DECENTRALIZED NONCONVEX OPTIMIZATION UNDER HEAVY-TAILED NOISE: NORMALIZATION AND OPTIMAL CONVERGENCE

Anonymous authors

000

001

002

004

006

008 009 010

011 012 013

014

015

016

018

019

021

023

024

025 026

028

029

035

037

040

041

042

043

044

046

047

048

051

052

Paper under double-blind review

ABSTRACT

Heavy-tailed noise in nonconvex stochastic optimization has garnered increasing research interest, as empirical studies, including those on training attention models, suggest it is a more realistic gradient noise condition. This paper studies first-order nonconvex stochastic optimization under heavy-tailed gradient noise in a decentralized setup, where each node can only communicate with its direct neighbors in a predefined graph. Specifically, we consider a class of heavytailed gradient noise that is zero-mean and has only p-th moment for $p \in (1,2]$. We propose GT-NSGDm, Gradient Tracking based Normalized Stochastic Gradient Descent with momentum, that utilizes normalization, in conjunction with gradient tracking and momentum, to cope with heavy-tailed noise on distributed nodes. We show that, when the communication graph admits primitive and doubly stochastic weights, GT-NSGDm guarantees, for the first time in the literature, that the expected gradient norm converges at an optimal non-asymptotic rate $O(1/T^{(p-1)/(3p-2)})$, which matches the lower bound in the centralized setup. When tail index p is unknown, GT-NSGDm attains a non-asymptotic rate $O(1/T^{(p-1)/(2p)})$ that is, for p < 2, topology independent and has a speedup factor $n^{1-1/p}$ in terms of the number of nodes n. Finally, experiments on nonconvex linear regression with tokenized synthetic data and decentralized training of language models on a real-world corpus demonstrate that GT-NSGDm is more robust and efficient than baselines.

1 Introduction

In this paper, we address the problem of nonconvex stochastic optimization under heavy-tailed gradient noise in the decentralized setup. Consider a graph with n nodes connected by a predefined topology $\mathcal{G}:=(\mathcal{V},\mathcal{E})$, where $\mathcal{V}:=\{1,\ldots,n\}$ is the set of node indices, and \mathcal{E} is the collection of directed pairs $(i,r), i,r\in\mathcal{V}$ such that node i can send information to the neighboring node r. Each node $i\in\mathcal{V}$ holds a local nonconvex differentiable cost function $f_i:\mathbb{R}^d\to\mathbb{R}$, and can access its stochastic gradient, subject to zero mean noise with a bounded p-th moment for some $p\in(1,2]$. Cooperatively, these nodes aim to solve $\min_{\boldsymbol{x}\in\mathbb{R}^d}f(\boldsymbol{x}):=(1/n)\sum_{i=1}^n f_i(\boldsymbol{x})$, through local computation and peer-to-peer communication only with their immediate neighbors.

Decentralized optimization in the above formulation has been studied for decades (Tsitsiklis et al., 1986), and has recently attracted growing research interest due to its advantages in scalability and privacy preservation across a wide range of distributed machine learning, signal processing, and control tasks over networks (Nedić et al., 2018; Li et al., 2020; Kairouz et al., 2021). For instance, in privacy-sensitive applications such as those in the medical domain (Brisimi et al., 2018), training data are often distributed across n nodes due to privacy constraints. In such cases, each f_i represents an empirical risk function, e.g., a neural network, defined over the local dataset on node i, and all nodes collaboratively train a global predictive model via peer-to-peer communication without sharing raw data. Moreover, decentralized optimization is also employed in data centers to reduce communication bottlenecks associated with the central node in traditional centralized training paradigms (Lian et al., 2017).

056

057

058

060

061

062

063

064

065

066

067

068

069

071

072

073

074

075

076

077

079

081

082

083

084

085

087

880

089

090

091

092

093

094

095 096

097

098 099

100

101 102

103

104 105

106

107

In decentralized optimization, first-order methods are widely favored for their simplicity and scalability (Xin et al., 2020b). However, computing the exact gradient of each local objective function f_i at every iteration can be computationally expensive, particularly in large-scale settings where each node holds a substantial volume of local data. To alleviate this computational burden, decentralized stochastic gradient methods, which approximate exact gradients, have been extensively studied. Most existing approaches, including decentralized stochastic (sub)gradient descent (Sundhar Ram et al., 2010; Koloskova et al., 2020; Wang & Joshi, 2021), variance reduction techniques (Yuan et al., 2018), and gradient tracking-based schemes (Di Lorenzo & Scutari, 2016; Pu & Nedić, 2021), typically assume that stochastic gradient noise has a *finite variance*. Nevertheless, recent empirical and theoretical evidence indicates that, when optimizing certain neural network architectures, especially attention-based models such as Transformers (Vaswani et al., 2017), the gradient noise often follows a heavy-tailed distribution¹ with significantly large or even infinite variance (Simsekli et al., 2019; Zhang et al., 2020; Gorbunov et al., 2020; Gurbuzbalaban et al., 2021; Ahn et al., 2024; Kunstner et al., 2024). The presence of heavy-tailed gradient noise poses substantial challenges for existing methods. Empirically, some stochastic gradient descent (SGD) based methods can suffer from instability and even dramatic drop of training accuracies (Zhang et al., 2020; Charles et al., 2021; Yang et al., 2022), particularly in distributed large-cohort training. Theoretically, unbounded variance renders many established analyses invalid, and in centralized settings it necessitates the use of nonlinear adaptive techniques such as clipping, sign, and normalization (Zhang et al., 2020; Sadiev et al., 2023; Compagnoni et al., 2025b; Hübler et al., 2024; Liu & Zhou, 2025; Armacki et al., 2025) to combat the strong noise. However, incorporating such adaptive strategies in decentralized algorithms introduces inherent nonlinearity into the algorithmic dynamics associated with the average-sum structured function f, making the design and analysis of decentralized algorithms under heavy-tailed noise significantly more challenging.

Decentralized optimization under heavy-tailed gradient noise remains underexplored. To the best of our knowledge, only recent studies Sun & Chen (2024); Yu et al. (2023) have attempted to address this problem under restrictive assumptions. Specifically, Sun & Chen (2024) considers zero-mean gradient noise with bounded p-th central moment $(p \in (1,2])$ similar to our setting but assumes a compact domain or bounded gradients. Their proposed decentralized gradient descent method with ℓ_2 gradient clipping achieves almost sure convergence for strongly convex local functions. However, the restrictive compact domain or gradients assumption in Sun & Chen (2024) limits its practical applicability, and the convergence rate is not explicitly provided. Another work, Yu et al. (2023), also assumes strongly convex local objectives and develops a decentralized gradient method with smoothed clipping and error feedback under gradient noise that is zero-mean, symmetric, and has bounded first absolute moment, showing an *in-expectation* convergence rate of $1/t^{\delta}$ for some $\delta \in (0, 2/5)$. Although the noise assumption in Yu et al. (2023) is weaker than ours (as it requires only a first-moment bound), the additional assumptions of noise symmetry and the dependence of the rate exponent δ on both the problem dimension and condition number restrict its general applicability. Moreover, both works Sun & Chen (2024); Yu et al. (2023) assume strong convexity, whereas many practical optimization problems involving heavy-tailed noise, particularly in modern machine learning, are inherently nonconvex. Further, the convergence rates in Sun & Chen (2024); Yu et al. (2023) are either unclear or sub-optimal, even compared to the optimal iteration complexity bound $O(1/T^{(p-1)/(3p-2)})$ for general nonconvex functions. In this work, we relax these restrictive assumptions and address the following question:

Can we design a decentralized algorithm for **nonconvex** optimization under general zero-mean gradient noise with only a finite p-th moment for $p \in (1, 2]$ with **optimal iteration complexity**?

1.1 CONTRIBUTIONS

We answer this question affirmatively through the following key contributions:

• We develop a decentralized method, called GT-NSGDm, using normalization, coupled with momentum variance reduction, to combat heavy-tailed noise, and using gradient tracking

 $^{^1}$ A random variable X is called heavy-tailed if it exhibits a heavier tail than any exponential distribution; formally, for any constant a>0, $\limsup_{x\to\infty}\mathbb{P}(X>x)e^{ax}=\infty$ (Nair et al., 2022). While some heavy-tailed distributions, such as log-normal and Weibull, still have bounded variance, this paper also considers the sub-class of heavy tailed gradient noise that may have unbounded (infinite) variance such as α -stable noise.

to handle cross-node heterogeneity. To further shed light on the design of GT-NSGDm, we provide a negative result for a vanilla variant of normalized decentralized SGD that employs no gradient tracking nor momentum.

- For general nonconvex and smooth local functions f_i 's that are bounded from below, we show that GT-NSGDm converges in expectation at a rate $O(1/T^{(p-1)/(3p-2)})$, which matches the lower bound in centralized setting and is order-optimal. Our convergence rate significantly improves upon related works (Sun & Chen, 2024; Yu et al., 2023), which assume strong convexity and lack an explicit rate exponent.
- When the tail index p is unknown, GT-NSGDm achieves a rate of $O(1/T^{(p-1)/(2p)})$, matching the best-known rate in the centralized setting without requiring knowledge of p. Notably, for $p \in (1,2)$ and sufficiently large T, this rate is *independent* of the network topology and exhibits a *speedup* in the number of nodes, with a factor of $n^{1-1/p}$.
- We test our theoretical findings in nonconvex linear regression models on a synthetic dataset that is built to simulate language tokens under controlled heavy-tailed noise injections. We also test GT-NSGDm on distributed training of decoder-only Transformer models on Multi30k datasets (Elliott et al., 2016). Experiments on multiple variants of network topologies show that GT-NSGDm is more robust to injected and empirical heavy-tailed noise and converges faster.

1.2 RELATED WORK

108

110

111

112

113

114

115

116

117

118

119

121

122

123

124

126 127

128 129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

156

157

159

161

Heavy-tailed gradient noise. Recent empirical studies suggest that the distribution of gradient noise in training various deep learning models resembles heavy-tailed distributions, such as Lévy's α -stable distribution (Simsekli et al., 2019; Zhang et al., 2020; Barsbey et al., 2021; Battash et al., 2024). For instance, the work Zhang et al. (2020) demonstrates that the empirical distribution of gradient norm samples during BERT pre-training closely aligns with an α -stable distribution, rather than a Gaussian one (see their Figure 1). The presence of heavy-tailed gradient noise is also supported by theoretical insights (Simsekli et al., 2019; Peluchetti et al., 2020; Gurbuzbalaban et al., 2021; Barsbey et al., 2021). In particular, Simsekli et al. (2019) leverages generalized central limit theorems to show that the gradient noise in SGD can converge to an α -stable random variable.

Adaptive methods. Under heavy-tailed noise, vanilla SGD based methods are shown to suffer from slower convergence or even model collapses in *centralized* settings (Zhang et al., 2020) as well as distributed settings with a central server (Yang et al., 2022; Lee et al., 2025), and adaptive methods such as clipping and normalization are introduced to stabilize training dynamics. In *centralized* settings, the work Zhang et al. (2020) provides lower bounds for both nonconvex and strongly convex smooth functions, showing that SGD with gradient clipping achieves in-expectation upper bounds matching lower bounds. In Sadiev et al. (2023); Liu et al. (2023); Nguyen et al. (2023); Chezhegov et al. (2024), the authors show that when equipped with gradient clipping, SGD, accelerated methods, AdaGrad (Duchi et al., 2011), and Adam (Kingma & Ba, 2014) can achieve (near-)optimal high-probability convergence under various function assumptions. Besides, the work Compagnoni et al. (2025b) shows that signSGD is also robust to heavy-tailed noise through the lens of stochastic differential equations. Further, SGD with gradient normalization, which advantageously requires less hyper-parameter tuning than clipping, is shown to achieve optimal in-expectation convergence (Hübler et al., 2024; Liu & Zhou, 2025; Sun et al., 2024). Our method incorporates the same normalization and variance reduction approach as Liu & Zhou (2025). Notably, in another line of works Jakovetić et al. (2023); Armacki et al. (2025; 2024), the authors conduct a unified convergence analyses for generic nonlinear methods including clipping, sign, and normalization under symmetric noise with positive probability mass around zero without assuming any noise moment bound or only assuming first absolute noise moment bound. In distributed settings with a server, the work Gorbunov et al. (2024) proposes an algorithm that incorporates an error feedback mechanism, wherein clipping is applied to the discrepancy between a local gradient estimator and a stochastic gradient, and establishes optimal high-probability bounds. Moreover, the work Compagnoni et al. (2025a) shows that distributed signSGD converges to an asymptotic neighborhood depending on the 'fatness' of noise tail. When multiple local updates are permitted between communication rounds, the authors of Yang et al. (2022) show that clipping per local step achieves order-optimal in-expectation convergence, albeit under a restrictive bounded gradient assumption. More recently, the work Lee et al. (2025) introduces the TailOPT framework, which adaptively leverages gradient geometry by applying clipping operators during local updates on distributed nodes and utilizing adaptive optimizers for global updates at the server, achieving in-expectation sublinear convergence rates that are independent of the moment parameter p.

Nonlinearities in decentralized optimization. Extending existing methods that are robust to heavy-tailed noise, whether developed for *centralized* settings or *distributed settings with a server*, to decentralized environments is highly nontrivial, primarily due to the *nonlinearities* introduced to peer-to-peer communication. This difficulty is reflected in that existing decentralized methods incorporating nonlinear adaptive techniques for other purposes often impose restrictive conditions (Yu & Kar, 2023; Li & Chi, 2025). For example, to achieve differential privacy through gradient clipping, the work Li & Chi (2025) establishes convergence in decentralized setups under the assumption of either bounded gradient or a stringent similarity condition, namely $\|\nabla f_i(x) - \nabla f(x)\| \le (1/12) \|\nabla f(x)\|$ for all $i \in [n]$ and all x. Similarly, to attain adversarial robustness against gradient attacks, the authors of Yu & Kar (2023) employ gradient clipping with momentum, assuming that all local functions are convex, share a *common minimizer*, and that $\sum_{i=1}^n f_i$ is strongly convex. In this work, we significantly relax these conditions and demonstrate the effective use of nonlinearity (specifically, normalization) in decentralized optimization, thereby motivating broader applications of nonlinear techniques in this setting.

1.3 NOTATION

We denote by \mathbb{N}_+ , \mathbb{R} , \mathbb{R}_+ and \mathbb{R}^d , respectively, the set of positive natural numbers, real numbers, nonnegative real numbers, and the d-dimensional Euclidean space. We use lowercase normal letters for scalars, lowercase boldface letters for vectors, and uppercase boldface letters for matrices. Further, we denote by $\mathbf{1}_k$ and $\mathbf{0}_k$ the all-ones and all-zeros vectors of size k, respectively, and by \mathbf{I}_k the $k \times k$ identity matrix. We let $\|\mathbf{x}\|$ denote the ℓ_2 norm of \mathbf{x} , and $\|\mathbf{A}\|_2$ denote the operator norm of \mathbf{A} . For functions p(t) and q(t) in t, we write p(t) = O(q(t)) if $\limsup_{t \to \infty} p(t)/q(t) < \infty$. Finally, we use $\mathbb E$ to denote expectation over random quantities.

2 PROBLEM FORMULATION

We consider a graph with n nodes, where each node holds a local and private function $f_i : \mathbb{R}^d \to \mathbb{R}$, and the nodes collectively minimize the unconstrained global objective $f(\boldsymbol{x}) := (1/n) \sum_{i=1}^n f_i(\boldsymbol{x})$ through peer-to-peer communication. We now present some standard assumptions on the problem.

