

BIQ: Bisection Interval Quantization for Communication-efficient Federated Learning

Luyang Gai^{1, 2}, Shusen Yang^{1, 2*}, Xuebin Ren^{1, 3}, Zihao Zhou^{1, 2}

¹National Engineering Laboratory for Big Data Analytics, Xi'an Jiaotong University, China

²School of Mathematics and Statistics, Xi'an Jiaotong University, China

³School of Computer Science and Technology, Xi'an Jiaotong University, China

wnt0954@gmail.com, shusenyang@mail.xjtu.edu.cn, xuebinren@mail.xjtu.edu.cn, zihaozhou19@gmail.com

Abstract

Quantization is a pivotal technique for enhancing communication efficiency in Federated Learning (FL). Traditional quantization methods often set uniform intervals, may fail to adequately characterize non-uniform data distributions, thus leading to substantial estimation errors and degraded model performance. Non-uniform quantization can better solve the problem. However, when applied to FL, it would bring additional communication overheads for the alignment of parameter distributions among distributed models. To address this issue, we propose Bisection Interval Quantization (BIQ), a novel non-uniform quantization framework for FL with great communication efficiency. In particular, BIQ works by optimizing the interval selection through recursive bisection among distributed clients without extra parameter communication. For scenarios involving amounts of boundary inputs, we further design Weighted Bisection Interval Quantization (WBIQ), which incorporates maximum likelihood estimation to refine boundary value reconstruction to enhance the estimation quality of boundary inputs. Our theoretical analysis rigorously establishes, for the first time under biased quantization conditions, that both BIQ and WBIQ achieve tighter error bounds and enhanced stability. Extensive experiments validate that both BIQ and WBIQ significantly accelerate the convergence of FL training when compared to the state-of-the-art quantizers under both convex and non-convex settings.

Code and Extended version —

<https://github.com/GaiLuyang/BIQ>

Introduction

Federated Learning (FL) (McMahan et al. 2017) has emerged as a promising paradigm for distributed machine learning, particularly in privacy-sensitive domains such as healthcare (Pati et al. 2022; Boscarino et al. 2022) and finance (Cui et al. 2021). FL enables collaborative model training across distributed clients while preserving data privacy. However, its practical deployment is hindered by communication bottlenecks caused by frequent parameter exchanges. To address it, quantization (Elgabli et al. 2020; Reiszadeh et al. 2020; Sun et al. 2022) has become a cornerstone for reducing communication overhead, and most al-

gorithms predominantly rely on uniform quantizers including Stochastic Quantization (SQ) (Elgabli et al. 2020; Reiszadeh et al. 2020; Gupta et al. 2015; Alistarh et al. 2017) and Rounding Quantization (RQ) (Sun et al. 2022; Gupta et al. 2015; Bai, Wang, and Liberty 2018; Nagel et al. 2022). However, these methods enforce equal intervals that fail to adapt to most neural networks with non-uniform parameter distributions (Han, Mao, and Dally 2015; Zhang et al. 2018; Jung et al. 2019) (e.g., such as bell-shaped and long-tail distributions observed in weights and activations (Gongyo et al. 2024; Li, Dong, and Wang 2019)), resulting in substantial information loss and degraded convergence.

Non-uniform quantization (Zhang et al. 2018; Chen et al. 2023; Luqman, Qazi, and Khan 2024; Wang et al. 2022) offers superior accuracy compared to uniform methods by dynamically mapping the parameters to non-uniformly spaced quantized values to match non-uniform distributions. Prior works, such as loss-driven adaptive quantization (Jung et al. 2019), power-of-two quantization (Li, Dong, and Wang 2019), logarithmic quantization (Miyashita, Lee, and Murmann 2016) and iterative range-level updates (Luqman, Qazi, and Khan 2024), have demonstrated its effectiveness in centralized settings. However, non-uniform quantization poses an inherent deployment challenge to FL due to the nature of misalignment of adaptive quantization schemes across heterogeneous clients. One solution is that clients upload both quantized parameters and their respective quantization schemes (Chen et al. 2023). However, this approach introduces extra communication and computational costs, thus undermining the communication benefits of quantization. What's more, additional parameters may expose information such as data distribution of clients, resulting in privacy leakage. Therefore, there is a pressing need to design non-uniform quantization schemes for FL to improve model accuracy without incurring additional communication costs.

The key insight that inspires our work lies in the object of quantization. Quantization in FL has been largely inherited from model quantization, where the quantized outputs must retain semantic meaning because they are directly used in training and inference. However, in the context of FL, quantization serves solely as a communication compression mechanism, as illustrated in Figure 1. This inspires us to design novel encoding-decoding mechanisms tailored specifically for the communication process. As a result, quanti-

*Corresponding author.

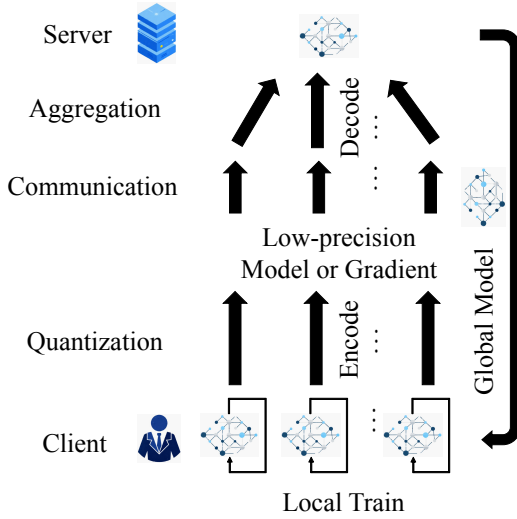


Figure 1: Quantification in FL. Quantization applies to client encoding and server decoding processes.

zation in FL opens up a broader design space compared to traditional model quantization. We revisit the fundamental interpretation of quantization levels: in most prior works, quantization levels serve as symbolic representations of data points. This naturally raises a question: Can we assign quantization a new meaning?

To address the problem, we try to quantize the process of quantization. Building on this, we propose Bisection Interval Quantization (BIQ), a novel non-uniform quantizer that recursively bisects intervals using binary codes. Each bit in BIQ encodes a client’s local bisection path, eliminating the need for explicit interval alignment. For boundary-critical scenarios (e.g., extreme gradient values), we further design Weighted BIQ (WBIQ), which refines boundary reconstruction through maximum likelihood estimation based on bit-count statistics. Crucially, both methods operate without increasing communication overhead.

However, in terms of theoretical analysis, BIQ and WBIQ sacrifice unbiasedness for higher quantization accuracy. Theoretical guarantees for biased quantizers in FL have remained elusive, as prior convergence analyses (Elgabri et al. 2020; Reisizadeh et al. 2020) strictly require unbiased quantization. We bridge this gap by establishing the first convergence bounds for biased quantizers under both strongly convex and non-convex objectives.

