algorithms on non-i.i.d MNIST.

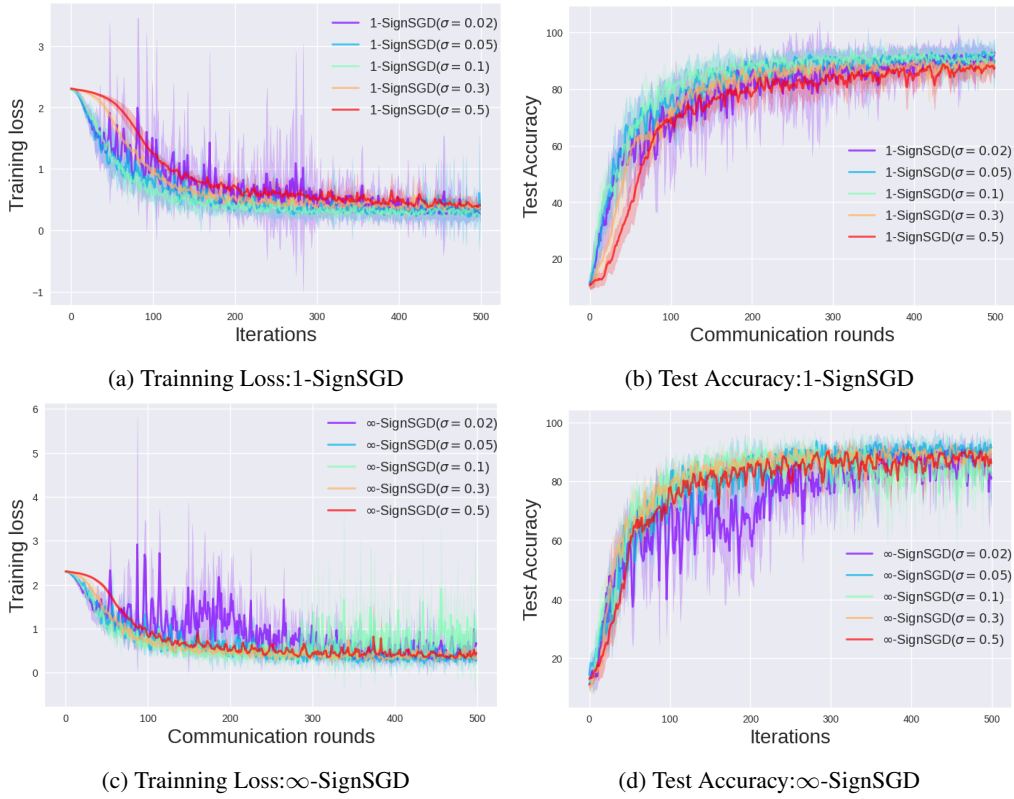


Figure 5: ALG 1 under different noise scales on non-i.i.d MNIST

Algorithm	Server stepsize	Noise scale
1-SignFedAvg	0.03	0.01
∞ -SignFedAvg	0.03	0.01

Table 3: Hyperparameters for tested Algorithms on EMNIST.

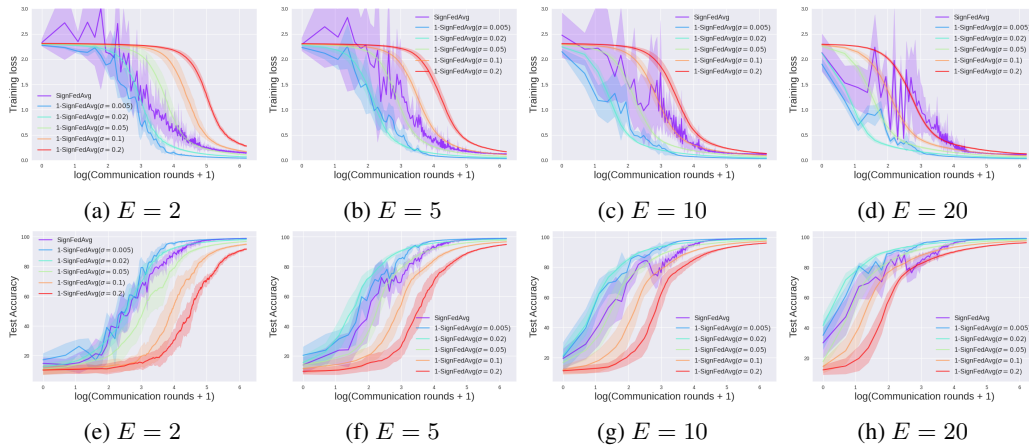


Figure 6: 1-SignFedAvg under different noise scales and local steps

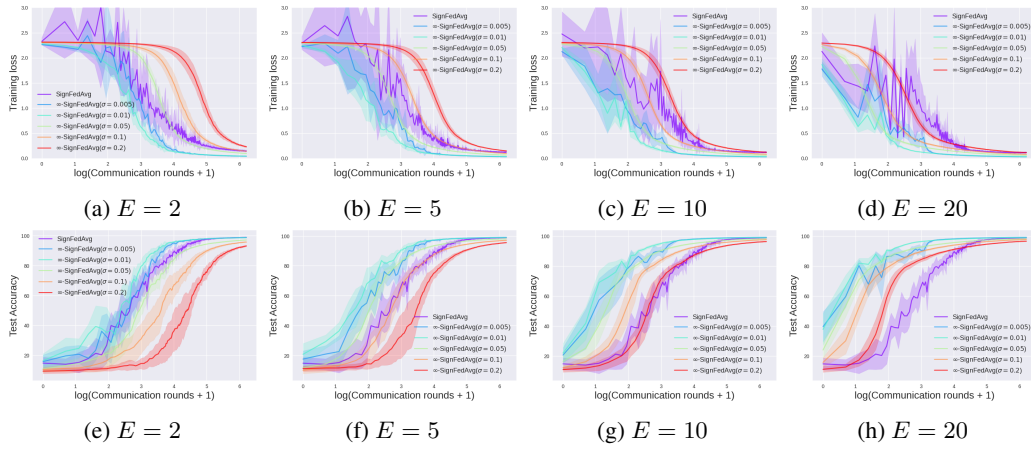


Figure 7: ∞ -SignFedAvg under different noise scales and local steps