Assumption 1 (Finite lower bound). There exists some $f_* := \inf_{x \in \mathbb{R}^d} f(x) > -\infty$.

Assumption 2 (*L*-smoothness). The local function f_i at each node $i \in [n]$ is differentiable and *L*-smooth, i.e., $\forall \boldsymbol{x}, \boldsymbol{y} \in \mathbb{R}^d, \|\nabla f_i(\boldsymbol{x}) - \nabla f_i(\boldsymbol{y})\| \leq L\|\boldsymbol{x} - \boldsymbol{y}\|.$

We next introduce the heavy-tailed noise model. For each node $i \in \mathcal{V}$, at t-th iteration with query \boldsymbol{x}_i^t , the stochastic first-order oracle returns the gradient estimator $g_i(\boldsymbol{x}_i^t, \boldsymbol{\xi}_i^t)$, where $\boldsymbol{\xi}_i^t$ denotes the random sample. Let Ω, \emptyset denote the universe, empty set, respectively. We use the following natural filtration, i.e., an increasing family of sub- σ -algebras, to denote the past history up to iteration t:

$$\mathcal{F}_{-1} := \{\Omega, \emptyset\}, \quad \mathcal{F}_t := \sigma(\{\boldsymbol{\xi}_i^0, \dots, \boldsymbol{\xi}_i^{t-1} : i \in [n]\}), \forall t \ge 0.$$

We then assume this stochastic first-order oracle have the following properties.

Assumption 3 (Heavy-tailed noise). For any \mathcal{F}_t -measurable random vectors $\boldsymbol{x} \in \mathbb{R}^d$, we have the following: $\forall i \in [n], \forall t \geq 0$, (1) $\mathbb{E}[\boldsymbol{g}_i(\boldsymbol{x}, \boldsymbol{\xi}_i^t) \mid \mathcal{F}_t] = \nabla f_i(\boldsymbol{x})$; (2) There exist $p \in (1, 2]$, some constant $\sigma \geq 0$ such that $\mathbb{E}[\|\boldsymbol{g}_i(\boldsymbol{x}, \boldsymbol{\xi}_i^t) - \nabla f_i(\boldsymbol{x})\|^p \mid \mathcal{F}_t] \leq \sigma^p$; (3) The family $\{\boldsymbol{\xi}_i^t : \forall t \geq 0, i \in [n]\}$ of random samples is independent.

Remark 1 (Heavy-tailed distributions). Assumption 3 covers a broad class of heavy-tailed distributions, including Lévy's α -stable distributions, Student's t-distributions, and Pareto distributions. Note that we do not assume noise symmetry as in Yu et al. (2023), and when p=2, Assumption 3 reduces to the standard bounded variance condition commonly assumed in the literature.

Remark 2 (Empirical evidence). Similar to Zhang et al. (2020); Yang et al. (2022), we investigate the empirical distribution of the gradient noise norm $\|g(x,\xi) - \nabla f(x)\|$ in a centralized setting by training a GPT model (Radford et al., 2018) with 3M parameters on the Multi30k dataset (Elliott

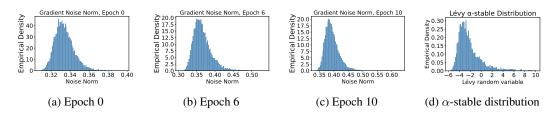


Figure 1: Comparisons of the empirical density of gradient noise norm in different epochs of training a Transformer model with a synthetic Lévy α -stable distribution.

et al., 2016), where g(x) denotes the mini-batch stochastic gradient and $\nabla f(x)$ denotes the full-batch gradient. We train the model for 12 epochs using SGD and plot the empirical density of the noise norm at the beginning of epochs 0, 6, and 10. As shown in Figure 1, as training progresses, the tail of the empirical gradient noise norm distribution becomes heavier (and longer) and increasingly resembles that of a synthetic α -stable distribution.

For peer-to-peer communication in decentralized settings, we need to specify a mixing matrix W on graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$.

Assumption 4 (Weight matrix). The nonnegative weight matrix W, whose (i, r)-th component of W, denoted as w_{ir} , is positive if and only if $(i, r) \in \mathcal{E}$ or i = r, is primitive and doubly stochastic, i.e., $\mathbf{1}_n W = \mathbf{1}_n$ and $W \mathbf{1}_n = \mathbf{1}_n$.

Assumption 4 is standard in the decentralized optimization literature (Xin et al., 2020a), and it guarantees that there exists some nonnegative λ , i.e., spectral gap, such that

$$\|\mathbf{W} - \mathbf{1}_n \mathbf{1}^\top / n\|_2 := \lambda < 1.$$

The assumed weight matrix W can be constructed on undirected and connected graphs (Olshevsky, 2014), and also on some directed and strongly connected graphs that are weight-balanced (Gharesifard & Cortés, 2012). For instance, the family of directed exponential graphs, is weight-balanced and serves as an important topology configuration in decentralized training (Assran et al., 2019).

3 ALGORITHM DEVELOPMENT: GT-NSGDM

We now describe the proposed Algorithm GT-NSGDm and discuss the intuition of its construction. We use \boldsymbol{x}_i^t to denote the estimate of a stationary point for the global cost function f at node i and t-th iteration, and recall that $\boldsymbol{g}_i(\boldsymbol{x}_i^t, \boldsymbol{\xi}_i^t)$ denotes the corresponding stochastic gradient returned from local first-order oracle. Motivated by the error-feedback approach in Yu et al. (2023), which serves as a momentum-type of variance reduction after applying a nonlinear operator to handle heavy-tailed noise, we also employ local momentum variance reduction

$$\boldsymbol{v}_i^t = \beta \boldsymbol{v}_i^{t-1} + (1 - \beta) \boldsymbol{g}_i(\boldsymbol{x}_i^t, \boldsymbol{\xi}_i^t), \tag{1}$$

where $\beta \in [0,1)$ serves as the momentum coefficient. Then, we use gradient tracking (Di Lorenzo & Scutari, 2016) to handle heterogeneous local functions $\{f_i\}_{i=1}^n$. Specifically, we use an estimator y_i^t to track global gradient

$$\mathbf{y}_{i}^{t} = \sum_{r=1}^{n} w_{ir} (\mathbf{y}_{r}^{t-1} + \mathbf{v}_{r}^{t} - \mathbf{v}_{r}^{t-1}).$$
 (2)

It is known that gradient tracking helps eliminate the dependence on heterogeneity among local functions $\{f_i\}_{i=1}^n$, such as the requirement of bounded gradient similarity. Furthermore, similar to the approach in Liu & Zhou (2025), which uses normalization to address heavy-tailed noise in centralized settings, we avoid applying normalization in the recursive updates of the local gradient estimator v_i^t in (1) and the global gradient tracker y_i^t . Instead, normalization is applied only during the update of x_i^t , with step size α , and nonnegative mixing weights $\{w_{ir}\}$ where $w_{ir} > 0$ only when $(i,r) \in \mathcal{E}$ or i=r,

$$\boldsymbol{x}_{i}^{t+1} = \sum_{r=1}^{n} w_{ir} \left(\boldsymbol{x}_{r}^{t} - \alpha \frac{\boldsymbol{y}_{r}^{t}}{\|\boldsymbol{y}_{r}^{t}\|} \right). \tag{3}$$

We combine the local updates (1)(2)(3) on node $i \in \mathcal{V}$ and call it GT-NSGDm, Gradient Tracking based Normalized Stochastic Gradient Descent with momentum. When taking $\beta = 0$, this simplifies to momentum-free gradient tracking with normalization in step 3. However, our analysis shows that GT-NSGDm performs optimal for some $\beta \in (0,1)$, making GT-NSGDm a non-trivial and optimal algorithmic design for the considered problem class. We provide a tabular description for GT-NSGDm in Algorithm 1, where all $\{\boldsymbol{x}_i^0\}$ are initialized from the same point $\bar{\boldsymbol{x}}^0$ for simplicity.

Algorithm 1 GT-NSGDm at each node *i*

Remark 3 (Why vanilla gradient normalization fails?). Although vanilla normalization is successfully used in *centralized* settings to robustify SGD against heavy-tailed noise (Hübler et al., 2024), its direct extension to the *decentralized* settings fails. Suppose we run a vanilla decentralized normalized (noiseless) gradient descent, i.e., in parallel $\forall i \in \mathcal{V}$,

$$\boldsymbol{x}_{i}^{t+1} = \sum_{i=1}^{n} w_{ir} \left(\boldsymbol{x}_{r}^{t} - \alpha \frac{\nabla f_{r}(\boldsymbol{x}_{r}^{t})}{\|\nabla f_{r}(\boldsymbol{x}_{r}^{t})\|} \right). \tag{4}$$

Then, the global average \bar{x}^t would update in the negative direction of the sum of normalized local gradients: $\bar{x}^{t+1} = \bar{x}^t - \frac{\alpha}{n} \sum_{r=1}^n \frac{\nabla f_r(x_r^t)}{\|\nabla f_r(x_r^t)\|}$. Let for some $t, \forall r \in \mathcal{V}, x_r^t = x_* = \arg\min\sum_{i=1}^n f_i(x)$, i.e., all nodes hold the optimal global solution that $\sum_{r=1}^n \nabla f_r(x_*) = 0$. Since $\|\nabla f_r(x_*)\|$ can be different quantities for $r=1,\ldots,n$, due to function heterogeneity, then \bar{x}^{t+1} will move away from x_* . Therefore, vanilla gradient normalization adds some intrinsic errors from heterogeneous local normalizations. By incorporating gradient tracking, we expect that y_r^t would converge to its global average \bar{y}^t , and \bar{y}^t would converge to $(1/n)\sum_{i=r}^n \nabla f_i(x_r^t)$. In this way, x_r^t would move along the direction of the normalized sum of local gradients, and thus emulating the centralized setting.

In the following claim, we further demonstrate that vanilla gradient normalization can cause the iterates x_i^t to remain arbitrarily far from the optimal solution (see Appendix A for a proof).

Claim 1. Consider algorithm 4. For any even n, for any $B \geq 1$, there exist $\{f_i\}_{i=1}^n$ satisfying Assumptions 1-2, a gradient oracle satisfying Assumption 3, a mixing matrix satisfying Assumption 4, and an initialization \mathbf{x}_0 , such that the associated parameters f_* , L, σ , \mathbf{W} , \mathbf{x}_0 , are independent of B. Then, $\forall T \geq 1, \forall \alpha > 0$, it holds that $\frac{1}{nT} \sum_{t=0}^{T-1} \sum_{i=1}^n \mathbb{E}[\|\nabla f(\mathbf{x}_i^t)\|] \geq B$.

We next break the limitations of vanilla gradient normalization in Claim 1 by incorporating gradient tracking and momentum variance reduction. This enables the successful use of normalization to suppress heavy-tailed noise while maintaining optimal convergence despite the added nonlinearity.

4 MAIN RESULTS

We present the main convergence results of GT-NSGDm and discuss their implications. The detailed analyses are deferred to the Appendix B. We first consider the case where the tail index p is known.

Theorem 1. Let Assumptions 1, 2, 3, 4 hold. Denote $f(\bar{x}^0) - f_* = \Delta_0, [\nabla f_1(\bar{x}^0), \dots, \nabla f_n(\bar{x}^0)]^\top = \nabla F(\mathbf{1}_n \otimes \bar{x}^0)$. Take

$$\alpha = \min\left(1, \sqrt{\frac{\Delta_0(1-\beta)(1-\lambda)}{4LT}}, \sqrt{\frac{\Delta_0(1-\lambda)}{3.5LT}}, \sqrt{\frac{(1-\lambda)^2 \Delta_0}{2n^{\frac{1}{2}}LT}}\right),\tag{5}$$

and $1-\beta=1/T^{\frac{p}{3p-2}}$. Assume $\beta\geq 1/10$, then the sequence generated by GT-NSGDm satisfies that

$$\begin{split} \frac{1}{nT} \sum_{t=0}^{T-1} \sum_{i=1}^{n} \mathbb{E} \big[\| \nabla f(\boldsymbol{x}_{i}^{t}) \| \big] &= O \Big(\frac{\sigma}{n^{1 - \frac{1}{p}} T^{\frac{p-1}{3p-2}}} + \frac{1}{T^{\frac{p-1}{3p-2}}} \sqrt{\frac{L\Delta_{0}}{1 - \lambda}} + \frac{\| \nabla f(\bar{\boldsymbol{x}}^{0}) \|}{T^{\frac{2p-2}{3p-2}}} + \sqrt{\frac{3.5 L\Delta_{0}}{(1 - \lambda)T}} \\ &+ \sqrt{\frac{n^{\frac{1}{2}} L\Delta_{0}}{(1 - \lambda)^{2} T}} + \frac{\sigma n^{\frac{1}{2}}}{(1 - \lambda)^{\frac{1}{p}} T^{\frac{p}{3p-2}}} + \frac{\| \nabla F(\mathbf{1}_{n} \otimes \bar{\boldsymbol{x}}^{0}) \|}{(1 - \lambda) n^{\frac{1}{2}} T^{\frac{p}{3p-2}}} + \frac{\sigma}{1 - \lambda} \frac{n^{\frac{1}{2}}}{T^{\frac{2p-1}{3p-2}}} + \frac{\Delta_{0}}{T} \Big). \end{split}$$

Remark 4 (Order-optimal rate). Theorem 1 establishes a non-asymptotic upper bound on the mean ℓ_2 norm stationary gap of GT-NSGDm over any finite time horizon T. The $O(\cdot)$ here only absorbs universal constants and preserves all problem parameters. It achieves the *optimal* $O(1/T^{\frac{p-1}{3p-2}})$ convergence rate in terms of T as it matches the lower bound proved in Zhang et al. (2020). This optimal guarantee is achieved in decentralized settings for the first time.

Remark 5 (Speedup in n). We discuss the asymptotic speedup in number of nodes n. For sufficiently large T (or sufficiently small target optimality gap), the upper bound in Theorem 1 is dominated by the leading terms $(1/T^{\frac{p-1}{3p-2}})(\sigma/n^{1-1/p}+\sqrt{L\Delta_0/(1-\lambda)})$. In the high-noise regime $\sigma\gg n^{1-1/p}\sqrt{L\Delta_0/(1-\lambda)}$, the upper bound has a speedup factor $n^{1-1/p}$. In practice, the noise scale (measured by σ) in training attention models or in other high-dimensional problems can be very large, and the speedup in n contributes as a noise reduction.

When the tail index p is unknown in advance, we establish the following convergence rate.