Our Contribution. In summary, our contributions are as follows: First, we propose a novel BIQ algorithm for FL, which achieves higher quantization accuracy without sacrificing communication cost. Second, we propose its variant WBIQ to further improve quantization accuracy. WBIQ incorporates maximum likelihood estimation to better capture input values near the boundaries of the quantization range. Third, we establish the first convergence analysis for biased quantizers in FL. Under appropriate control of bias and noise, we demonstrate that BIQ and WBIQ converge

at rates of $\mathcal{O}\left(\frac{1}{T}\right)$ under strong convexity assumptions and $\mathcal{O}\left(\frac{1}{\sqrt{T}}\right)$ in non-convex settings, respectively. Fourth, extensive experiments on four datasets validate that BIQ and WBIQ achieve more accurate quantization, faster convergence and higher model accuracy compared with existing quantization methods under the same resolution.

Our work redefines the role of quantization and establishes a general analytical paradigm that accommodates diverse quantization strategies in FL. Our results unlock the potential of high-precision FL in bandwidth-limited scenarios.

Problem Definition

We first introduce FL and present the quantization-based communication compression method.

Federated Learning (FL). In FL, a server collaborates with clients $\mathcal{N} = \{1, 2, \dots, N\}$ to jointly train a model by solving the following Empirical Risk Minimization (ERM) problem:

$$\min_{\theta} f(\theta) \triangleq \min_{\theta} \sum_{n=1}^N \frac{M_n}{N} f_n(\theta; \mathcal{M}_n), \quad (1)$$

where $\theta \in R^d$ denotes the model parameter. \mathcal{M}_n is the local dataset with M_n samples. f_n represents the loss function of the n -th client.

SGD-Based Optimizer in FL with Quantization. In FL, the server distributes the initialized model θ^0 to each client. To enhance local training efficiency, the n -th client, when participating in the k -th communication round, randomly samples a mini-batch $\hat{\mathcal{M}}_n$ from its local dataset \mathcal{M}_n and updates its local model according to the following procedure

$$\theta_n^{k,t+1} = \theta_n^{k,t} - \alpha \hat{\nabla} f_n(\theta_n^{k,t}; \hat{\mathcal{M}}_n), \quad (2)$$

where $t = 0, 1, \dots, \tau - 1$ denotes the local update rounds. After completing these updates, the client uploads the difference $\theta_n^{k,\tau} - \theta_n^{k,t}$ to the server for aggregation. The communication bottleneck occurs during the update transmission process, as limited bandwidth can lead to congestion when transmitting full updates. We quantize the uploaded data to reduce the communication overhead. Additionally, resource heterogeneity poses another challenge in FL. Some clients may not respond promptly to the server’s update requests (Reisizadeh et al. 2020; Sun et al. 2022), resulting in inefficient training and amplifying the impact of limited communication bandwidth. To address these issues, at each communication round, we select a subset of clients $S_k \subset \mathcal{N}$. These clients quantize the model updates by quantizer $Q(\cdot)$ and upload the quantized models for aggregation, which efficiently alleviates the communication bottleneck and improves training efficiency. The aggregated process is as follows

$$\theta^{k+1} = \theta^k + \frac{1}{s} \sum_{n \in S_k} Q(\theta_n^{k,\tau} - \theta_n^{k,t}), \quad (3)$$

where $|S_k| = s$. Equation (3) directly impacts the effectiveness of model aggregation. At lower quantization resolutions, quantization significantly reduces communication

Algorithm 1: BIQ encoder process.

Input: $\mathbf{x} \in \mathbb{R}^d, b, R$

Output: B

```

1:  $q_L = -R \cdot \mathbf{1}, q_R = R \cdot \mathbf{1}$ 
2: for  $j = 1, 2, \dots, d$  do
3:   for  $i = b-1, b-2, \dots, 0$  do
4:     if  $x^{(j)} \leq \frac{q_L^{(j)} + q_R^{(j)}}{2}$  then
5:        $q_R^{(j)} = \frac{q_L^{(j)} + q_R^{(j)}}{2}, B_i^{(j)} = 0$ 
6:     else
7:        $q_L^{(j)} = \frac{q_L^{(j)} + q_R^{(j)}}{2}, B_i^{(j)} = 1$ 
8:     end if
9:   end for
10:   $B^{(j)} = B_{b-1}^{(j)} B_{b-2}^{(j)} \dots B_0^{(j)}$ 
11: end for
```

costs at the expense of increased errors. Therefore, designing high-precision quantizers at low quantization resolutions is an effective approach to achieve a better trade-off between communication efficiency and utility.

Bisection Interval Quantization

In this section, we first propose BIQ and WBIQ. Then, we present error analysis and convergence analysis.

BIQ Encoding Process

BIQ encodes the selection process of quantized values using binary system. Specifically, given a quantization resolution b , a quantization range R and an input $x \in \mathbb{R}^d$ with each coordinate satisfying $x^{(j)} \in [-R, R]$, BIQ recursively partitions the interval $[-R, R]$ into subintervals by bisection. At each step, a bit ("0" or "1") indicates whether the input lies in the left or right subinterval, respectively. BIQ performs the following process when determining the encoded value for element $x^{(j)}$: if the input lies in the left half of the current quantization interval, the corresponding bit is set to $B_j^{(j)} = 0$; if it lies in the right half, the bit is set to $B_j^{(j)} = 1$. Let $B = (B_{b-1}^{(1)} B_{b-2}^{(1)} \dots B_0^{(1)}, B_{b-1}^{(2)} B_{b-2}^{(2)} \dots B_0^{(2)}, \dots, B_{b-1}^{(d)} B_{b-2}^{(d)} \dots B_0^{(d)})$ represent the quantized binary code generated by BIQ. Algorithm 1 describes the encoding mechanism, where $\mathbf{1} = (1, 1, \dots, 1) \in \mathbb{R}^d$.

BIQ Decoding Process

The client transmits the binary code B to the server, which decodes each coordinate $B^{(j)}$ according to its encoding scheme. Specifically, the server processes the bits in $B^{(j)}$ from left to right, where a "0" indicates that the left half of the current interval is retained, and a "1" indicates that the right half is retained. The decoding process of BIQ results in an interval $[q_{B^{(j)}}^L, q_{B^{(j)}}^R]$. However, the aggregation step requires a deterministic value. To minimize quantization error, BIQ selects the midpoint of the interval as the final output:

$$q_{B^{(j)}} = \frac{q_{B^{(j)}}^L + q_{B^{(j)}}^R}{2}.$$

Algorithm 2: FedBIQ.

Input: N clients with datasets $\{\mathcal{M}_n\}_{n=1}^N$, initialize server model parameters θ^0 , communication rounds K , quantization resolution b , learning rate α .