Theorem 2. Let Assumptions 1, 2, 3, 4 hold and take α as in (5). Take $1 - \beta = 1/\sqrt{T}$ and assume $\beta \ge 1/10$. Then GT-NSGDm guarantees that

$$\begin{split} \frac{1}{nT} \sum_{t=0}^{T-1} \sum_{i=1}^{n} \mathbb{E} \big[\| \nabla f(\boldsymbol{x}_{i}^{t}) \| \big] &\leq O \Big(\frac{\sigma}{n^{1 - \frac{1}{p}} T^{\frac{p-1}{2p}}} + \frac{1}{T^{\frac{1}{4}}} \sqrt{\frac{L\Delta_{0}}{1 - \lambda}} + \frac{\| \nabla f(\bar{\boldsymbol{x}}^{0}) \|}{\sqrt{T}} + \sqrt{\frac{3.5 L\Delta_{0}}{(1 - \lambda)T}} \\ &+ \frac{\sigma n^{\frac{1}{2}}}{(1 - \lambda)^{\frac{1}{p}} \sqrt{T}} + \frac{\| \nabla F(\mathbf{1}_{n} \otimes \bar{\boldsymbol{x}}^{0}) \|}{(1 - \lambda) n^{\frac{1}{2}} \sqrt{T}} + \sqrt{\frac{n^{\frac{1}{2}} L\Delta_{0}}{(1 - \lambda)^{2} T}} + \frac{\sigma n^{\frac{1}{2}}}{(1 - \lambda) T^{\frac{2p-1}{2p}}} + \frac{\Delta_{0}}{T} \Big). \end{split}$$

Theorem 2 establishes an upper bound of $O(1/T^{\frac{p-1}{2p}})$ when the tail index p is unknown, matching the best-known rate in the *centralized* setting where algorithm parameters do not rely on p (Liu & Zhou, 2025). While the convergence rate in Yu et al. (2023) is also independent of the knowledge of p, it is only for strongly convex functions and its exact rate exponent remains unspecified.

Remark 6 (Speedup in n and topology independent rate). Consider $p \in (1,2)$, i.e., the heavy-tailed case with unbounded variance this paper focuses on. When T is sufficiently large (as required to achieve sufficiently small target optimality gap), the upper bound in Theorem 2 is dominated by $\frac{\sigma}{n^{1-1/p}} \cdot \frac{1}{T^{(p-1)/2p}}$. Significantly, this upper bound is *independent* of network topology (λ) and exhibits a speedup factor $n^{1-1/p}$ in all regimes.

5 EXPERIMENTS

We assess the performance of GT-NSGDm through numerical experiments. We first conduct studies on synthetic datasets that mimic language modeling under controlled heavy-tailed noise injection, following Lee et al. (2025). We also present experiments on decentralized training of a decoder-only Transformer (GPT) model with 3M parameters on the Multi30k dataset.

Baselines. We compare GT-NSGDm with four decentralized baselines: DSGD(Nedic & Ozdaglar, 2009), GT-DSGD (Xin et al., 2020b), DSGD-Clip (Sun & Chen, 2024), and SClip-EF-Network (Yu et al., 2023). DSGD and GT-DSGD handle regular stochastic noise with bounded variance. DSGD-Clip converges for strongly convex functions under bounded domains or gradients (Sun & Chen, 2024). SClip-EF-Network achieves convergence under symmetric noise with bounded $\mathbb{E}\left[\|\boldsymbol{\xi}_i^t\|^p \mid \mathcal{F}_{t-1}\right]$ for p=1. All methods are initialized identically and tuned via grid search. Detailed baseline descriptions appear in Table 2 (Appendix C.1).

Graph topology. We consider three graph topologies: undirected ring, directed exponential, and complete graphs (see Lian et al. (2017); Nedić et al. (2018); Assran et al. (2019)). Weight matrices

use Metropolis weights (Xiao et al., 2005). For synthetic experiments, we set the number of nodes to n=20, we obtain $\lambda=0.904,\,0.714$, and 0 for the ring, exponential, and complete graphs, respectively. For Transformer training with n=8, we have corresponding $\lambda=0.804,\,0.6$, and 0.

5.1 ROBUST LINEAR REGRESSION ON SYNTHETIC TOKENIZED DATA

We use this synthetic experiment to test our convergence rates under controlled heavy-tailed noise. We consider nonconvex regularized linear regression on synthetic data mimicking language tokens. In language modeling, token frequencies exhibit heavy-tailed distributions: few tokens appear frequently, while most are rare but contextually important. We construct the following synthetic dataset \boldsymbol{X} of 1k samples of dimension d=20. The first two features simulate frequent tokens, sampled from Bernoulli distributions Bern(0.9) and Bern(0.5), respectively. The remaining 8 features represent rare tokens, each sampled from Bern(0.1). The optimal weight \boldsymbol{w}_* is Gaussian-sampled, with labels $\boldsymbol{y}=\boldsymbol{X}\boldsymbol{w}_*$. The synthetic dataset $(\boldsymbol{X},\boldsymbol{y})$ is evenly distributed over n=20 nodes, where each node i holds a sub-dataset $(\boldsymbol{X}_i,\boldsymbol{y}_i)$, estimate \boldsymbol{w}_i , and a linear regression model with nonconvex robust Tukey's biweight loss function (Beaton & Tukey, 1974) to estimate \boldsymbol{w}_* . We inject three different zero-mean noises, Gaussian noise $(\mathcal{N}(0,3\boldsymbol{I}_d))$, Student's t noise (degrees of freedom 1.5, scale 1.0), and Lévy α -stable noise (stability parameter 1.5, skewness parameter 0.5, scale 1.0, nonsymmetric, multiplied by 0.1) into exact gradient, using corrupted stochastic gradients for updates. See Appendix C.2 for additional details.

In Figure 2, we evaluate GT-NSGDm against baselines on ring graphs under various gradient noise. DSGD and GT-DSGD converge under Gaussian noise but become unstable under heavy-tailed noise. DSGD-Clip remains stable but fails to reach optimum. Both GT-NSGDm and SClip-EF-Network exhibit robust convergence and near-optimal performance across all scenarios, consistent with their theoretical guarantees under heavy-tailed noise.

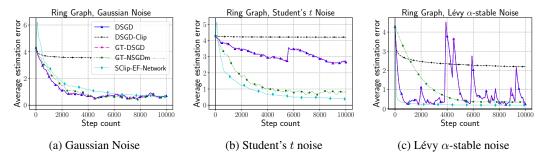


Figure 2: Comparison of performance on a ring graph under various types of injected stochastic gradient noise, measured by the average estimation error $(1/n)\sum_{i=1}^{n}\|\boldsymbol{w}_{i}^{t}-\boldsymbol{w}_{*}\|$ over step count t.

In Figure 3, we test GT-NSGDm's dependence on connectivity (λ), noise level (σ), and the number of nodes (n), varying each while fixing others. In Figure 3(a), we inject Lévy α -stable noise and test the performance of GT-NSGDm on ring, directed exponential, undirected exponential, and complete graphs with $\lambda=0.904,0.714,0.6,0$, respectively. GT-NSGDm achieves comparable final errors under weak connectivity (i.e., large λ) versus complete graphs, showing favorable dependence on network connectivity under heavy-tailed noise. In Figure 3(b), we evaluate GT-NSGDm's performance under different noise levels on a directed exponential graph. Under Gaussian noise with scale 1 (unit variance), GT-NSGDm reaches the best optimality; the final error increases as σ grows, as observed under both Gaussian and Lévy α -stable noise. In Figure 3(c), we inject Lévy α -stable noise on complete graphs ($\lambda=0$ for all n) with varying number of nodes. As n increases from 2 to 40, convergence speed improves with final errors [0.4,0.35,0.29,0.20,0.21], demonstrating speedup over certain n ranges, supporting theoretical discussions in Remarks 5 and 6.

5.2 DECENTRALIZED TRAINING OF TRANSFORMERS

We evaluate GT-NSGDm's empirical performance on language modeling using a 3M-parameter GPT model (Radford et al., 2018) for auto-regressive modeling on Multi30k (29k sentences, 4.4M to-kens). We assess performance using validation log-perplexity. On 8-node graphs with three topolo-

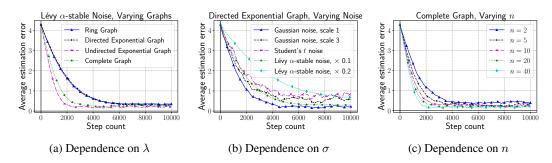


Figure 3: Empirical studies on GT-NSGDm's dependence on problem parameters λ, σ, n .

gies, we distribute training data evenly and initialize identical GPT models per node. We introduce three additional baselines: DSGD-GClip (DSGD with constant step size and ℓ_2 gradient clipping level), DSGD-CClip (DSGD with constant step size and component-wise gradient clipping level), QG-DSGDm (Lin et al., 2021) and GT-Adam (Carnevale et al., 2022) (all without theoretical guarantees under heavy-tailed noise; see Table 2 in Appendix C.1). We run all methods for 12 epochs with batch size 64. See Appendix C.3 for model and hyperparameter details.

Table 1 presents average validation loss and standard deviation over five independent runs for each algorithm across three topologies. Results show that GT-NSGDm nearly matches the best baseline DSGD-GClip (which lacks theoretical guarantees under heavy-tailed gradient noise) while significantly outperforming the other two theoretically-guaranteed baselines across all topologies. We note that this decentralized training experiment is simulated to demonstrate algorithm effectiveness and has practical limitations.

Table 1: Topologies ring, directed exponential (Exp.), and complete (Comp.) graphs. Algorithms are grouped by theoretical (Theo.) guarantees under heavy-tailed noise: with (w/) or without (w/o).

Algorithms	Theo.	Ring	Exp.	Comp.
DSGD	w/o	$5.633_{\pm 0.008}$	$5.632_{\pm 0.007}$	$5.635_{\pm 0.007}$
DSGD-GClip	w/o	$0.253_{\pm 0.007}$	$0.249_{\pm 0.010}$	$0.267_{\pm 0.010}$
DSGD-CClip	w/o	$2.725_{\pm 3.179}$	$5.058_{\pm 2.388}$	$8.225_{\pm 1.695}$
GT-DSGD	w/o	$5.362_{\pm 0.002}$	$5.632_{\pm 0.002}$	$5.631_{\pm 0.002}$
GT-Adam	w/o	$0.520_{\pm 0.038}$	$0.587_{\pm 0.096}$	$0.524_{\pm 0.045}$
QG-DSGDm	w/o	$0.394_{\pm 0.007}$	$0.388_{\pm0.013}$	$0.353_{\pm 0.011}$
SClip-EF-Network	w/	$5.653_{\pm 0.012}$	$5.632_{\pm 0.003}$	$5.636_{\pm 0.004}$
DSGD-Clip	w/	$5.633_{\pm 0.004}$	$5.659_{\pm 0.013}$	$5.661_{\pm 0.006}$
GT-NSGDm	w/	$0.258_{\pm 0.007}$	$0.261_{\pm 0.007}$	$0.282_{\pm 0.009}$

6 CONCLUSION AND FUTURE WORK

In this paper, we have proposed GT-NSGDm for solving decentralized nonconvex smooth optimization to address heavy-tailed noise. The key idea is to leverage normalization, together with momentum variance reduction, to combat heavy-tailed noise, and use gradient tracking to handle cross-node heterogeneity and the nonlinearity brought by normalization. Theoretical analyses establish that GT-NSGDm attains optimal convergence rate when the tail index p is known, and a rate that matches the best centralized one when p is unknown. Extensive experiments on nonconvex linear regression and decentralized Transformer training show that GT-NSGDm is robust and efficient under heavy-tailed noise across various topologies, and achieves a speedup in n. Future directions include extending the current analysis to other nonlinearities, such as sign and clipping (Zhang et al., 2020), and generalizing GT-NSGDm to handle objective functions under relaxed smoothness conditions (Liu & Zhou, 2025).

REPRODUCIBLE STATEMENT

We provide detailed proofs for Claim 1 in Appendix A and for our main theoretical results, Theorems 1 and 2, in Appendix B. In Appendix 5, we provide detailed hardware configurations, algorithm descriptions, and hyperparameter settings for our numerical experiments.

REFERENCES

- Kwangjun Ahn, Xiang Cheng, Minhak Song, Chulhee Yun, Ali Jadbabaie, and Suvrit Sra. Linear attention is (maybe) all you need (to understand transformer optimization). In *The Twelfth International Conference on Learning Representations*, 2024.
- Aleksandar Armacki, Shuhua Yu, Dragana Bajovic, Dusan Jakovetic, and Soummya Kar. Large deviations and improved mean-squared error rates of nonlinear sgd: Heavy-tailed noise and power of symmetry. *arXiv preprint arXiv:2410.15637*, 2024.
- Aleksandar Armacki, Shuhua Yu, Pranay Sharma, Gauri Joshi, Dragana Bajovic, Dusan Jakovetic, and Soummya Kar. High-probability convergence bounds for online nonlinear stochastic gradient descent under heavy-tailed noise. In *The 28th International Conference on Artificial Intelligence and Statistics*, 2025.
- Mahmoud Assran, Nicolas Loizou, Nicolas Ballas, and Mike Rabbat. Stochastic gradient push for distributed deep learning. In *International Conference on Machine Learning*, pp. 344–353. PMLR, 2019.
- Melih Barsbey, Milad Sefidgaran, Murat A Erdogdu, Gael Richard, and Umut Simsekli. Heavy tails in sgd and compressibility of overparametrized neural networks. *Advances in neural information processing systems*, 34:29364–29378, 2021.
- Barak Battash, Lior Wolf, and Ofir Lindenbaum. Revisiting the noise model of stochastic gradient descent. In *International Conference on Artificial Intelligence and Statistics*, pp. 4780–4788. PMLR, 2024.
- Albert E Beaton and John W Tukey. The fitting of power series, meaning polynomials, illustrated on band-spectroscopic data. *Technometrics*, 16(2):147–185, 1974.
- Theodora S Brisimi, Ruidi Chen, Theofanie Mela, Alex Olshevsky, Ioannis Ch Paschalidis, and Wei Shi. Federated learning of predictive models from federated electronic health records. *International journal of medical informatics*, 112:59–67, 2018.
- Guido Carnevale, Francesco Farina, Ivano Notarnicola, and Giuseppe Notarstefano. Gtadam: Gradient tracking with adaptive momentum for distributed online optimization. *IEEE Transactions on Control of Network Systems*, 10(3):1436–1448, 2022.
- Zachary Charles, Zachary Garrett, Zhouyuan Huo, Sergei Shmulyian, and Virginia Smith. On large-cohort training for federated learning. *Advances in neural information processing systems*, 34: 20461–20475, 2021.
- Savelii Chezhegov, Yaroslav Klyukin, Andrei Semenov, Aleksandr Beznosikov, Alexander Gasnikov, Samuel Horváth, Martin Takáč, and Eduard Gorbunov. Gradient clipping improves adagrad when the noise is heavy-tailed. *arXiv preprint arXiv:2406.04443*, 2024.
- Enea Monzio Compagnoni, Rustem Islamov, Frank Norbert Proske, and Aurelien Lucchi. Unbiased and sign compression in distributed learning: Comparing noise resilience via SDEs. In *The 28th International Conference on Artificial Intelligence and Statistics*, 2025a.
- Enea Monzio Compagnoni, Tianlin Liu, Rustem Islamov, Frank Norbert Proske, Antonio Orvieto, and Aurelien Lucchi. Adaptive methods through the lens of SDEs: Theoretical insights on the role of noise. In *The Thirteenth International Conference on Learning Representations*, 2025b.
- Paolo Di Lorenzo and Gesualdo Scutari. Next: In-network nonconvex optimization. *IEEE Transactions on Signal and Information Processing over Networks*, 2(2):120–136, 2016.

- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(7), 2011.
- Desmond Elliott, Stella Frank, Khalil Sima'an, and Lucia Specia. Multi30k: Multilingual englishgerman image descriptions. *arXiv preprint arXiv:1605.00459*, 2016.
 - Bahman Gharesifard and Jorge Cortés. Distributed strategies for generating weight-balanced and doubly stochastic digraphs. *European Journal of Control*, 18(6):539–557, 2012.
 - Eduard Gorbunov, Marina Danilova, and Alexander Gasnikov. Stochastic optimization with heavy-tailed noise via accelerated gradient clipping. *Advances in Neural Information Processing Systems*, 33:15042–15053, 2020.
 - Eduard Gorbunov, Abdurakhmon Sadiev, Marina Danilova, Samuel Horváth, Gauthier Gidel, Pavel Dvurechensky, Alexander Gasnikov, and Peter Richtárik. High-probability convergence for composite and distributed stochastic minimization and variational inequalities with heavy-tailed noise, 2024. URL https://openreview.net/forum?id=qOFLnOpMoe.
 - Mert Gurbuzbalaban, Umut Simsekli, and Lingjiong Zhu. The heavy-tail phenomenon in sgd. In *International Conference on Machine Learning*, pp. 3964–3975. PMLR, 2021.
 - Florian Hübler, Ilyas Fatkhullin, and Niao He. From gradient clipping to normalization for heavy tailed sgd. *arXiv preprint arXiv:2410.13849*, 2024.
 - Dusan Jakovetić, Dragana Bajović, Anit Kumar Sahu, Soummya Kar, Nemanja Milosević, and Dusan Stamenković. Nonlinear gradient mappings and stochastic optimization: A general framework with applications to heavy-tail noise. *SIAM Journal on Optimization*, 33(2):394–423, 2023.
 - Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. Advances and open problems in federated learning. *Foundations and Trends® in Machine Learning*, 14(1–2):1–210, 2021.
 - Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
 - Anastasia Koloskova, Nicolas Loizou, Sadra Boreiri, Martin Jaggi, and Sebastian Stich. A unified theory of decentralized sgd with changing topology and local updates. In *International Conference on Machine Learning*, pp. 5381–5393. PMLR, 2020.
 - Frederik Kunstner, Alan Milligan, Robin Yadav, Mark Schmidt, and Alberto Bietti. Heavy-tailed class imbalance and why adam outperforms gradient descent on language models. *Advances in Neural Information Processing Systems*, 37:30106–30148, 2024.
 - Su Hyeong Lee, Manzil Zaheer, and Tian Li. Efficient distributed optimization under heavy-tailed noise. *arXiv preprint arXiv:2502.04164*, 2025.
 - Boyue Li and Yuejie Chi. Convergence and privacy of decentralized nonconvex optimization with gradient clipping and communication compression. *IEEE Journal of Selected Topics in Signal Processing*, 2025.
 - Tian Li, Anit Kumar Sahu, Ameet Talwalkar, and Virginia Smith. Federated learning: Challenges, methods, and future directions. *IEEE signal processing magazine*, 37(3):50–60, 2020.
 - Xiangru Lian, Ce Zhang, Huan Zhang, Cho-Jui Hsieh, Wei Zhang, and Ji Liu. Can decentralized algorithms outperform centralized algorithms? a case study for decentralized parallel stochastic gradient descent. *Advances in neural information processing systems*, 30, 2017.
 - Tao Lin, Sai Praneeth Karimireddy, Sebastian Stich, and Martin Jaggi. Quasi-global momentum: Accelerating decentralized deep learning on heterogeneous data. In *International Conference on Machine Learning*, pp. 6654–6665. PMLR, 2021.
 - Zijian Liu and Zhengyuan Zhou. Nonconvex stochastic optimization under heavy-tailed noises: Optimal convergence without gradient clipping. In *The Thirteenth International Conference on Learning Representations*, 2025.

- Zijian Liu, Jiawei Zhang, and Zhengyuan Zhou. Breaking the lower bound with (little) structure: Acceleration in non-convex stochastic optimization with heavy-tailed noise. In *The Thirty Sixth Annual Conference on Learning Theory*, pp. 2266–2290. PMLR, 2023.
 - Jayakrishnan Nair, Adam Wierman, and Bert Zwart. *The fundamentals of heavy tails: Properties, emergence, and estimation*, volume 53. Cambridge University Press, 2022.
 - Angelia Nedic and Asuman Ozdaglar. Distributed subgradient methods for multi-agent optimization. *IEEE Transactions on Automatic Control*, 54(1):48–61, 2009.
 - Angelia Nedić, Alex Olshevsky, and Michael G Rabbat. Network topology and communication-computation tradeoffs in decentralized optimization. *Proceedings of the IEEE*, 106(5):953–976, 2018.
 - Ta Duy Nguyen, Thien H Nguyen, Alina Ene, and Huy Nguyen. Improved convergence in high probability of clipped gradient methods with heavy tailed noise. *Advances in Neural Information Processing Systems*, 36:24191–24222, 2023.
 - Alex Olshevsky. Linear time average consensus on fixed graphs and implications for decentralized optimization and multi-agent control. *arXiv* preprint arXiv:1411.4186, 2014.
 - Stefano Peluchetti, Stefano Favaro, and Sandra Fortini. Stable behaviour of infinitely wide deep neural networks. In *International Conference on Artificial Intelligence and Statistics*, pp. 1137–1146. PMLR, 2020.
 - Shi Pu and Angelia Nedić. Distributed stochastic gradient tracking methods. *Mathematical Programming*, 187(1):409–457, 2021.
 - Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. 2018.
 - Abdurakhmon Sadiev, Marina Danilova, Eduard Gorbunov, Samuel Horváth, Gauthier Gidel, Pavel Dvurechensky, Alexander Gasnikov, and Peter Richtárik. High-probability bounds for stochastic optimization and variational inequalities: the case of unbounded variance. In *International Conference on Machine Learning*, pp. 29563–29648. PMLR, 2023.
 - Egor Shulgin, Sarit Khirirat, and Peter Richtárik. Smoothed normalization for efficient distributed private optimization. *arXiv preprint arXiv:2502.13482*, 2025.
 - Umut Simsekli, Levent Sagun, and Mert Gurbuzbalaban. A tail-index analysis of stochastic gradient noise in deep neural networks. In *International Conference on Machine Learning*, pp. 5827–5837. PMLR, 2019.
 - Chao Sun and Bo Chen. Distributed stochastic strongly convex optimization under heavy-tailed noises. In 2024 IEEE International Conference on Cybernetics and Intelligent Systems (CIS) and IEEE International Conference on Robotics, Automation and Mechatronics (RAM), pp. 150–155. IEEE, 2024.
 - Tao Sun, Xinwang Liu, and Kun Yuan. Gradient normalization provably benefits nonconvex sgd under heavy-tailed noise. *arXiv preprint arXiv:2410.16561*, 2024.
 - S Sundhar Ram, Angelia Nedić, and Venugopal V Veeravalli. Distributed stochastic subgradient projection algorithms for convex optimization. *Journal of optimization theory and applications*, 147:516–545, 2010.
 - John Tsitsiklis, Dimitri Bertsekas, and Michael Athans. Distributed asynchronous deterministic and stochastic gradient optimization algorithms. *IEEE transactions on automatic control*, 31(9): 803–812, 1986.
 - Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

- Jianyu Wang and Gauri Joshi. Cooperative sgd: A unified framework for the design and analysis of local-update sgd algorithms. *Journal of Machine Learning Research*, 22(213):1–50, 2021.
- Lin Xiao, Stephen Boyd, and Sanjay Lall. A scheme for robust distributed sensor fusion based on average consensus. In *IPSN 2005. Fourth International Symposium on Information Processing in Sensor Networks*, 2005., pp. 63–70. IEEE, 2005.
- Ran Xin, Usman A Khan, and Soummya Kar. Variance-reduced decentralized stochastic optimization with accelerated convergence. *IEEE Transactions on Signal Processing*, 68:6255–6271, 2020a.
- Ran Xin, Shi Pu, Angelia Nedić, and Usman A Khan. A general framework for decentralized optimization with first-order methods. *Proceedings of the IEEE*, 108(11):1869–1889, 2020b.
- Haibo Yang, Peiwen Qiu, and Jia Liu. Taming fat-tailed ("heavier-tailed" with potentially infinite variance) noise in federated learning. *Advances in Neural Information Processing Systems*, 35: 17017–17029, 2022.
- Shuhua Yu and Soummya Kar. Secure distributed optimization under gradient attacks. *IEEE Transactions on Signal Processing*, 2023.
- Shuhua Yu, Dusan Jakovetic, and Soummya Kar. Smoothed gradient clipping and error feedback for decentralized optimization under symmetric heavy-tailed noise. *arXiv preprint arXiv:2310.16920*, 2023.
- Kun Yuan, Bicheng Ying, Jiageng Liu, and Ali H Sayed. Variance-reduced stochastic learning by networked agents under random reshuffling. *IEEE Transactions on Signal Processing*, 67(2): 351–366, 2018.
- Jingzhao Zhang, Sai Praneeth Karimireddy, Andreas Veit, Seungyeon Kim, Sashank Reddi, Sanjiv Kumar, and Suvrit Sra. Why are adaptive methods good for attention models? Advances in Neural Information Processing Systems, 33:15383–15393, 2020.

APPENDIX

A Proof of Claim 1

Proof. Consider n scalar functions that for each $i \in \mathcal{V}$, $f_i(x) = (1/2)(x - a_i)^2$ for some a_i , and complete graph with $\mathbf{W} = (1/n)\mathbf{1}_n\mathbf{1}_n^{\top}$. Let $a_i = a, \forall i = 1, \ldots, n/2$, and $a_i = b, \forall i = n/2+1, \ldots, n$, and b-a > 2B+1. Let $x_i^0 = a+0.5, \forall i \in \mathcal{V}$. Then, $\forall i \in \mathcal{V}$, vanilla normalization reduces to

$$\begin{split} x_i^1 &= \frac{1}{n} \sum_{r=1}^n \left(x_r^0 - \alpha \mathrm{sign}(x_r^0 - a_r) \right) \\ &= x_r^0 - \frac{\alpha}{n} \sum_{r=1}^n \mathrm{sign}(x_r^0 - a_r) \\ &= x_r^0 - \frac{\alpha}{n} \sum_{r=1}^{n/2} \mathrm{sign}(0.5) - \frac{\alpha}{n} \sum_{r=n/2+1}^n \mathrm{sign}(0.5 - (b-a)) \\ &= x_r^0. \end{split}$$

- Therefore, $x_r^t = a + 0.5, \forall r \in \mathcal{V}, \forall t \geq 0$. Since the optimal solution to the original problem is $\frac{a+b}{2}$, the optimality gap is $\frac{b-a}{2} 0.5 \geq B$.
- Remark 7. Note that the proof above can be further extended to the case where the gradient oracle admits almost surely bounded gradient noise. We can use the noise bound to adapt the choices of a, b, ε such that all signs still get canceled. Similar examples have been used to show divergence results in Shulgin et al. (2025).

B PROOFS OF THEOREMS

B.1 PRELIMINARIES

We define some stacked long vectors,

$$F(oldsymbol{x}^t) := [f_1(oldsymbol{x}_1^t), \dots, f_n(oldsymbol{x}_n^t)]^ op, \
abla F(oldsymbol{x}^t) := [
abla_1(oldsymbol{x}_1^t)^ op, \dots,
abla_f(oldsymbol{x}_n^t)^ op]^ op, \
abla G(oldsymbol{x}^t, oldsymbol{\xi}^t) := [oldsymbol{g}_1(oldsymbol{x}_1^t, oldsymbol{\xi}_1^t)^ op, \dots, oldsymbol{g}_n(oldsymbol{x}_n^t, oldsymbol{\xi}_n^t)^ op]^ op, \
abla G(oldsymbol{x}^t) := [oldsymbol{y}_1^t]^ op, \dots, oldsymbol{y}_n^t]^ op, \
abla G(oldsymbol{x}_n^t) := [oldsymbol{x}_1^t]^ op, \dots, oldsymbol{y}_n^t]^ op, \
abla G(oldsymbol{x}_n^t)^ op]^ op.$$

Then, Algorithm 1 can be rewritten in the compact long-vector form:

$$\boldsymbol{v}^{t} = \beta \boldsymbol{v}^{t-1} + (1 - \beta) \boldsymbol{g}(\boldsymbol{x}^{t}, \boldsymbol{\xi}^{t}); \tag{6}$$

$$\mathbf{y}^t = (\mathbf{W} \otimes \mathbf{I}_d)(\mathbf{y}^{t-1} + \mathbf{v}^t - \mathbf{v}^{t-1}), \tag{7}$$

$$\boldsymbol{x}^{t+1} = (\boldsymbol{W} \otimes \boldsymbol{I}_d)(\boldsymbol{x}^t - \alpha \mathcal{N}(\boldsymbol{y}^t)). \tag{8}$$

We define the following averages over network:

$$\bar{\boldsymbol{v}}^t = \frac{1}{n} \sum_{i=1}^n \boldsymbol{v}_i^t, \quad \bar{\boldsymbol{y}}^t = \frac{1}{n} \sum_{i=1}^n \boldsymbol{y}_i^t, \quad \tilde{\boldsymbol{y}}^t = \frac{1}{n} \sum_{i=1}^n \frac{\boldsymbol{y}_i^t}{\|\boldsymbol{y}_i^t\|}, \quad \bar{\boldsymbol{x}}^t = \frac{1}{n} \sum_{i=1}^n \boldsymbol{x}_i^t, \quad \overline{\nabla} F(\boldsymbol{x}^t) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\boldsymbol{x}_i^t).$$
(9)

From the doubly stochasticity of W, the global average updates as

$$\bar{\boldsymbol{x}}^{t+1} = \bar{\boldsymbol{x}}^t - \frac{\alpha}{n} \sum_{r=1}^n \frac{\boldsymbol{y}_r^t}{\|\boldsymbol{y}_r^t\|} = \bar{\boldsymbol{x}}^t - \alpha \tilde{\boldsymbol{y}}^t.$$
(10)

B.2 Intermediate Lemmas

We first present some standard useful relations to be used in our analyses.

Lemma 1. The following relations hold:

1.
$$\bar{\boldsymbol{y}}^t = \bar{\boldsymbol{v}}^t$$
;

2.
$$W - \mathbf{1}_n \mathbf{1}_n^{\top} / n = (W - \mathbf{1}_n \mathbf{1}_n^{\top} / n) (I_n - \mathbf{1}_n \mathbf{1}_n^{\top} / n) = (I_n - \mathbf{1}_n \mathbf{1}_n^{\top} / n) (W - \mathbf{1}_n \mathbf{1}_n^{\top} / n);$$

3.
$$\mathbf{W}^k - \mathbf{1}_n \mathbf{1}_n^{\top} / n = (\mathbf{W} - \mathbf{1}_n \mathbf{1}_n^{\top} / n)^k, \forall k \in \mathbb{N}_+;$$

4.
$$(1/\sqrt{n}) \sum_{i=1}^{n} \|\boldsymbol{a}_i\| \le \|\boldsymbol{a}\| \le \sum_{i=1}^{n} \|\boldsymbol{a}_i\|, \forall \boldsymbol{a} = [\boldsymbol{a}_1^\top, \dots, \boldsymbol{a}_n^\top]^\top \in \mathbb{R}^{nd},$$

5.
$$\sum_{i=1}^m a_i^p \le \left(\sum_{i=1}^m a_i\right)^p \le m^{p-1} \sum_{i=1}^m a_i^p, \forall m \in \mathbb{N}_+, \forall a_i \in \mathbb{R}_+.$$

We then present a standard decent lemma for L-smooth functions.