Output: Global model θ^*

```

1: Server broadcasts  $\theta^0$  to all clients.
2: for  $k = 0, 1, \dots, K-1$  do
3:   Server broadcasts  $\theta^k$  to clients  $S^k$  who are selected uniformly at random.
4:   for client  $n \in S^k$  in parallel do
5:     for  $t = 0, 1, \dots, \tau-1$  do
6:       Local update via Equation (2)
7:     end for
8:      $B_n^k \leftarrow \text{BIQ}(\theta_n^{k,\tau} - \theta^k, b, R_n^k)$ 
9:     Send  $B_n^k, R_n^k$  to Server.
10:   end for
11:   Server decodes quantized signals and aggregates via Equation (3).
12: end for
```

BIQ Variant: WBIQ. Since BIQ uses the midpoint of the interval as its output, it introduces gaps when representing inputs that are close to the boundaries of the quantization range. To address this limitation while maintaining the same communication cost, we propose an enhanced scheme called WBIQ. WBIQ aims to more accurately handle the quantization of inputs near the range boundaries without increasing the communication overhead.

Definition 1 (Counting). For any $c = c_{b-1}c_{b-2}\dots c_0, c_i \in \{0, 1\}$, we define the functions $\delta_j : \{0, 1\}^b \rightarrow \mathbb{Z}_{\geq 0}, j \in \{0, 1\}$. δ_j count the number of j in c . Specifically, $\delta_0(c) = b - \sum_{i=1}^b c_i, \delta_1(c) = \sum_{i=1}^b c_i$.

Our improvement is inspired by the principle of maximum likelihood estimation. In this context, the bits "0" and "1" represent the number of times the left and right intervals were selected during the bisection process, respectively. If "1" appears more frequently than "0" in $B^{(j)}$, then according to maximum likelihood estimation, the input value is more likely to lie closer to the right end of the interval. This insight motivates our design of WBIQ, where we incorporate the number of "0"s and "1"s in $B^{(j)}$ to refine the quantized output.

In WBIQ, the quantization interval $[q_{B^{(j)}}^L, q_{B^{(j)}}^R]$ is determined in the same manner as in standard BIQ. However, WBIQ further leverages the maximum likelihood estimation principle by performing a weighted average of the interval endpoints. The weights are determined by the relative frequencies of "0" and "1". Specifically, the final output is given by $q_{B^{(j)}} = \frac{\delta_0(B^{(j)})}{b} q_{B^{(j)}}^L + \frac{\delta_1(B^{(j)})}{b} q_{B^{(j)}}^R$. This maximum likelihood-based approach improves the quantization accuracy for inputs near the boundaries of the range, thereby enhancing the overall performance of the algorithm.

Error Analysis

In this section, we analyze the quantization error introduced by BIQ and WBIQ. Let C, b denote the quantization range

and the quantization resolution respectively. After decoding, the server reconstructs the interval $[q_{B(j)}^L, q_{B(j)}^R]$. The true input value x_j can lie anywhere within this interval. Assuming without loss of generality that x_j is uniformly distributed over $[q_{B(j)}^L, q_{B(j)}^R]$, then the quantization error is $\eta_j^{BIQ} = x_j - BIQ(x_j) \sim \mathcal{U}(-\frac{q_{B(j)}^R - q_{B(j)}^L}{2}, \frac{q_{B(j)}^R - q_{B(j)}^L}{2})$. The variance of the quantization error is $D\eta_j^{BIQ} = \frac{C^2}{12 \cdot 2^{2b}}$. In contrast, for SQ, $D\eta_j^{SQ} = \frac{C^2}{6 \cdot (2^b - 1)^2}$. Comparing the two error variances yields:

$$\frac{D\eta_j^{BIQ}}{D\eta_j^{SQ}} = \frac{6 \cdot (2^b - 1)^2}{12 \cdot 2^{2b}} \leq 0.5, \quad \forall b.$$

Under the same quantization resolution b , the variance of the quantization error for BIQ is less than half of SQ, indicating that BIQ produces more stable and reliable quantized outputs. For WBIQ, we have $D\eta_j^{WBIQ} \leq D\eta_j^{SQ}$, which still ensures greater stability in the quantized results compared to SQ. In terms of error bounds, the range of the quantization error satisfies $|\eta_j^{BIQ}| \leq \frac{C}{2^b}$, $|\eta_j^{WBIQ}| \leq \frac{2C}{2^b}$. While for SQ, the bound is $|\eta_j^{SQ}| \leq \frac{2C}{2^b - 1}$. These bounds demonstrate that BIQ achieves a tighter error bound than SQ, resulting in more accurate quantized values. To validate these theoretical findings, we conduct sampling experiments with various quantizers. We show sampling results for uniform, bell-shaped and long-tail data distributions. As shown in Section **Sampling** (please refer to Other Experiments in supplementary materials), the empirical results align closely with our theoretical analysis, confirming the superior performance of BIQ and WBIQ in practical settings.

Convergence Analysis

BIQ is a biased quantizer, thus failing to satisfy the unbiased property required by most convergence analysis frameworks (Elgabli et al. 2020; Reiszadeh et al. 2020). To address this challenge, we develop a novel convergence analysis for biased quantizers under both strongly convex and non-convex settings. Proofs of both theorems are provided in supplementary materials. Our analysis builds on the following assumptions.

Assumption 1. Let $Q(\cdot)$ be a quantizer such that $\mathbb{E}[Q(\mathbf{x})|\mathbf{x}] = \mathbf{x} + \eta$, $\eta \in \mathbb{R}^d$, $\|\eta\| \leq G$, $G > 0$. The quantization variance is bounded as follows: $\mathbb{E}[\|Q(\mathbf{x}) - \mathbb{E}[Q(\mathbf{x})]\|^2|\mathbf{x}] \leq q\|\mathbf{x}\|^2$, $q > 0$.

Assumption 2. f_i is L -smooth if there exists a constant $L > 0$ such that $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, $\|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|$.

Assumption 3. The stochastic gradient $\hat{\nabla} f_i(\mathbf{x})$ is unbiased, i.e. $\mathbb{E}_\zeta[\hat{\nabla} f_i(\mathbf{x})] = \nabla f_i(\mathbf{x})$ with bounded variance, i.e. $\mathbb{E}_\zeta[\|\hat{\nabla} f_i(\mathbf{x}) - \nabla f_i(\mathbf{x})\|^2] \leq \sigma^2$.

Assumption 4. f_i is μ -strongly convex if there exists a constant $\mu > 0$ such that for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, $\langle \nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq \mu\|\mathbf{x} - \mathbf{y}\|^2$.

Theorem 1 (Convergence of Strongly Convex Settings). Suppose Assumptions 1 to 4 hold. Define $A_0 :=$

$\frac{2qL}{\mu} + \frac{4L^2(N-s)(q+4)}{s(N-1)\mu}$. Set a constant k_0 that satisfies $k_0 \geq \max\{\frac{48\tau-1}{\tau}, \frac{48N-\mu^2\tau}{\mu^2\tau^2}, \frac{48L^2A_0\tau-\mu^3\tau+6L^2\mu}{\mu^2\tau^2-L^2\mu\tau}, \frac{4L^2-\mu^2}{\mu^2\tau}, \frac{4L\sqrt{\tau(\tau-1)-\mu}}{\mu\tau}\}$. Then for any $k \geq k_0$, if the learning rate satisfies $\alpha_k = \min\{\frac{\mu}{L^2}, \frac{1}{L\sqrt{\tau(\tau-1)}}, \frac{1}{\mu\tau}\}$, Algorithm 2 holds

$$\begin{aligned} \mathbb{E}\|\theta^k - \theta^*\|^2 &\leq \left(\frac{k_0\tau+1}{k\tau+1}\right)^2 \mathbb{E}\|\theta^{k_0} - \theta^*\|^2 + A_2 \frac{1}{\tau} \cdot \frac{\tau}{k\tau+1} \\ &\quad + A_3 \cdot \frac{\tau}{k\tau+1} + A_4 \frac{g^{k^2}}{\tau^2} \cdot \frac{\tau}{k\tau+1}. \end{aligned} \quad (4)$$

where $\{g^k\}$ is a bounded positive sequence. The constants in Theorem 1 are defined as

$$\begin{aligned} A_2 &:= \frac{32\sigma^2}{\mu^2} \left[\frac{e}{N} (L^2 + 1 + 2qN) + \frac{4e(N-s)(q+4)}{s(N-1)} \right], \\ A_3 &:= \frac{32\sigma^2}{\mu^2} \left[\frac{1}{N} + 2q + \frac{4(N-s)(q+4)}{s(N-1)} \right], \\ A_4 &:= \frac{3sN + 12N - 15s}{s(N-1)}. \end{aligned}$$

Theorem 1 establishes convergence guarantees under strongly convex assumption. The bias η introduced by the biased quantizer affects the convergence rate through its influence on g^k . In FedBIQ, by choosing the quantization range as $R_n^k \leq (\frac{12}{d} \cdot 2^{2b} \|\theta_n^{k,\tau} - \theta_n^k\|^2)^{\frac{1}{2}}$, we ensure that $G_n^k \leq \frac{R_n^k}{2^b}$. Let us find $g^k = \max_{n \in [N]} \{(k\tau+1)G_n^k\}$. This implies that the sequence g^k is bounded, thereby ensuring that FedBIQ achieves a convergence rate of $\mathcal{O}(\frac{1}{T})$. For FedWBIQ, we choose the quantization range as $R_n^k \leq (\frac{48}{d} \cdot 2^{2b} \|\theta_n^{k,\tau} - \theta_n^k\|^2)^{\frac{1}{2}}$ to reach the same convergence rate.

Theorem 2 (Convergence of Non-convex Settings). Under Assumptions 1 to 3 and the update method of Algorithm 2, we define the following two constants: $F_1 = 2L^2\tau(\tau-1)$, $F_2 = L + \frac{N+q+3}{N}\tau L + \frac{4L(N-s)}{s(N-1)}(q+4)\tau$. When the constraint $T \geq \frac{F_1^2}{(\sqrt{F_2^2+2F_1-F_2})^2}$ is satisfied and the learning rate is $\alpha_k = \frac{1}{\sqrt{T}}$, the following first-order stability condition holds

$$\begin{aligned} \frac{1}{T} \sum_{k=0}^{K-1} \sum_{t=0}^{\tau-1} \mathbb{E}\|\nabla f(\bar{\theta}^{k,t})\|^2 &\leq \frac{2[f(\theta^0) - f^*]}{\sqrt{T}} + H_1 \frac{\tau-1}{T} \\ &\quad + H_2 \frac{\tau}{\sqrt{T}} + H_3 \frac{1}{\sqrt{T}}. \end{aligned} \quad (5)$$

where $g^k \leq \max_{n \in [N]} \{G_n^k\}$, $\|\eta_n^k\| \leq G_n^k$,

$$\begin{aligned} H_1 &:= \frac{L\sigma^2(N+1)}{N}, \\ H_2 &:= L\sigma^2 \left[\frac{(N+q+3)}{N} + \frac{4(N-s)}{s(N-1)}(q+4) + \frac{1}{N\tau} \right], \\ H_3 &:= \sum_{k=0}^{K-1} \left[\frac{12L(N-s)}{s(N-1)} g^k + \frac{(3N+1)L}{N} g^k \right. \\ &\quad \left. + \|\nabla f(\bar{\theta}^{k,\tau})\| \right] g^k. \end{aligned}$$

Theorem 2 provides convergence guarantees under non-convex assumption. In FedBIQ, to satisfy the condition $T \geq \frac{F_1^2}{(\sqrt{F_2^2 + 2F_1 - F_2})^2}$, we select the quantization range as $R_n^k \leq \{12 \cdot 2^{2b} \cdot \frac{1}{d}[\tau L + \frac{4L(N-s)\tau}{S(N-1)}]^{-1}[\sqrt{T} - \frac{1}{2\sqrt{T}}F_1 - L - \frac{N+3}{N}\tau L - \frac{16L(N-s)\tau}{s(N-1)}]\|\theta_n^{k,r} - \theta_n^k\|^2\}^{\frac{1}{2}}$, which ensures that $G_n^k \leq \frac{R_n^k}{2^b}$. Let us find $g^k = \max_{n \in [N]} \{G_n^k\}$. This implies that g^k is bounded. We set $\frac{1}{2}R_n^k$ as the quantization range for FedWBIQ. FedBIQ and FedWBIQ converges at a rate of $\mathcal{O}\left(\frac{1}{\sqrt{T}}\right)$. Notably, when employing biased quantizers, the convergence rate can be preserved as long as the bias is appropriately controlled. Moreover, if the quantizer is unbiased, i.e. $\eta = 0$, then $g^k = 0, \forall k$. The results in Theorems 1 and 2 recover the conclusions in (Reisizadeh et al. 2020), showing that our analysis generalizes existing theory.

Experiments

In this section, we conducted extensive experiments on BIQ and WBIQ. First, we compare six algorithms on four datasets, demonstrating the superior performance of BIQ in terms of cost, communication efficiency and utility. Finally, we perform quantizer sampling experiments to evaluate the accuracy of different quantization methods under low-resolution settings in supplementary materials.