Lemma 2 (Decent lemma for L-smooth functions). Let Assumption 2 hold. For any $x, y \in \mathbb{R}^d$, there holds

$$f(\boldsymbol{y}) \leq f(\boldsymbol{x}) + \nabla f(\boldsymbol{x})^{\top} (\boldsymbol{y} - \boldsymbol{x}) + \frac{L}{2} \|\boldsymbol{x} - \boldsymbol{y}\|^{2}.$$

We next present the main descent lemma on the network average.

Lemma 3 (Decent lemma for network average). Let Assumption 2 hold. Let $\epsilon^t = \bar{y}^t - \nabla f(\bar{x}^t)$. We have

$$\sum_{t=0}^{T-1} \alpha \|\nabla f(\bar{\boldsymbol{x}}^t)\| \le f(\bar{\boldsymbol{x}}^0) - f_* + \sum_{t=0}^{T-1} 2\alpha \|\boldsymbol{\epsilon}^t\| + \sum_{t=0}^{T-1} \frac{\alpha}{n} \sum_{i=1}^n \|\bar{\boldsymbol{y}}^t - \boldsymbol{y}_i^t\| + \sum_{t=0}^{T-1} \frac{L}{2} \alpha^2.$$

Proof. Since $\|\bar{x}^{t+1} - \bar{x}^t\| = \alpha \|\tilde{y}^t\| = \alpha$, applying Lemma 2 on \bar{x}^{t+1}, \bar{x}^t gives that

$$f(\bar{\boldsymbol{x}}^{t+1}) \leq f(\bar{\boldsymbol{x}}^t) + \nabla f(\bar{\boldsymbol{x}}^t)^{\top} (\bar{\boldsymbol{x}}^{t+1} - \bar{\boldsymbol{x}}^t) + \frac{L}{2} \|\bar{\boldsymbol{x}}^{t+1} - \bar{\boldsymbol{x}}^t\|^2$$

$$\stackrel{(i)}{\leq} f(\bar{\boldsymbol{x}}^t) - \alpha (\bar{\boldsymbol{y}}^t - \boldsymbol{\epsilon}^t)^{\top} \tilde{\boldsymbol{y}}^t + \frac{L}{2} \alpha^2$$

$$\stackrel{(ii)}{\leq} f(\bar{\boldsymbol{x}}^t) - \alpha (\bar{\boldsymbol{y}}^t)^{\top} \tilde{\boldsymbol{y}}^t + \alpha \|\boldsymbol{\epsilon}^t\| + \frac{L}{2} \alpha^2, \tag{11}$$

where we used the definitions (9)(10) in (i), and used Cauchy-Schwartz inequality followed by $\|\tilde{y}^t\| \le 1$ in (ii). Next,

$$-(\bar{\boldsymbol{y}}^{t})^{\top}\tilde{\boldsymbol{y}}^{t} = -(\bar{\boldsymbol{y}}^{t})^{\top} \left[\frac{\bar{\boldsymbol{y}}^{t}}{\|\bar{\boldsymbol{y}}^{t}\|} + \frac{1}{n} \sum_{i=1}^{n} \boldsymbol{y}_{i}^{t} \left(\frac{1}{\|\boldsymbol{y}_{i}^{t}\|} - \frac{1}{\|\bar{\boldsymbol{y}}^{t}\|} \right) \right]$$

$$\leq -\|\bar{\boldsymbol{y}}^{t}\| + \|\frac{1}{n} \sum_{i=1}^{n} \boldsymbol{y}_{i}^{t} \left(\frac{\|\bar{\boldsymbol{y}}^{t}\|}{\|\boldsymbol{y}_{i}^{t}\|} - 1 \right) \|$$

$$\stackrel{(i)}{\leq} -\|\nabla f(\bar{\boldsymbol{x}}^{t})\| + \|\boldsymbol{\epsilon}^{t}\| + \frac{1}{n} \sum_{i=1}^{n} \|\|\bar{\boldsymbol{y}}^{t}\| - \|\boldsymbol{y}_{i}^{t}\| \|$$

$$\stackrel{(ii)}{\leq} -\|\nabla f(\bar{\boldsymbol{x}}^{t})\| + \|\boldsymbol{\epsilon}^{t}\| + \frac{1}{n} \sum_{i=1}^{n} \|\bar{\boldsymbol{y}}^{t} - \boldsymbol{y}_{i}^{t}\|, \tag{12}$$

where we used $\|\bar{\boldsymbol{y}}^t\| = \|\nabla f(\bar{\boldsymbol{x}}^t) + \boldsymbol{\epsilon}^t\| \ge \|\nabla f(\bar{\boldsymbol{x}}^t)\| - \|\boldsymbol{\epsilon}^t\|$, and Cauchy-Schwartz inequality in (i), and $\|\boldsymbol{a}\| - \|\boldsymbol{b}\| \le \|\boldsymbol{a} - \boldsymbol{b}\|$ for any $\boldsymbol{a}, \boldsymbol{b} \in \mathbb{R}^d$ in (ii). Plugging in (12) into (11), and summing over $t = 0, \dots, T-1$, we have

$$f(\bar{\boldsymbol{x}}^T) \leq f(\bar{\boldsymbol{x}}^0) - \sum_{t=0}^{T-1} \alpha \|\nabla f(\bar{\boldsymbol{x}}^t)\| + \sum_{t=0}^{T-1} 2\alpha \|\boldsymbol{\epsilon}^t\| + \sum_{t=0}^{T-1} \frac{\alpha}{n} \sum_{i=1}^n \|\bar{\boldsymbol{y}}^t - \boldsymbol{y}_i^t\| + \sum_{t=0}^{T-1} \frac{L}{2} \alpha^2.$$

Using $f(\bar{x}^T) \geq f_*$ and rearranging terms above give the desired result.

With Lemma 3, it remains to bound the gradient estimation error $\|\epsilon^t\|$ and the consensus error $y_i^t - \bar{y}^t$. Let us decompose the gradient estimation error as follows:

$$\boldsymbol{\epsilon}^{t} = \bar{\boldsymbol{y}}^{t} - \nabla f(\bar{\boldsymbol{x}}^{t}) = \bar{\boldsymbol{v}}^{t} - \nabla f(\bar{\boldsymbol{x}}^{t}) = \underbrace{\bar{\boldsymbol{v}}^{t} - \overline{\nabla} F(\boldsymbol{x}^{t})}_{:=\boldsymbol{\epsilon}_{1}^{t} \in \mathbb{R}^{d}} + \underbrace{\overline{\nabla} F(\boldsymbol{x}^{t}) - \nabla f(\bar{\boldsymbol{x}}^{t})}_{:=\boldsymbol{\epsilon}_{2}^{t} \in \mathbb{R}^{d}}.$$
 (13)

It is clear that ϵ_1^t is the gradient estimation error, and ϵ_2^t , exploiting the smoothness property in 2, can be bounded by the consensus error $\boldsymbol{x}_i^t - \bar{\boldsymbol{x}}^t$. Since the consensus error is also used in bounding ϵ_1^t , we need to first bound the consensus errors $\boldsymbol{x}_i^t - \bar{\boldsymbol{x}}^t$ and $\boldsymbol{y}_i^t - \bar{\boldsymbol{y}}^t$.

Lemma 4 (Consensus errors of $\{x_i^t\}$). We have for all t = 0, ..., T,

$$\frac{1}{n} \sum_{i=1}^{n} \|\boldsymbol{x}_{i}^{t} - \bar{\boldsymbol{x}}^{t}\| \le \frac{\alpha \lambda}{1 - \lambda}.$$
(14)

Proof. Using the relation 4 in Lemma 1 we have

$$\frac{1}{n} \sum_{i=1}^{n} \| \boldsymbol{x}_{i}^{t} - \bar{\boldsymbol{x}}^{t} \| \leq \frac{1}{\sqrt{n}} \| \boldsymbol{x}^{t} - \mathbf{1}_{n} \otimes \bar{\boldsymbol{x}}^{t} \|.$$
 (15)

From the compact form update in (8), we have

$$m{x}^t = ig(m{W} \otimes m{I}_dig)m{x}^0 - lpha \sum_{k=0}^{t-1} (m{W} \otimes m{I}_d)^{t-k} \mathcal{N}(m{y}^k).$$

It follows that

$$\|\boldsymbol{x}^{t} - \boldsymbol{1}_{n} \otimes \bar{\boldsymbol{x}}^{t}\|$$

$$= \|(\boldsymbol{I}_{nd} - \frac{1}{n} \boldsymbol{1}_{n} \boldsymbol{1}_{n}^{\top} \otimes \boldsymbol{I}_{d}) \boldsymbol{x}^{t}\|$$

$$\stackrel{(i)}{=} \alpha \| \sum_{k=0}^{t-1} (\boldsymbol{I}_{nd} - \frac{1}{n} \boldsymbol{1}_{n} \boldsymbol{1}_{n}^{\top} \otimes \boldsymbol{I}_{d}) (\boldsymbol{W} \otimes \boldsymbol{I}_{d})^{t-k} \mathcal{N}(\boldsymbol{y}^{k}) \|$$

$$\leq \alpha \| \sum_{k=0}^{t-1} \|\boldsymbol{W}^{t-k} - \frac{1}{n} \boldsymbol{1}_{n} \boldsymbol{1}_{n}^{\top} \|_{2} \|\mathcal{N}(\boldsymbol{y}^{k})\|$$

$$\stackrel{(ii)}{\leq} \alpha \| \sum_{k=0}^{t-1} \|\boldsymbol{W} - \frac{1}{n} \boldsymbol{1}_{n} \boldsymbol{1}_{n}^{\top} \|_{2}^{t-k} \|\mathcal{N}(\boldsymbol{y}^{k})\|$$

$$\leq \alpha \sqrt{n} \sum_{k=0}^{t-1} \lambda^{t-k}$$

$$\stackrel{(iii)}{\leq} \frac{\alpha \sqrt{n} \lambda}{1 - \lambda}.$$

$$(16)$$

where we used the double stochasticity of W and $x_i^0 = \bar{x}^0, \forall i \in [n]$ in (i), the relation 3 in Lemma 1 in (ii), and Assumption 4 in (iii). Substituting (17) into (15) gives the desired bound in (14).

Before proceeding to bound consensus errors for $\{y_i^t\}$, we present the following bound on vector-valued martingale difference sequence from Liu & Zhou (2025).

Lemma 5. Given a sequence of random vectors $\mathbf{d}_t \in \mathbb{R}^d$, $\forall t$ such that $\mathbb{E}[\mathbf{d}_t \mid \mathcal{F}_{t-1}] = \mathbf{0}$ where $\mathcal{F}_t = \sigma(\mathbf{d}_1, \dots, \mathbf{d}_t)$ is the natural filtration, then for any $p \in [1, 2]$, there is

$$\mathbb{E}\Big[\|\sum_{t=1}^T \boldsymbol{d}_t\|\Big] \leq 2\sqrt{2}\mathbb{E}\Big[\Big(\sum_{t=1}^T \|\boldsymbol{d}_t\|^p\Big)^{\frac{1}{p}}\Big], \forall T \geq 0.$$

Lemma 6 (Consensus errors for $\{y_i^t\}$). We have for all t = 0, ..., T,

$$\begin{split} &\frac{1}{n}\mathbb{E}\big[\sum_{i=1}^{n}\|\boldsymbol{y}_{i}^{t}-\bar{\boldsymbol{y}}^{t}\|\big]\\ &\leq 2\sqrt{2}n^{\frac{1}{2}}\big(\frac{1}{\beta}-1\big)\Big(\sum_{k=0}^{t}\lambda^{(t-k+1)p}\Big)^{\frac{1}{p}}\sigma+\frac{1}{\sqrt{n}}\big(\frac{1}{\beta}-1\big)\sum_{k=0}^{t}\lambda^{t-k+1}\mathbb{E}\big[\|\nabla F(\boldsymbol{x}^{k})-\boldsymbol{v}^{k}\|\big]. \end{split}$$

Proof. Similar to (15), we have

$$\frac{1}{n}\sum_{i=1}^{n}\|\boldsymbol{y}_{i}^{t}-\bar{\boldsymbol{y}}^{t}\| \leq \frac{1}{\sqrt{n}}\|\boldsymbol{y}^{t}-\boldsymbol{1}_{n}\otimes\bar{\boldsymbol{y}}^{t}\|.$$
(18)

Following from (7),

$$\mathbf{y}^{t} - \mathbf{1}_{n} \otimes \bar{\mathbf{y}}^{t} \qquad (19)$$

$$\stackrel{(7)}{=} \left(\mathbf{I}_{nd} - \frac{1}{n} \mathbf{1}_{n} \mathbf{1}_{n}^{\top} \otimes \mathbf{I}_{d}\right) \left(\mathbf{W} \otimes \mathbf{I}_{d}\right) \left(\mathbf{y}^{t-1} + \mathbf{v}^{t} - \mathbf{v}^{t-1}\right)$$

$$= \left(\mathbf{W} \otimes \mathbf{I}_{d} - \frac{1}{n} \mathbf{1}_{n} \mathbf{1}_{n}^{\top} \otimes \mathbf{I}_{d}\right) \mathbf{y}^{t-1} + \left(\mathbf{W} \otimes \mathbf{I}_{d} - \frac{1}{n} \mathbf{1}_{n} \mathbf{1}_{n}^{\top} \otimes \mathbf{I}_{d}\right) \left(\mathbf{v}^{t} - \mathbf{v}^{t-1}\right)$$

$$\stackrel{(i)}{=} \left(\mathbf{W} \otimes \mathbf{I}_{d} - \frac{1}{n} \mathbf{1}_{n} \mathbf{1}_{n}^{\top} \otimes \mathbf{I}_{d}\right) \left(\mathbf{I}_{nd} - \frac{1}{n} \mathbf{1}_{n} \mathbf{1}_{n}^{\top} \otimes \mathbf{I}_{d}\right) \mathbf{y}^{t-1} + \left(\mathbf{W} \otimes \mathbf{I}_{d} - \frac{1}{n} \mathbf{1}_{n} \mathbf{1}_{n}^{\top} \otimes \mathbf{I}_{d}\right) \left(\mathbf{v}^{t} - \mathbf{v}^{t-1}\right)$$

$$\stackrel{(ii)}{=} \sum_{t=0}^{t} \left(\mathbf{W} \otimes \mathbf{I}_{d} - \frac{1}{n} \mathbf{1}_{n} \mathbf{1}_{n}^{\top} \otimes \mathbf{I}_{d}\right)^{t-k+1} \left(\mathbf{v}^{t} - \mathbf{v}^{t-1}\right),$$
(20)

where we used relation 3 in Lemma 1 in (i) and used $\mathbf{y}_i^0 = \mathbf{0}_d, \forall i \in [n]$ in (ii). From the update in (6), we have

$$v^{t} - v^{t-1} = (\beta - 1)v^{t-1} + (1 - \beta)g(x^{t}, \xi^{t}) = (1 - \beta)(v^{t} - v^{t-1}) + (1 - \beta)(g(x^{t}, \xi^{t}) - v^{t}).$$

Then, there holds,

$$v^{t} - v^{t-1} = (\frac{1}{\beta} - 1)(g(x^{t}, \xi^{t}) - v^{t}) = (\frac{1}{\beta} - 1)(g(x^{t}, \xi^{t}) - \nabla F(x^{t}) + \nabla F(x^{t}) - v^{t}).$$
(21)