Data and Models. Our experiments are conducted in a computing environment equipped with an AMD Ryzen 7 9700X CPU and an RTX 4070 Ti Super GPU. We evaluate our methods on four benchmark datasets: CDC Diabetes Health Indicators (abbreviated as CDC for convenience) (abbreviated as CDC for convenience), MNIST (Lecun et al. 1998), CIFAR-10 (Krizhevsky, Nair, and Hinton 2009) and Tiny ImageNet. CDC contains 253,680 samples consisting of 21 features and 2 labels, 70% of which are divided for training and the rest for evaluation. MNIST consists of 60,000 grayscale training images of size 28*28 pixels and 10,000 test samples, while CIFAR-10 contains 50,000 images of size 32*32 pixels across three channels, along with an equal-sized test set. MNIST and CIFAR-10 are labeled over 10 classes. To assess the practical effectiveness of BIQ on real-world datasets, we conducted comprehensive experiments on Tiny ImageNet, which comprises 100K training and 10K test images spanning 200 classes of 64*64 RGB images.

To study the impact of data heterogeneity, we adopt two data partitioning strategies: (1) I.I.D. setting: The training data is uniformly distributed across clients, ensuring balanced class distributions; (2) Non-I.I.D. setting: Data is partitioned in a heterogeneous manner using a Dirichlet distribution with concentration parameter $\alpha = 0.6$. The global test set is centrally maintained by the server to ensure a consistent and standardized evaluation of model performance across all experimental configurations.

For strongly convex settings, we build logistic regression models to solve the binary classification problem on CDC. For non-convex settings, we employ a two-layer CNN model on MNIST and CIFAR-10. Each layer of the CNN model

	CDC	MNIST	CIFAR-10	Tiny ImageNet
CR	30	30	2000	50
LU	15	15	15	30
LR	0.01	0.03	0.05	0.007
MO	0.5	0.5	0.5	0.0002
BS	32	32	32	128
OP	SGD	SGD	SGD	Adam

Table 1: Experiment Settings. CR: communication rounds; LU: the number of rounds in a local update; LR: learning rate; MO: momentum; BS: batch size; OP: optimizer.

consists of convolutional layers, ReLU activation, and max-pooling, followed by two fully connected layers. For Tiny ImageNet, we employed Resnet50.

Baselines. Our work conducts a systematic comparison among BIQ, WBIQ and existing quantization-based FL methods. The baseline algorithms include: FedAvg (McMahan et al. 2017), FedSQ which incorporates stochastic quantization (Gupta et al. 2015), FedRQ which applies rounding quantization (Gupta et al. 2015), FedPAQ (Reisizadeh et al. 2020) which is based on QSGD quantization (Alistarh et al. 2017), FedNQFL (Chen et al. 2023) which incorporates nonuniform quantization, FedSQFL (Marnissi, El Hammouti, and Bergou 2024) which integrates quantization and sparsification.

Criterion. In this experiment, we establish a multidimensional evaluation framework for BIQ and WBIQ. First, we develop a cost model to quantify algorithmic efficiency. We use the shifted-exponential model (Reisizadeh et al. 2020; Lee et al. 2018) to fit the client-side encode and gradient computation time as the local computation cost. Communication cost is defined as the ratio of the total number of bits received by the server per round to the available bandwidth (BW), where BW is set to 1.5MB/s for CDC and MNIST, and 3MB/s for CIFAR-10 and Tiny ImageNet. The total cost combines both local computation and communication costs to represent the overall iteration overhead. For utility assessment, we focus on global model performance, evaluating convergence and generalization through two metrics: cross-entropy loss and classification accuracy on the test set. This dual-metric approach provides a comprehensive measure of model effectiveness. Our experiments assess algorithm optimization performance from both resource efficiency and model utility perspectives.

Numerical Results and Discussions

End-to-end Comparison. We conducted experiments on CDC, MNIST, CIFAR-10 and Tiny ImageNet. In each round, 15 clients were randomly sampled from a total of 80 for training. The specific settings are presented in Table 1.

FedAvg was implemented using 32-bit floating-point precision, while other methods were fixed to a quantization resolution of $b = 3$ bits. Figure 2 shows the variation of test loss ($Loss$) and accuracy (Acc) with the total cost ($Time$) on CDC with 1-sigma error bars, averaged over 5 runs. Table 2 report the average training cost per communication round

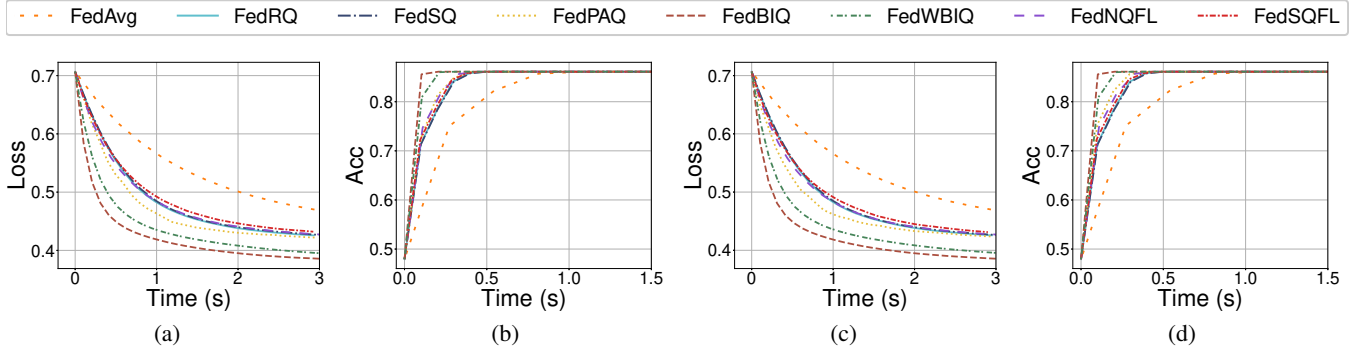


Figure 2: Loss and classification accuracy evolution on CDC under I.I.D. (left) and Non-I.I.D. (right) data distributions.

Dataset	Algorithm	I.I.D. Distribution			Non-I.I.D. Distribution		
		$cost_{avg}$	$Loss$	Acc (%)	$cost_{avg}$	$Loss$	Acc (%)
MNIST	FedAvg	$100.06 \pm 3e-3$	$0.15 \pm 6e-4$	$95.58 \pm 1e-3$	$100.06 \pm 2e-3$	$0.17 \pm 2e-3$	$94.86 \pm 9e-4$
	FedRQ	$9.46 \pm 2e-3$	$0.38 \pm 1e-3$	$90.32 \pm 1e-3$	$9.46 \pm 1e-3$	$0.58 \pm 9e-3$	$86.82 \pm 8e-3$
	FedSQ	$9.47 \pm 4e-3$	$0.38 \pm 1e-3$	$90.30 \pm 3e-4$	$9.47 \pm 1e-3$	$0.57 \pm 6e-3$	$87.10 \pm 5e-3$
	FedPAQ	$9.47 \pm 1e-3$	$0.31 \pm 1e-2$	$91.23 \pm 5e-3$	$9.47 \pm 2e-3$	$0.43 \pm 1e-2$	$87.57 \pm 8e-3$
	FedNQFL	$12.05 \pm 1e-3$	$0.17 \pm 3e-3$	$94.96 \pm 1e-3$	$11.80 \pm 3e-3$	$0.20 \pm 4e-3$	$94.14 \pm 2e-3$
	FedSQFL	$5.73 \pm 2e-3$	$0.30 \pm 1e-2$	$91.00 \pm 5e-3$	$5.73 \pm 2e-3$	$0.35 \pm 2e-2$	$89.85 \pm 6e-3$
	FedBIQ	$9.47 \pm 3e-3$	$0.16 \pm 1e-3$	$95.22 \pm 7e-4$	$9.47 \pm 2e-3$	$0.19 \pm 2e-3$	$94.39 \pm 2e-3$
	FedWBIQ	$9.48 \pm 3e-3$	$0.16 \pm 3e-3$	$95.37 \pm 7e-4$	$9.48 \pm 2e-3$	$0.18 \pm 2e-3$	$94.58 \pm 4e-4$
CIFAR-10	FedAvg	$615.00 \pm 9e-4$	$0.63 \pm 7e-3$	$78.95 \pm 4e-3$	$615.00 \pm 1e-3$	$0.70 \pm 6e-3$	$76.51 \pm 5e-3$
	FedRQ	$57.76 \pm 2e-3$	$0.88 \pm 4e-3$	$69.86 \pm 2e-3$	$57.76 \pm 9e-4$	$0.95 \pm 3e-3$	$67.03 \pm 1e-3$
	FedSQ	$57.77 \pm 8e-4$	$0.85 \pm 3e-3$	$70.67 \pm 2e-3$	$57.76 \pm 2e-3$	$0.93 \pm 4e-3$	$67.96 \pm 2e-3$
	FedPAQ	$57.76 \pm 1e-3$	$0.85 \pm 2e-3$	$70.83 \pm 2e-3$	$57.77 \pm 3e-3$	$0.93 \pm 3e-3$	$67.96 \pm 9e-4$
	FedNQFL	$60.29 \pm 3e-4$	$0.68 \pm 4e-3$	$77.05 \pm 3e-3$	$60.09 \pm 4e-4$	$0.75 \pm 4e-3$	$74.48 \pm 3e-3$
	FedSQFL	$23.96 \pm 1e-3$	$0.97 \pm 3e-3$	$66.56 \pm 1e-3$	$23.96 \pm 8e-4$	$1.03 \pm 3e-3$	$63.68 \pm 2e-3$
	FedBIQ	$57.76 \pm 3e-4$	$0.67 \pm 1e-2$	$77.49 \pm 2e-3$	$57.77 \pm 1e-3$	$0.75 \pm 1e-2$	$74.58 \pm 3e-3$
	FedWBIQ	$57.77 \pm 3e-3$	$0.65 \pm 5e-3$	$78.06 \pm 2e-3$	$57.77 \pm 9e-4$	$0.72 \pm 1e-2$	$75.56 \pm 4e-3$

Table 2: End-to-end comparison results on MNIST and CIFAR-10.

($cost_{avg}$), $Loss$ and Acc over 5 independent runs. For Tiny ImageNet, we report the top-1 accuracy (Acc_1) and top-5 accuracy (Acc_5) in Table 3. All of the results are presented in the form of mean \pm standard deviation. Our supplementary materials presents the results of the Wilcoxon signed-rank test at a significance level of 0.05.

Figure 2 presents the experimental results under strongly convex settings. It illustrates that FedBIQ and FedWBIQ significantly outperform other FL algorithms in terms of convergence speed for loss reduction and accuracy improvement, across both I.I.D. and non-I.I.D. data distributions. Specifically, these methods demonstrate accelerated convergence, achieving lower loss and higher accuracy more rapidly than other methods. For non-convex settings, both FedBIQ and FedWBIQ achieve significantly lower loss (0.16) under the I.I.D. distribution on MNIST, with FedWBIQ’s accuracy being only 0.21% lower than that of FedAvg. Under non-I.I.D. distribution, FedWBIQ maintain the lowest loss (0.18) and high accuracy (94.58%), indicating the robustness of the quantizer. On CIFAR-10 and Tiny ImageNet, FedBIQ and FedWBIQ also exhibit markedly higher accuracy compared to other quantization algorithms, further validating their effectiveness in more complex scenarios.

geNet, FedBIQ and FedWBIQ also exhibit markedly higher accuracy compared to other quantization algorithms, further validating their effectiveness in more complex scenarios.

FedSQFL combines quantization with sparsification, substantially reducing training overhead; however, the dual sources of noise introduced by both techniques degrade model performance. FedNQFL introduces additional parameter overhead and computational cost, resulting in approximately 27% higher training overhead compared to FedBIQ and FedWBIQ while achieving comparable model accuracy on MNIST. FedNQFL relies on the assumption that the distribution of local gradient vectors tend to a Gaussian distribution, which is a potential factor contributing to its lower accuracy on Tiny ImageNet.

Although FedBIQ and FedWBIQ incur slightly higher training costs compared to other quantization methods due to the additional communication overhead of transmitting the quantization range, they still achieve effective cost control. The communication costs of FedBIQ and FedWBIQ remain significantly lower than those of unquantized base-

Dataset	Algorithm	I.I.D. Distribution			Non-I.I.D. Distribution		
		$cost_{avg}$	Acc_1 (%)	Acc_5 (%)	$cost_{avg}$	Acc_1 (%)	Acc_5 (%)
Tiny ImageNet	FedAvg	54750.17±1e-1	21.54±5e-3	44.30±7e-3	54750.24±2e-1	21.38±7e-3	43.84±9e-3
	FedRQ	5138.81±2e-1	17.31±1e-3	39.14±4e-3	5138.89±2e-1	17.54±2e-3	39.54±3e-3
	FedSQ	5140.28±1e-1	16.96±4e-3	38.37±4e-3	5140.12±4e-1	17.09±2e-3	38.54±3e-3
	FedPAQ	5140.51±3e-1	18.19±2e-3	40.69±4e-3	5140.24±3e-2	18.45±3e-3	41.06±3e-3
	FedNQFL	5141.03±1e-1	5.92±4e-3	17.38±9e-3	5141.17±2e-1	6.07±2e-3	17.65±5e-3
	FedSQFL	1717.27±6e-2	14.31±1e-3	33.51±2e-3	1717.30±1e-1	14.45±2e-3	33.67±2e-3
	FedBIQ	5139.31±1e-2	20.29±2e-3	43.11±5e-3	5139.51±3e-1	20.38±3e-3	43.30±1e-3
	FedWBIQ	5140.19±1e-2	20.39±2e-3	43.11±4e-3	5140.12±4e-1	20.32±8e-4	43.22±2e-3

Table 3: End-to-end comparison results on Tiny ImageNet.

lines, and the increase over other quantization methods is no more than 0.21%. Remarkably, FedBIQ and FedWBIQ achieve superior accuracy under low-bit settings, rendering the proposed quantizers highly competitive for communication-constrained FL systems.