Putting the relation above into (20) and applying (20) recursively, from $y_i^0 = \bar{y}^0$, we have

$$\|\boldsymbol{y}^{t} - \boldsymbol{1}_{n} \otimes \bar{\boldsymbol{y}}^{t}\| \leq \left(\frac{1}{\beta} - 1\right) \| \sum_{k=0}^{t} \left(\boldsymbol{W} \otimes \boldsymbol{I}_{d} - \frac{1}{n} \boldsymbol{1}_{n} \boldsymbol{1}_{n}^{\top} \otimes \boldsymbol{I}_{d}\right)^{t-k+1} \left(\boldsymbol{g}(\boldsymbol{x}^{k}, \boldsymbol{\xi}^{k}) - \nabla F(\boldsymbol{x}^{k})\right) \|$$

$$+ \left(\frac{1}{\beta} - 1\right) \| \sum_{k=0}^{t} \left(\boldsymbol{W} \otimes \boldsymbol{I}_{d} - \frac{1}{n} \boldsymbol{1}_{n} \boldsymbol{1}_{n}^{\top} \otimes \boldsymbol{I}_{d}\right)^{t-k+1} \left(\nabla F(\boldsymbol{x}^{k}) - \boldsymbol{v}^{k}\right) \|$$
(22)

We note that the first half of the right hand side above can be addressed by Lemma 5:

$$\mathbb{E}\left[\|\sum_{k=0}^{t} \left(\boldsymbol{W} \otimes \boldsymbol{I}_{d} - \frac{1}{n} \boldsymbol{1}_{n} \boldsymbol{1}_{n}^{\top} \otimes \boldsymbol{I}_{d}\right)^{t-k+1} \left(\boldsymbol{g}(\boldsymbol{x}^{k}, \boldsymbol{\xi}^{k}) - \nabla F(\boldsymbol{x}^{k})\right)\|\right]$$

$$\leq 2\sqrt{2}\mathbb{E}\left[\left(\sum_{k=0}^{t} \lambda^{(t-k+1)p} \|\boldsymbol{g}(\boldsymbol{x}^{k}, \boldsymbol{\xi}^{k}) - \nabla F(\boldsymbol{x}^{k})\|^{p}\right)^{\frac{1}{p}}\right].$$
(23)

We observe that

$$2\sqrt{2}\mathbb{E}\left[\left(\sum_{k=0}^{t} \lambda^{(t-k+1)p} \| \boldsymbol{g}(\boldsymbol{x}^{k}, \boldsymbol{\xi}^{k}) - \nabla F(\boldsymbol{x}^{k}) \|^{p}\right)^{\frac{1}{p}} | \mathcal{F}_{t-1}\right]$$

$$\leq 2\sqrt{2}\mathbb{E}\left[\left(\sum_{k=0}^{t} \lambda^{(t-k+1)p} \left(\sum_{i=1}^{n} \| \boldsymbol{g}_{i}(\boldsymbol{x}_{i}^{k}, \boldsymbol{\xi}_{i}^{k}) - \nabla f_{i}(\boldsymbol{x}_{i}^{k}) \|\right)^{p}\right)^{\frac{1}{p}} | \mathcal{F}_{t-1}\right]$$

$$\stackrel{(i)}{\leq} 2\sqrt{2}\mathbb{E}\left[\left(\sum_{k=0}^{t} \sum_{i=1}^{n} \lambda^{(t-k+1)p} n^{p-1} \| \boldsymbol{g}_{i}(\boldsymbol{x}_{i}^{k}, \boldsymbol{\xi}_{i}^{k}) - \nabla f_{i}(\boldsymbol{x}_{i}^{k}) \|^{p}\right)^{\frac{1}{p}} | \mathcal{F}_{t-1}\right]$$

$$\stackrel{(ii)}{\leq} 2\sqrt{2}\left(\mathbb{E}\left[\sum_{i=1}^{n} \lambda^{p} n^{p-1} \| \boldsymbol{g}_{i}(\boldsymbol{x}_{i}^{t}, \boldsymbol{\xi}_{i}^{t}) - \nabla f_{i}(\boldsymbol{x}_{i}^{t}) \|^{p} | \mathcal{F}_{t-1}\right]$$

$$+ \sum_{k=0}^{t-1} \sum_{i=1}^{n} \lambda^{(t-k+1)p} n^{p-1} \| \boldsymbol{g}_{i}(\boldsymbol{x}_{i}^{k}, \boldsymbol{\xi}_{i}^{k}) - \nabla f_{i}(\boldsymbol{x}^{k}) \|^{p}\right)^{\frac{1}{p}}$$

$$\stackrel{(iii)}{\leq} 2\sqrt{2}\left(\lambda^{p} n^{p} \sigma^{p} + \sum_{i=1}^{t-1} \sum_{i=1}^{n} \lambda^{(t-k+1)p} n^{p-1} \| \boldsymbol{g}_{i}(\boldsymbol{x}_{i}^{k}, \boldsymbol{\xi}_{i}^{k}) - \nabla f_{i}(\boldsymbol{x}_{i}^{k}) \|^{p}\right)^{\frac{1}{p}},$$

where we used relation 5 from Lemma 1 in (i), Jensen's inequality in (ii), and Assumption 3 in (iii). From (23), taking expectations on both sides of (24), and applying the above arguments recursively from \mathcal{F}_{t-2} to \mathcal{F}_0 , we have

$$\mathbb{E}\Big[\|\sum_{k=0}^t \big(\boldsymbol{W} \otimes \boldsymbol{I}_d - \frac{1}{n}\boldsymbol{1}_n\boldsymbol{1}_n^\top \otimes \boldsymbol{I}_d\big)^{t-k+1} \big(\boldsymbol{g}(\boldsymbol{x}^k,\boldsymbol{\xi}^k) - \nabla F(\boldsymbol{x}^k)\big)\|\Big] \leq 2\sqrt{2}\Big(\sum_{k=0}^t \lambda^{(t-k+1)p}\Big)^{\frac{1}{p}} n\sigma.$$

Therefore, using the relation above, and (18), (22), we reach the desired relation.

We then bound average gradient estimation errors $\epsilon_1^t = \bar{\boldsymbol{v}}^t - \overline{\nabla} F(\boldsymbol{x}^t)$.

Lemma 7 (Average gradient estimation errors). For all t = 0, ..., T, we have

$$\mathbb{E}\big[\|\bar{\boldsymbol{v}}^t - \overline{\nabla}F(\boldsymbol{x}^t)\|\big] \leq \beta^{t+1}\|\nabla f(\bar{\boldsymbol{x}}^0)\| + \frac{2\sqrt{2}}{n^{1-\frac{1}{p}}}\Big(\sum_{k=0}^t \beta^{(t-k)p}(1-\beta)^p\Big)^{\frac{1}{p}}\sigma + \sum_{k=0}^t \beta^{t-k+1}\Big(\frac{2\alpha\lambda}{1-\lambda} + \alpha\Big)L.$$

Proof. Following from the step 4 in Algorithm 1, $\forall i \in [n]$,

$$\boldsymbol{v}_i^t - \nabla f_i(\boldsymbol{x}_i^t) = \beta(\boldsymbol{v}_i^{t-1} - \nabla f_i(\boldsymbol{x}_i^{t-1})) + (1 - \beta)(\boldsymbol{g}_i(\boldsymbol{x}_i^t, \boldsymbol{\xi}_i^t) - \nabla f_i(\boldsymbol{x}_i^t)) + \beta(\nabla f_i(\boldsymbol{x}_i^{t-1}) - \nabla f_i(\boldsymbol{x}_i^t)). \tag{25}$$

Averaging the above relation over i = 1, ..., n leads to that:

$$\boldsymbol{\epsilon}_1^t = \bar{\boldsymbol{v}}^t - \overline{\nabla} F(\boldsymbol{x}^t)$$

$$=\beta(\bar{\boldsymbol{v}}^{t-1}-\overline{\nabla}F(\boldsymbol{x}^{t-1}))+(1-\beta)\cdot\underbrace{\frac{1}{n}\sum_{i=1}^{n}(\boldsymbol{g}_{i}(\boldsymbol{x}_{i}^{t},\boldsymbol{\xi}_{i}^{t})-\nabla f_{i}(\boldsymbol{x}_{i}^{t}))}_{:=\boldsymbol{s}^{t}\in\mathbb{R}^{d}}+\beta\cdot\underbrace{\frac{1}{n}\sum_{i=1}^{n}(\nabla f_{i}(\boldsymbol{x}_{i}^{t-1})-\nabla f_{i}(\boldsymbol{x}_{i}^{t}))}_{:=\boldsymbol{z}^{t}\in\mathbb{R}^{d}}$$

$$= \beta^{t+1} \epsilon_1^{-1} + \sum_{k=0}^t \beta^{t-k} (1-\beta) s^k + \sum_{k=0}^t \beta^{t-k+1} z^k.$$

Taking Euclidean norms on both sides gives that

$$\|\boldsymbol{\epsilon}_{1}^{t}\| \leq \beta^{t+1}\|\boldsymbol{\epsilon}_{1}^{-1}\| + \|\sum_{k=0}^{t} \beta^{t-k}(1-\beta)\boldsymbol{s}^{k}\| + \|\sum_{k=0}^{t} \beta^{t-k+1}\boldsymbol{z}^{k}\|.$$
 (26)

We now bound the terms on the right hand side of (26) one by one. First,

$$\|\boldsymbol{\epsilon}_1^{-1}\| = \|\bar{\boldsymbol{v}}^{-1} - \frac{1}{n} \sum_{i=1}^n \nabla f_i(\bar{\boldsymbol{x}}^0)\| = \|\nabla f(\bar{\boldsymbol{x}}^0)\|.$$
 (27)

Second, notice that $\{\beta^{t-k}(1-\beta)(\boldsymbol{g}_i(\boldsymbol{x}_i^k,\boldsymbol{\xi}_i^k)-\nabla f_i(\boldsymbol{x}_i^k))\}$ is a martingale difference sequence that falls into the pursuit of Lemma 5, and thus we obtain

$$\mathbb{E}\left[\left\|\sum_{k=0}^{t} \beta^{t-k} (1-\beta) \boldsymbol{s}^{k}\right\|\right] \\
= \frac{1}{n} \mathbb{E}\left[\left\|\sum_{k=0}^{t} \sum_{i=1}^{n} \beta^{t-k} (1-\beta) (\boldsymbol{g}_{i}(\boldsymbol{x}_{i}^{k}, \boldsymbol{\xi}_{i}^{k}) - \nabla f_{i}(\boldsymbol{x}_{i}^{k}))\right\|\right] \\
\leq \frac{2\sqrt{2}}{n} \mathbb{E}\left[\left(\sum_{k=0}^{t} \sum_{i=1}^{n} \|\beta^{t-k} (1-\beta) (\boldsymbol{g}_{i}(\boldsymbol{x}_{i}^{k}, \boldsymbol{\xi}_{i}^{k}) - \nabla f_{i}(\boldsymbol{x}_{i}^{k}))\right\|^{p}\right]^{\frac{1}{p}}\right].$$
(28)

Note that

$$\frac{2\sqrt{2}}{n} \mathbb{E}\left[\left(\sum_{k=0}^{t} \sum_{i=1}^{n} \|\beta^{t-k} (1-\beta) (\boldsymbol{g}_{i}(\boldsymbol{x}_{i}^{k}, \boldsymbol{\xi}_{i}^{k}) - \nabla f_{i}(\boldsymbol{x}_{i}^{k}))\|^{p}\right)^{\frac{1}{p}} | \mathcal{F}_{t-1}\right] \\
\stackrel{(i)}{\leq} \frac{2\sqrt{2}}{n} \left(\mathbb{E}\left[\sum_{k=0}^{t} \sum_{i=1}^{n} \|\beta^{t-k} (1-\beta) (\boldsymbol{g}_{i}(\boldsymbol{x}_{i}^{k}, \boldsymbol{\xi}_{i}^{k}) - \nabla f_{i}(\boldsymbol{x}_{i}^{k}))\|^{p} | \mathcal{F}_{t-1}\right]\right)^{\frac{1}{p}} \\
\leq \frac{2\sqrt{2}}{n} \left(\mathbb{E}\left[\sum_{i=1}^{n} (1-\beta)^{p} \|(\boldsymbol{g}_{i}(\boldsymbol{x}_{i}^{t}, \boldsymbol{\xi}_{i}^{t}) - \nabla f_{i}(\boldsymbol{x}_{i}^{t}))\|^{p} | \mathcal{F}_{t-1}\right] \\
+ \sum_{k=0}^{t-1} \sum_{i=1}^{n} \|\beta^{t-k} (1-\beta) (\boldsymbol{g}_{i}(\boldsymbol{x}_{i}^{k}, \boldsymbol{\xi}_{i}^{k}) - \nabla f_{i}(\boldsymbol{x}_{i}^{k}))\|^{p}\right)^{\frac{1}{p}} \\
\stackrel{(ii)}{\leq} \frac{2\sqrt{2}}{n} \left(n(1-\beta)^{p} \sigma^{p} + \sum_{l=0}^{t-1} \sum_{i=1}^{n} \|\beta^{t-k} (1-\beta) (\boldsymbol{g}_{i}(\boldsymbol{x}_{i}^{k}, \boldsymbol{\xi}_{i}^{k}) - \nabla f_{i}(\boldsymbol{x}_{i}^{k}))\|^{p}\right)^{\frac{1}{p}}, \tag{29}$$

where we used Jensen's inequality in (i) and Assumption 3 in (ii). From (28), taking expectations on (29), and recursively applying the preceding arguments from \mathcal{F}_{t-2} to \mathcal{F}_0 , we have

$$\mathbb{E}\Big[\|\sum_{k=0}^{t} \beta^{t-k} (1-\beta) s^{k}\|\Big] \le \frac{2\sqrt{2}}{n^{1-\frac{1}{p}}} \Big(\sum_{k=0}^{t} \beta^{(t-k)p} (1-\beta)^{p}\Big)^{\frac{1}{p}} \sigma.$$
(30)

Third, $\|\sum^t \beta^{t-k+1} \boldsymbol{z}^k\|$ $\leq \sum_{i=1}^{t} \beta^{t-k+1} \| \frac{1}{n} \sum_{i=1}^{n} (\nabla f_i(\boldsymbol{x}_i^{k-1}) - \nabla f_i(\boldsymbol{x}_i^k)) \|$ $\leq \sum_{i=1}^{t} \beta^{t-k+1} \left(\frac{1}{n} \sum_{i=1}^{n} \|\nabla f_i(\bar{x}_i^{k-1}) - \nabla f_i(\bar{x}_i^{k-1})\| + \frac{1}{n} \sum_{i=1}^{n} \|\nabla f_i(\bar{x}_i^{k-1}) - \nabla f_i(\bar{x}_i^{k})\| \right)$ $+\frac{1}{n}\sum_{i=1}^{n}\left\|
abla f_{i}(ar{oldsymbol{x}}^{k}abla f_{i}(oldsymbol{x}_{i}^{k})
ight\|
ight)$

$$\overset{(i)}{\leq} \sum_{k=0}^{t} \beta^{t-k+1} \Big(L \cdot \frac{1}{n} \sum_{i=1}^{n} \| \boldsymbol{x}_{i}^{k-1} - \bar{\boldsymbol{x}}^{k-1} \| + L \cdot \frac{1}{n} \sum_{i=1}^{n} \| \bar{\boldsymbol{x}}^{k-1} - \bar{\boldsymbol{x}}^{k} \| + L \cdot \frac{1}{n} \sum_{i=1}^{n} \| \bar{\boldsymbol{x}}^{k} - \boldsymbol{x}_{i}^{k} \| \Big)$$

$$\stackrel{(ii)}{\leq} \sum_{k=0}^{t} \beta^{t-k+1} \left(\frac{2\alpha\lambda}{1-\lambda} + \alpha \right) L.$$

 (31)

where in (i) we used Assumption 2 and in (ii) we used (14) in Lemma 4. Putting relations (27)(30)(31) together leads to the final bound for this lemma.