FedBIQ vs. FedWBIQ. Theoretically, FedWBIQ admits an error upper bound that is twice that of FedBIQ. However, in practice, the two methods exhibit comparable empirical performance. This discrepancy arises from differences in the input data distribution. Our experiments reveal that input values tend to concentrate near the interval boundaries when quantization range is small. In such cases, FedWBIQ’s outputs are biased toward the endpoints, yielding more accurate quantization compared to FedBIQ. Consequently, FedWBIQ demonstrates superior performance when using a larger learning rate and a smaller quantization range.

Related Work

Communication Compression in FL. Addressing the communication bottleneck in FL has become a central focus of recent research. Prior work has explored various techniques to improve communication efficiency in distributed training, including sparsification (Aji and Heafield 2017; Lin et al. 2017), periodic aggregation (McMahan et al. 2017; Sun et al. 2022) and quantization (Alistarh et al. 2017; Yue et al. 2021; Bernstein et al. 2018). While sparsification reduces communication bits by transmitting only significant updates, most deterministic sparsification schemes suffer from inadequate performance guarantees (Sun et al. 2022). Periodic aggregation reduces the frequency of communication but does not decrease the per-round transmission cost. Quantization, which represents model updates using low-precision data, provides a scalable solution by significantly compressing communication bits. This approach has proven particularly effective in resource-constrained environments, such as wireless sensor networks (Msechu and Giannakis 2012).

Quantization. Quantization can be broadly categorized into uniform and non-uniform techniques. Uniform quantization, which includes RQ (Sun et al. 2022; Nagel et al. 2020; Lee, Kim, and Ham 2021; Nagel et al. 2022) and SQ (Elgabli et al. 2020; Reiszadeh et al. 2020; Alistarh et al. 2017; Tran et al. 2019), utilizes evenly spaced intervals, providing simplicity but often struggling to effectively represent non-uniform distributions. In contrast, non-uniform

quantization (Zhang et al. 2018; Jung et al. 2019; Gongyo et al. 2024; Li, Dong, and Wang 2019) dynamically allocates intervals based on distributions, thereby enhancing accuracy. However, its implementation in FL has been limited due to the challenges of aligning quantization schemes across clients without incurring extra communication costs. Recent hybrid approaches (Feng and Venkatasubramaniam 2024) and privacy-focused variants (Feng and Venkatasubramaniam 2025; Youn et al. 2023) have investigated trade-offs between accuracy and privacy.

Theoretical Guarantees. Prior work (Alistarh et al. 2017) analyzed the convergence of unbiased quantizers in distributed learning, focusing on efficient computation on GPUs rather than FL. Existing convergence analyses for quantized FL, such as those by (Elgabli et al. 2020) and (Reiszadeh et al. 2020), assume unbiased quantization errors. The assumption limits their applicability to biased quantizers. Prior work has argued that directly applying biased compression in FL leads to non-convergence (Li and Li 2023). Our work bridges this gap by dynamically controlling the quantization range and establishing the first convergence guarantees for biased quantizers under both strongly convex and non-convex settings. We generalizing prior theoretical framework and demonstrate the feasibility of high-precision, communication-efficient quantization in FL systems.

Conclusion

In our work, we address the communication bottleneck in FL by proposing two efficient non-uniform quantizers: BIQ and WBIQ. BIQ leverages a bisection-based encoding mechanism to achieve higher quantization accuracy without increasing communication costs. WBIQ further enhances robustness by improving the quantization of boundary values through maximum likelihood estimation. Additionally, we provide the first convergence analysis for biased quantizers within the FL framework, thereby extending theoretical guarantees beyond the limitations of unbiased quantization. Our experiments validate that BIQ and WBIQ outperform existing methods under both I.I.D. and non-I.I.D. distributions. The end-to-end training cost of BIQ and WBIQ is comparable to that of other quantization algorithms, while the model accuracy remains close to that of their full-precision counterparts. Our work offers a practical and theoretically grounded solution for resource-constrained FL, ensuring efficient communication and robust performance.