We next bound the stacked gradient estimation errors.

Lemma 8 (Stacked gradient estimation errors). For all t = 0, ..., T, we have

$$\mathbb{E}\big[\|\boldsymbol{v}^t - \nabla F(\boldsymbol{x}^t)\|\big]$$

$$\leq \beta^{t+1} \|\nabla F(\mathbf{1}_n \otimes \bar{\mathbf{x}}^0)\| + 2\sqrt{2} \Big(\sum_{k=0}^t \beta^{(t-k)p} (1-\beta)^p \Big)^{\frac{1}{p}} n\sigma + n \sum_{k=0}^t \beta^{t-k+1} \Big(\frac{2\alpha\lambda}{1-\lambda} + \alpha \Big) L.$$

Proof. Define $\tilde{\epsilon}_1^t := v^t - \nabla F(x^t) \in \mathbb{R}^{nd}$. Similar to (25), we have

$$\boldsymbol{v}^{t} - \nabla F(\boldsymbol{x}^{t}) = \beta(\boldsymbol{v}^{t-1} - \nabla F(\boldsymbol{x}^{t-1})) + (1 - \beta) \underbrace{(\boldsymbol{g}(\boldsymbol{x}^{t}, \boldsymbol{\xi}^{t}) - \nabla F(\boldsymbol{x}^{t}))}_{:=\tilde{\boldsymbol{s}}^{t} \in \mathbb{R}^{nd}} + \beta \underbrace{(\nabla F(\boldsymbol{x}^{t-1}) - \nabla F(\boldsymbol{x}^{t}))}_{:=\tilde{\boldsymbol{z}}^{t} \in \mathbb{R}^{nd}}$$

$$= \beta^{t+1} \tilde{\epsilon}_1^{-1} + \sum_{k=0}^t \beta^{t-k} (1-\beta) \tilde{s}^k + \sum_{k=0}^t \beta^{t-k+1} \tilde{z}^k.$$

Taking Euclidean norms on both sides gives that

$$\|\tilde{\boldsymbol{\epsilon}}_1^t\| \le \beta^{t+1} \|\tilde{\boldsymbol{\epsilon}}_1^{-1}\| + \|\sum_{k=0}^t \beta^{t-k} (1-\beta) \tilde{\boldsymbol{s}}^k\| + \|\sum_{k=0}^t \beta^{t-k+1} \tilde{\boldsymbol{z}}^k\|.$$

Similar to the analysis in Lemma 7, we bound the right hand side above term by term. First,

$$\|\tilde{\boldsymbol{\epsilon}}_1^{-1}\| = \|\nabla F(\mathbf{1}_n \otimes \bar{\boldsymbol{x}}^0)\|.$$

Second, notice also that $\{\beta^{t-k}(1-\beta)\tilde{s}^k\}$ is a martingale difference sequence and can be dealt with using Lemma 5. We have

$$\mathbb{E}\left[\|\sum_{k=0}^{t} \beta^{t-k} (1-\beta)\tilde{\boldsymbol{s}}^{k}\|\right]$$

$$= \mathbb{E}\left[\|\sum_{k=0}^{t} \beta^{t-k} (1-\beta)(\boldsymbol{g}(\boldsymbol{x}^{t}, \boldsymbol{\xi}^{t,b}) - \nabla F(\boldsymbol{x}^{t}))\|\right]$$

$$\leq 2\sqrt{2}\mathbb{E}\left[\left(\sum_{k=0}^{t} \beta^{(t-k)p} (1-\beta)^{p} \|\boldsymbol{g}(\boldsymbol{x}^{t}, \boldsymbol{\xi}^{t}) - \nabla F(\boldsymbol{x}^{t})\|^{p}\right)^{\frac{1}{p}}\right].$$
(32)

In addition,

$$2\sqrt{2}\mathbb{E}\Big[\Big(\sum_{k=0}^{t}\beta^{(t-k)p}(1-\beta)^{p}\|\boldsymbol{g}(\boldsymbol{x}^{t},\boldsymbol{\xi}^{t})-\nabla F(\boldsymbol{x}^{t})\|^{p}\Big)^{\frac{1}{p}}\mid\mathcal{F}_{t-1}\Big]$$

$$\stackrel{(i)}{\leq} 2\sqrt{2}\Big(\mathbb{E}\Big[\sum_{k=0}^{t}\beta^{(t-k)p}(1-\beta)^{p}\|\boldsymbol{g}(\boldsymbol{x}^{t},\boldsymbol{\xi}^{t})-\nabla F(\boldsymbol{x}^{t})\|^{p}\Big]\mid\mathcal{F}_{t-1}\Big)^{\frac{1}{p}}$$

$$\stackrel{(ii)}{\leq} 2\sqrt{2}\Big(\mathbb{E}\Big[\sum_{k=0}^{t}\beta^{(t-k)p}(1-\beta)^{p}\Big(\sum_{i=1}^{n}\|\boldsymbol{g}_{i}(\boldsymbol{x}_{i}^{t},\boldsymbol{\xi}_{i}^{t})-\nabla f_{i}(\boldsymbol{x}_{i}^{t})\|^{p}\Big]\mid\mathcal{F}_{t-1}\Big)^{\frac{1}{p}}$$

$$\stackrel{(iii)}{\leq} 2\sqrt{2}\Big(\mathbb{E}\Big[\sum_{k=0}^{t}\sum_{i=1}^{n}\beta^{(t-k)p}(1-\beta)^{p}n^{p-1}\|\boldsymbol{g}_{i}(\boldsymbol{x}_{i}^{t},\boldsymbol{\xi}_{i}^{t})-\nabla f_{i}(\boldsymbol{x}_{i}^{t})\|^{p}\Big]\mid\mathcal{F}_{t-1}\Big)^{\frac{1}{p}}$$

$$= 2\sqrt{2}\Big(\mathbb{E}\Big[\sum_{i=1}^{n}(1-\beta)^{p}n^{p-1}\|\nabla \boldsymbol{g}_{i}(\boldsymbol{x}_{i}^{t},\boldsymbol{\xi}_{i}^{t})-\nabla f_{i}(\boldsymbol{x}_{i}^{t})\|^{p}\mid\mathcal{F}_{t-1}\Big]$$

$$+\sum_{k=0}^{t-1}\sum_{i=1}^{n}\beta^{(t-k)p}(1-\beta)^{p}n^{p-1}\|\boldsymbol{g}_{i}(\boldsymbol{x}_{i}^{t},\boldsymbol{\xi}_{i}^{t})-\nabla f_{i}(\boldsymbol{x}_{i}^{t})\|^{p}\Big)^{\frac{1}{p}}$$

$$\leq 2\sqrt{2}\Big((1-\beta)^{p}n^{p}\sigma^{p}+\sum_{k=0}^{t-1}\sum_{i=1}^{n}\beta^{(t-k)p}(1-\beta)^{p}n^{p-1}\|\boldsymbol{g}_{i}(\boldsymbol{x}_{i}^{t},\boldsymbol{\xi}_{i}^{t})-\nabla f_{i}(\boldsymbol{x}_{i}^{t})\|^{p}\Big)^{\frac{1}{p}},$$

where in (i) we used Jensen's inequality, and in (ii), (iii) we used relations 4, and 5 in Lemma 1, respectively. Based on (32), taking expectations on both sides of (33), and applying the above arguments from \mathcal{F}_{t-2} to \mathcal{F}_0 , we obtain

$$\mathbb{E} \big[\| \sum_{k=0}^{t} \beta^{t-k} (1-\beta) \tilde{s}^k \| \big] \le 2\sqrt{2} \Big(\sum_{k=0}^{t} \beta^{(t-k)p} (1-\beta)^p \Big)^{\frac{1}{p}} n\sigma.$$

Third,

$$\begin{split} &\| \sum_{k=0}^{t} \beta^{t-k+1} \tilde{\boldsymbol{z}}^{k} \| \\ &\leq \sum_{k=0}^{t} \beta^{t-k+1} \| \nabla F(\boldsymbol{x}^{k-1}) - \nabla F(\boldsymbol{x}^{k}) \| \\ &\leq \sum_{k=0}^{t} \sum_{i=1}^{n} \beta^{t-k+1} \| \nabla f_{i}(\boldsymbol{x}_{i}^{k-1}) - \nabla f_{i}(\boldsymbol{x}_{i}^{k}) \| \\ &\leq \sum_{k=0}^{t} \sum_{i=1}^{n} \beta^{t-k+1} \| \nabla f_{i}(\boldsymbol{x}_{i}^{k-1}) - \nabla f_{i}(\bar{\boldsymbol{x}}^{k}) \| \\ &\leq \sum_{k=0}^{t} \sum_{i=1}^{n} \beta^{t-k+1} \Big(\| \nabla f_{i}(\boldsymbol{x}_{i}^{k-1}) - \nabla f_{i}(\bar{\boldsymbol{x}}^{k-1}) \| + \| \nabla f_{i}(\bar{\boldsymbol{x}}^{k-1}) - \nabla f_{i}(\bar{\boldsymbol{x}}^{k}) \| + \| \nabla f_{i}(\boldsymbol{x}_{i}^{k}) - f_{i}(\bar{\boldsymbol{x}}^{k}) \| \Big) \\ &\stackrel{(i)}{\leq} n \sum_{k=0}^{t} \beta^{t-k+1} \Big(\frac{2\alpha\lambda}{1-\lambda} + \alpha \Big) L. \end{split}$$

where (i) follows from similar arguments in (30).

Now we are ready to prove our main theorems.

Proof of Theorem 1. We observe that

$$\frac{1}{n} \sum_{t=0}^{T-1} \sum_{i=1}^{n} \mathbb{E}[\|\nabla f(\boldsymbol{x}_{i}^{t})\|] \leq \frac{1}{n} \sum_{t=0}^{T-1} \sum_{i=1}^{n} \mathbb{E}[\|\nabla f(\boldsymbol{x}_{i}^{t}) - \nabla f(\bar{\boldsymbol{x}}^{t})\| + \|\nabla f(\bar{\boldsymbol{x}}^{t})\|] \\
\leq T \cdot \frac{\alpha \lambda L}{1 - \lambda} + \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(\bar{\boldsymbol{x}}^{t})\|].$$
(34)

From Lemmas 3, (13), and Lemma 6,
$$\sum_{i=0}^{T-1} \alpha \|\nabla f(\vec{x}^i)\|$$

$$\leq f(\vec{x}^0) - f_* + \sum_{i=0}^{T-1} 2\alpha (\|\vec{e}_i^i\| + \|\nabla F(\vec{x}^i) - \nabla f(\vec{x}^i)\|) + \sum_{i=0}^{T-1} \frac{\alpha}{n} \sum_{i=1}^{n} \|\vec{y}^i - y_i^i\| + \sum_{i=0}^{T-1} \frac{L}{2}\alpha^2$$

$$\leq f(\vec{x}^0) - f_* + \sum_{i=0}^{T-1} 2\alpha (\|\vec{e}_i^i\| + \|\nabla F(\vec{x}^i) - \nabla f(\vec{x}^i)\|) + \sum_{i=0}^{T-1} \frac{\alpha}{n} \sum_{i=1}^{n} \|\vec{y}^i - y_i^i\| + \sum_{i=0}^{T-1} \frac{L}{2}\alpha^2$$

$$\leq f(\vec{x}^0) - f_* + \sum_{i=0}^{T-1} 2\alpha (\|\vec{e}_i^i\| + \|\nabla F(\vec{x}^i) - \nabla f(\vec{x}^i)\|) + \sum_{i=0}^{T-1} \frac{\alpha}{n} \sum_{i=0}^{n} \|\vec{y}^i - y_i^i\| + \sum_{i=0}^{T-1} \frac{L}{2}\alpha^2$$

$$+ \sum_{i=0}^{T-1} 2\alpha (\|\vec{g}^{i+1}\|\nabla f(\vec{x}^0)\| + \frac{2\sqrt{2}}{n^{1-\frac{1}{p}}} (\sum_{k=0}^{i} \beta^{(i-k)p}(1-\beta)^p)^{\frac{1}{p}} \sigma + \sum_{k=0}^{i} \beta^{i-k+1} \mathbb{E} \|\nabla F(\vec{x}^k) - \vec{v}^k\| \|]$$

$$+ \sum_{i=0}^{T-1} \alpha (\|\vec{g}^{i+1}\|\nabla f(\vec{x}^0)\| + \frac{\alpha}{1-\beta} + \frac{4\sqrt{2}\sigma}{n^{1-\frac{1}{p}}} \alpha (1-\beta)^{1-\frac{1}{p}} T + \frac{4L}{1-\lambda} \cdot \frac{\alpha^2T}{1-\beta} + \frac{2L}{1-\lambda} \cdot \alpha^2T$$

$$+ \sum_{i=0}^{2} \alpha^2 LT$$

$$+ \frac{2\sqrt{2}\sigma n_j^2}{(1-\lambda)^{\frac{1}{p}}} \cdot (\frac{1}{\beta} - 1)\alpha T + \frac{1}{2}L \cdot \alpha^2T$$

$$+ \frac{1}{\sqrt{n}} (\frac{1}{\beta} - 1)\alpha \sum_{i=0}^{T-1} \sum_{i=1}^{i} \lambda^{i-k+1} (\beta^{i+1}\|\nabla F(1_n \otimes \vec{x}^0)\| + 2\sqrt{2}n\sigma(1-\beta)^{1-\frac{1}{p}} + \frac{2nL}{1-\lambda} \cdot \frac{\alpha\beta}{1-\beta})$$
where in (i) we used $\beta \leq 1, \lambda < 1$ and Lemma 8. Denote $f(\vec{x}^0) - f_* = \Delta_0$. Dividing αT from the above relation on both sides, and putting it into (34), then rearranging terms leads to

$$\frac{1}{nT} \sum_{i=0}^{T-1} \sum_{i=1}^{n} \mathbb{E} \|\nabla f(\vec{x}^i)\| + 4\sqrt{2}\sigma \cdot \frac{(1-\beta)^{1-\frac{1}{p}}}{n^{1-\frac{1}{p}}} + \frac{4L}{1-\lambda} \cdot \frac{\alpha}{1-\beta} + \frac{3L}{1-\lambda} \cdot \frac{1}{n^2} - \frac{1}{1-\lambda} \cdot \frac{2L}{1-\lambda}$$

$$+ \frac{2\sqrt{2}\sigma}{(1-\lambda)^{\frac{1}{p}}} \cdot \frac{1}{\beta} (1-\beta) + \frac{1}{n^{1-\frac{1}{p}}} + \frac{4L}{1-\lambda} \cdot \frac{\alpha}{1-\beta} + \frac{2\sqrt{2}\sigma}{1-\lambda} \cdot \frac{n^{\frac{1}{2}}}{(1-\beta)^{1-\frac{1}{p}}} + \frac{2L}{(1-\lambda)^{2}} \cdot \frac{n^{\frac{1}{2}}}{(1-\lambda)^{2}} \cdot \frac{n^{\frac{1}{2}}}{(1-\lambda)^{2}} \cdot \frac{n^{\frac{1}{2}}}{(1-\lambda)^{2}} + \frac{2\sqrt{2}\sigma}{(1-\lambda)^{2}} \cdot \frac{n^{\frac{1}{2}}}{(1-\lambda)^{2}} \cdot \frac{n^{\frac{1}{2}}}{(1-\lambda)^{2}} \cdot \frac{n^{\frac{1}{2}}}{n^{\frac{1}{2}}} + \frac{n^{\frac{1}{2}}}{(1-\lambda)^{2}} \cdot \frac{n^{\frac{1}{2}}}{(1-\lambda)^{2}} \cdot \frac{n^{\frac{1}{2}}}{(1-\lambda)^{2}} \cdot \frac{n^{\frac{1}{2}}}{(1-\lambda)^{2}} \cdot \frac{n^{\frac{1}{2}}}{(1-\lambda)^{2}} \cdot \frac{n^{\frac{1}{2}}}{(1-\lambda)^{2}} \cdot \frac{n^{\frac{1}{2}}}{(1-\lambda)$$

(35)

where in (i) we take $\beta \geq 1/10$, in (ii) we used

$$\alpha = \min\left(1, \sqrt{\frac{\Delta_0(1-\beta)(1-\lambda)}{4LT}}, \sqrt{\frac{\Delta_0(1-\lambda)}{3.5LT}}, \sqrt{\frac{(1-\lambda)^2 \Delta_0}{2n^{\frac{1}{2}}LT}}\right),$$
(36)

and in (iii) we used
$$1 - \beta = \frac{1}{T^{\frac{p}{3p-2}}}$$
.

Proof of Theorem 2. Note that (35)(ii) still holds under the same choice of α in (36) and $\beta \geq 1/10$. Continuing with $1 - \beta = 1/\sqrt{T}$, we have

$$\begin{split} &\frac{1}{nT} \sum_{t=0}^{T-1} \sum_{i=1}^{n} \mathbb{E} \big[\| \nabla f(\boldsymbol{x}_{i}^{t}) \| \big] \\ & \leq O \Big(\frac{\Delta_{0}}{T} + \frac{\| \nabla f(\bar{\boldsymbol{x}}^{0}) \|}{\sqrt{T}} + \frac{\sigma}{n^{1 - \frac{1}{p}}} \cdot \frac{1}{T^{\frac{p-1}{2p}}} + \frac{1}{T^{\frac{1}{4}}} \sqrt{\frac{L\Delta_{0}}{1 - \lambda}} + \sqrt{\frac{3.5L\Delta_{0}}{(1 - \lambda)T}} + \frac{\sigma n^{\frac{1}{2}}}{(1 - \lambda)^{\frac{1}{p}}} \frac{1}{\sqrt{T}} + \\ & \frac{\| \nabla F(\mathbf{1}_{n} \otimes \bar{\boldsymbol{x}}^{0}) \|}{(1 - \lambda)n^{\frac{1}{2}}} \cdot \frac{1}{\sqrt{T}} + \frac{\sigma n^{\frac{1}{2}}}{1 - \lambda} \cdot \frac{1}{T^{\frac{2p-1}{2p}}} + \sqrt{\frac{n^{\frac{1}{2}}L\Delta_{0}}{(1 - \lambda)^{2}T}} \Big). \end{split}$$

Rearranging above terms leads to the desired upper bound.

C ADDITIONAL EXPERIMENT DETAILS

C.1 BASELINE DESCRIPTIONS

Please see Table 2 for detailed descriptions of baselines.

C.2 Additional details for synthetic experiments

Loss function. Let $(X_{i,k}, y_{i,k})$ denote the k-th sample of sub-dataset (X_i, y_i) on node i. The loss function of the considered nonconvex linear regression model on this sample is $\ell(y_{i,k} - X_{i,k} w_i^t)$, where the

$$\ell(r) = \begin{cases} \frac{c^2}{6} \left(1 - \left[1 - \left(\frac{r}{c}\right)^2\right]^3\right) & \text{if } |r| \leq c, \\ \frac{c^2}{6} & \text{otherwise} \end{cases},$$

and we use the suggested value c = 4.6851 in the robust statistics literature.

Hyperparameter tuning. Please see Table 3 for hyperparameter searching ranges for this experiment.

Hardware. We ran this experiment on Mac OS X 15.3, CPU M4 10 Cores, RAM 16GB.

C.3 Additional details for decentralized training of Transformers

Transformer architecture. We consider the following decoder-only Transformer model (GPT): vocabulary size is 10208, context length is 64, embedding size is 128, number of attention heads is 4, number of attention layers is 2, the linear projection dimension within attention block is 512, and LayerNorm is applied after the 2nd attention block. The total number of parameters of this model is 3018240.

Hyperparameter tuning. See Table 4 for our grid search range for algorithm hyperparameters.

Hardware. We simulate the distributed training on one NVIDIA H100 GPU, using PyTorch 3.2 with CUDA 12. The total hyperparameter search and training procedure took around 100 GPU hours.

Table 2: Summary of Baseline Methods

Method	Parallel update on node i	Hyper-parameters
DSGD	$oldsymbol{x}_i^{t+1} = \sum_{r=1}^n w_{ir} ig(oldsymbol{x}_r^t - lpha g_r(oldsymbol{x}_r^t, oldsymbol{\xi}_r^t) ig)$	α : constant stepsize
DSGD-GClip	$m{x}_i^{t+1} = \sum_{r=1}^n w_{ir} m{x}_r^t - lpha \operatorname{clip}(g_i(m{x}_i^t, m{\xi}_i^t), au)$	α, τ : stepsize α , and ℓ_2 clipping levels τ
DSGD-CClip	$oldsymbol{x}_i^{t+1} = \sum_{r=1}^n w_{ir} oldsymbol{x}_r^t - lpha \operatorname{clip}(g_i(oldsymbol{x}_i^t, oldsymbol{\xi}_i^t), au)$	α, τ : stepsize α , and component-wise clipping levels τ
DSGD-Clip	$m{x}_i^{t+1} = \sum_{r=1}^n w_{ir} m{x}_r^t - lpha_t \operatorname{Clip}(g_i(m{x}_i^t, m{\xi}_i^t), au_t)$	$lpha, au$: stepsize $lpha_t = lpha/(t+1),$ and ℓ_2 clipping levels $ au_t = au(t+1)^{2/5}$
GT-DSGD	$egin{aligned} oldsymbol{y}_i^{t+1} &= \sum_{r=1}^n w_{ir} ig(oldsymbol{y}_r^t + g_r(oldsymbol{x}_r^t, oldsymbol{\xi}_r^t) - g_r(oldsymbol{x}_r^{t-1}, oldsymbol{x}_r^{t+1}) \ oldsymbol{x}_i^{t+1} &= \sum_{r=1}^n w_{ir} ig(oldsymbol{x}_i^t - lpha oldsymbol{y}_r^{t+1}ig) \end{aligned}$	(ξ_r^{t-1}) α : constant stepsize
	$\boldsymbol{m}_{i}^{t+1} = \beta_{1} \boldsymbol{m}_{i}^{t} + (1 - \beta_{1}) \boldsymbol{s}_{i}^{t}$	
GT-Adam	$egin{aligned} m{v}_i^{t+1} &= \min \left(eta_2 m{v}_i^t + (1-eta_2) m{s}_i^t \odot m{s}_i^t, G ight) \ m{x}_i^{t+1} &= \sum_{r=1}^n w_{ir} m{x}_r^t - lpha rac{m{m}_i^{t+1}}{\sqrt{m{v}_i^{t+1} + \epsilon}} \ m{g}_i^{t+1} &= abla f_i(m{x}_i^{t+1}) \ m{s}_i^{t+1} &= \sum_{r=1}^n w_{ir} m{s}_r^t + m{g}_i^{t+1} - m{g}_i^t \end{aligned}$	α, G : constant stepsize α , and upper bound G , stabilization factor ϵ
QG-DSGDm	$egin{aligned} oldsymbol{m}_i^{t+1} &= eta \hat{oldsymbol{m}}_i^t + g_i(oldsymbol{x}_i^t, oldsymbol{\xi}_i^t) \ oldsymbol{x}_i^{t+1} &= \sum_{r=1}^n w_{ir} ig(oldsymbol{x}_i^t - \eta oldsymbol{m}_i^{t+1} ig) \ oldsymbol{d}_i^t &= (oldsymbol{x}_i^{t+1} - oldsymbol{x}_i^t)/\eta \ oldsymbol{\hat{m}}_i^{t+1} &= \mu \hat{oldsymbol{m}}_i^t + (1-\mu) oldsymbol{d}_i^t \end{aligned}$	η, β, μ : constant stepsize η , momentum parametes β, μ
SClip-EF-Net	$\begin{aligned} & \boldsymbol{m}_i^{t+1} = \beta_t \boldsymbol{m}_i^t + (1 - \beta_t) \boldsymbol{\Psi}_t(g_i(\boldsymbol{x}_i^t, \boldsymbol{\xi}_i^t) - \boldsymbol{m}_i^t) \\ & \boldsymbol{x}_i^{t+1} = \sum_{r=1}^n w_{ir} \big(\boldsymbol{x}_r^t - \alpha_t \boldsymbol{m}_r^{t+1} \big) \end{aligned}$	$c_{\varphi}, \tau, \alpha, \beta$: Component-wise smooth clipping operator: $\Psi_t(y) = \frac{c_{\varphi}}{\sqrt{t+1}} \frac{y}{\sqrt{y^2+\tau(t+1)^{3/4}}}$ stepsize $\alpha_t = \alpha/(t+1)^{1/5}$, momentum stepsize

1242 1243 1244

Table 3: Hyperparameter grid search in synthetic experiments

Method	Hyperparameter search set
DSGD	$\begin{array}{l}\alpha \in \{10^{-5}, 5*10^{-5}, 10^{-4}, 5*10^{-4}, 10^{-3}, 5*\\10^{-3}, 10^{-2}, 5*10^{-2}, 10^{-1}, 0.5, 1, 5, 10\}\end{array}$
DSGD-Clip	$\alpha \in \{10^{-5}, 5*10^{-5}, 10^{-4}, 5*10^{-4}, 10^{-3}, 5*10^{-3}, 10^{-2}, 5*10^{-2}, 10^{-1}, 0.5, 1, 5, 10\}, \tau \in \{10^{-3}, 5*10^{-3}, 10^{-2}, 5*10^{-2}, 5*10^{-2}, 10^{-1}, 0.5, 1, 5, 10, 50, 10^2\}$
GT-DSGD	$\alpha \in \{10^{-5}, 5*10^{-5}, 10^{-4}, 5*10^{-4}, 10^{-3}, 5*10^{-3}, 10^{-2}, 5*10^{-2}, 10^{-1}, 0.5, 1, 5, 10\}$
GT-NSGDm	$\alpha \in \{10^{-5}, 5*10^{-5}, 10^{-4}, 5*10^{-4}, 10^{-3}, 5*10^{-3}, 10^{-2}, 5*10^{-2}, 10^{-1}, 0.5, 1, 5, 10\}, \beta \in \{0.01, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 0.99\}$
SClip-EF-Netwo	$\begin{array}{l} \text{Drk} \in \{10^{-3}, 10^{-2}, 0.1, 1, 10, 30\}, \beta \in \\ \{10^{-2}, 0.1, 0.5, 0.8, 0.99\}, c_{\varphi} \in \\ \{1, 5, 10, 20, 30, 50\}, \tau \in \{0.1, 1, 10, 50, 100\} \end{array}$

126312641265

1266

12671268

Table 4: Hyperparameter grid search in decentralized training of Transformers

1	270
1	271
1	272

1283

1284

1290

Method Hyperparameter search set $\alpha \in \{10^{-4}, 5 * 10^{-4}, 10^{-3}, 5 * 10^{-3}, 10^{-2}, 5 * 10^{-3}, 10^{-2}, 5 * 10^{-3}, 10^{-2}$ DSGD $10^{-2}, 10^{-1}, 0.5, 1$ $\alpha \in \{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1, 10, 10^{2}\}, \tau \in$ DSGD-GClip $\{10^{-3}, 10^{-2}, 10^{-1}, 1, 10, 10^2\}$ $\alpha \in \{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1, 10, 10^{2}\}, \tau \in$ DSGD-CClip $\{10^{-3}, 10^{-2}, 10^{-1}, 1, 10, 10^2\}$ $\alpha \in \{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1, 10, 10^{2}\}, \tau \in$ DSGD-Clip $\{10^{-3}, 10^{-2}, 10^{-1}, 1, 10, 10^2\}$ $\alpha \in \{10^{-4}, 5*10^{-4}, 10^{-3}, 5*10^{-3}, 10^{-2}, 5*$ GT-DSGD $10^{-2}, 10^{-1}, 0.5, 1$ $\alpha \in \{5 * 10^{-5}, 10^{-4}, 5 * 10^{-4}, 10^{-3}, 5 * 10^{-4}, 10^{-3}, 10^{-4}, 10$ GT-Adam $10^{-3}, 10^{-2}, 5*10^{-2}, 10^{-1}, 0.5, 1, 5, 10\}, G \in$ $\{10^{-3},10^{-2},10^{-1},1,10\},\epsilon=10^{-8}$ $\eta \in \{5 * 10^{-5}, 10^{-4}, 5 * 10^{-4}, 10^{-3}, 5 * 10^{-4}, 10^{-3}, 5 * 10^{-4}, 10^{-3}, 10^{-4}, 10^{-4}, 10^{-3}, 10^{-4}$ QG-DSGDm 10^{-3} , 10^{-2} , $5*10^{-2}$, 10^{-1} , 0.5, 1, 5, 10}, $\beta = \mu \in$ $\{0.01, 0.2, 0.4, 0.6, 0.8, 0.99\}$ $\alpha \in \{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1, 10\}, \beta \in$ GT-NSGDm $\{0.01, 0.2, 0.4, 0.6, 0.8, 0.99\}$ SClip-EF-Network $\in \{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 10^{0}, 10^{1}, 10^{2}\}, \beta \in \{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 10^{0}, 10^{1}, 10^{2}\}, \beta \in \{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 10^{0}, 10^{1}, 10^{2}\}, \beta \in \{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 10^{0}, 10^{1}, 10^{2}\}, \beta \in \{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 10^{0}, 10^{1}, 10^{2}\}, \beta \in \{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 10^{0}, 10^{1}, 10^{2}\}, \beta \in \{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 10^{0}, 10^{1}, 10^{2}\}, \beta \in \{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 10^{0}, 10^{1}, 10^{2}\}, \beta \in \{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 10^{0}, 10^{1}, 10^{2}\}, \beta \in \{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 10^{0}, 10^{1}, 10^{2}\}, \beta \in \{10^{-4}, 10^{-2}, 10^{-1}, 10^{0}, 10^{1}, 10^{2}\}, \beta \in \{10^{-4}, 10^{-2}, 10^{-1}, 10^{0}, 10^{1}, 10^{2}\}, \beta \in \{10^{-4}, 10^{-2}, 10^{-2}, 10^{-1}, 10^{0}, 10^{1}, 10^{2}\}, \beta \in \{10^{-4}, 10^{-2}, 10^{-2}, 10^{-2}, 10^{-2}, 10^{-2}, 10^{-2}, 10^{-2}, 10^{-2}, 10^{-2}, 10^{$ $\{0.01, 0.1, 1, 10\}$