References

- Aji, A. F.; and Heafield, K. 2017. Sparse communication for distributed gradient descent. *arXiv preprint arXiv:1704.05021*.
- Alistarh, D.; Grubic, D.; Li, J.; Tomioka, R.; and Vojnovic, M. 2017. QSGD: Communication-efficient SGD via gradient quantization and encoding. *Advances in neural information processing systems*, 30.
- Bai, Y.; Wang, Y.-X.; and Liberty, E. 2018. Proxquant: Quantized neural networks via proximal operators. *arXiv preprint arXiv:1810.00861*.
- Bernstein, J.; Wang, Y.-X.; Azizadenesheli, K.; and Anandkumar, A. 2018. signSGD: Compressed optimisation for non-convex problems. In *International Conference on Machine Learning*, 560–569. PMLR.
- Boscarino, N.; Cartwright, R. A.; Fox, K.; and Tsosie, K. S. 2022. Federated learning and Indigenous genomic data sovereignty. *Nature machine intelligence*, 4(11): 909–911.
- Chen, G.; Xie, K.; Tu, Y.; Song, T.; Xu, Y.; Hu, J.; and Xin, L. 2023. Nqfl: Nonuniform quantization for communication efficient federated learning. *IEEE Communications Letters*, 28(2): 332–336.
- Cui, L.; Qu, Y.; Xie, G.; Zeng, D.; Li, R.; Shen, S.; and Yu, S. 2021. Security and privacy-enhanced federated learning for anomaly detection in IoT infrastructures. *IEEE Transactions on Industrial Informatics*, 18(5): 3492–3500.
- Elgabli, A.; Park, J.; Bedi, A. S.; Bennis, M.; and Aggarwal, V. 2020. Q-GADMM: Quantized Group ADMM for Communication Efficient Decentralized Machine Learning. In *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 8876–8880.
- Feng, C.; and Venkitasubramaniam, P. 2024. RQP-SGD: Differential Private Machine Learning through Noisy SGD and Randomized Quantization. *arXiv preprint arXiv:2402.06606*.
- Feng, C.; and Venkitasubramaniam, P. 2025. Randomized Quantization for Privacy in Resource Constrained Machine Learning at-the-edge and Federated Learning. *IEEE Transactions on Machine Learning in Communications and Networking*, 1–1.
- Gongyo, S.; Liang, J.; Ambai, M.; Kawakami, R.; and Sato, I. 2024. Learning Non-uniform Step Sizes for Neural Network Quantization. In *Proceedings of the Asian Conference on Computer Vision*, 4385–4402.
- Gupta, S.; Agrawal, A.; Gopalakrishnan, K.; and Narayanan, P. 2015. Deep learning with limited numerical precision. In *International conference on machine learning*, 1737–1746. PMLR.
- Han, S.; Mao, H.; and Dally, W. J. 2015. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. *arXiv preprint arXiv:1510.00149*.
- Jung, S.; Son, C.; Lee, S.; Son, J.; Han, J.-J.; Kwak, Y.; Hwang, S. J.; and Choi, C. 2019. Learning to quantize deep networks by optimizing quantization intervals with task loss. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 4350–4359.
- Krizhevsky, A.; Nair, V.; and Hinton, G. 2009. Learning Multiple Layers of Features from Tiny Images. Technical report, University of Toronto. CIFAR-10 dataset.
- Lecun, Y.; Bottou, L.; Bengio, Y.; and Haffner, P. 1998. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11): 2278–2324.
- Lee, J.; Kim, D.; and Ham, B. 2021. Network quantization with element-wise gradient scaling. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 6448–6457.
- Lee, K.; Lam, M.; Pedarsani, R.; Papailiopoulos, D.; and Ramchandran, K. 2018. Speeding Up Distributed Machine Learning Using Codes. *IEEE Transactions on Information Theory*, 64(3): 1514–1529.
- Li, X.; and Li, P. 2023. Analysis of error feedback in federated non-convex optimization with biased compression: Fast convergence and partial participation. In *International Conference on Machine Learning*, 19638–19688. PMLR.
- Li, Y.; Dong, X.; and Wang, W. 2019. Additive powers-of-two quantization: An efficient non-uniform discretization for neural networks. *arXiv preprint arXiv:1909.13144*.
- Lin, Y.; Han, S.; Mao, H.; Wang, Y.; and Dally, W. J. 2017. Deep gradient compression: Reducing the communication bandwidth for distributed training. *arXiv preprint arXiv:1712.01887*.
- Luqman, A.; Qazi, K.; and Khan, I. 2024. Post-Training Non-Uniform Quantization for Convolutional Neural Networks. *arXiv preprint arXiv:2412.07391*.
- Marnissi, O.; El Hammouti, H.; and Bergou, E. H. 2024. Adaptive sparsification and quantization for enhanced energy efficiency in federated learning. *IEEE Open Journal of the Communications Society*.
- McMahan, B.; Moore, E.; Ramage, D.; Hampson, S.; and y Arcas, B. A. 2017. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, 1273–1282. PMLR.
- Miyashita, D.; Lee, E. H.; and Murmann, B. 2016. Convolutional neural networks using logarithmic data representation. *arXiv preprint arXiv:1603.01025*.
- Msechu, E. J.; and Giannakis, G. B. 2012. Sensor-Centric Data Reduction for Estimation With WSNs via Censoring and Quantization. *IEEE Transactions on Signal Processing*, 60(1): 400–414.
- Nagel, M.; Amjad, R. A.; Van Baalen, M.; Louizos, C.; and Blankevoort, T. 2020. Up or down? adaptive rounding for post-training quantization. In *International conference on machine learning*, 7197–7206. PMLR.
- Nagel, M.; Fournarakis, M.; Bondarenko, Y.; and Blankevoort, T. 2022. Overcoming oscillations in quantization-aware training. In *International Conference on Machine Learning*, 16318–16330. PMLR.
- Pati, S.; Baid, U.; Edwards, B.; Sheller, M.; Wang, S.-H.; Reina, G. A.; Foley, P.; Gruzdev, A.; Karkada, D.; Davatzikos, C.; et al. 2022. Federated learning enables big data

for rare cancer boundary detection. *Nature communications*, 13(1): 7346.

Reisizadeh, A.; Mokhtari, A.; Hassani, H.; Jadbabaie, A.; and Pedarsani, R. 2020. Fedpaq: A communication-efficient federated learning method with periodic averaging and quantization. In *International conference on artificial intelligence and statistics*, 2021–2031. PMLR.

Sun, J.; Chen, T.; Giannakis, G. B.; Yang, Q.; and Yang, Z. 2022. Lazily Aggregated Quantized Gradient Innovation for Communication-Efficient Federated Learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(4): 2031–2044.

Tran, N. H.; Bao, W.; Zomaya, A.; Nguyen, M. N.; and Hong, C. S. 2019. Federated learning over wireless networks: Optimization model design and analysis. In *IEEE INFOCOM 2019-IEEE conference on computer communications*, 1387–1395. IEEE.

Wang, L.; Dong, X.; Wang, Y.; Liu, L.; An, W.; and Guo, Y. 2022. Learnable lookup table for neural network quantization. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 12423–12433.

Youn, Y.; Hu, Z.; Ziani, J.; and Abernethy, J. 2023. Randomized quantization is all you need for differential privacy in federated learning. *arXiv preprint arXiv:2306.11913*.

Yue, L.; Cai, Y.; Zhu, M.; Wang, P.; Zhang, L.; Sun, M.; Liang, S.; Lei, M.; Zhang, J.; Hua, B.; Tian, L.; Zou, Y.; and Li, A. 2021. Improving Performance of Direct-Detection Terahertz Communication System based on k-Means Adaptive Vector Quantization. In *2021 19th International Conference on Optical Communications and Networks (ICOON)*, 1–3.

Zhang, D.; Yang, J.; Ye, D.; and Hua, G. 2018. Lq-nets: Learned quantization for highly accurate and compact deep neural networks. In *Proceedings of the European conference on computer vision (ECCV)*, 365–382.

Appendix

This work addresses the theoretical limitations of existing convergence analyses for quantization algorithms in FL. Building upon the unbiased quantizer framework introduced in (Reisizadeh et al. 2020), we make a key innovation by introducing a bounded bias assumption on the quantization error. Under this assumption, we establish the first convergence rate upper bounds for biased quantizers within the FedAvg framework. This theoretical advance relaxes the strict reliance on unbiasedness in prior work, not only providing a formal justification for the feasibility of biased quantizers such as BIQ, but more importantly, establishing a general analytical paradigm that accommodates diverse quantization strategies. Our analysis fundamentally extends the theoretical foundation of communication-efficient optimization in resource-constrained FL scenarios.

A Code Availability

Our code is available at <https://github.com/GaiLuyang/BIQ>.

B Proof of Theorem 1

In this section, we analyze the convergence of strongly convex assumption.

B.1 Proof of Lemma 1

Lemma 1. *Under Assumption 1 and the update rules specified in Algorithm 2 the gap between the global model θ^{k+1} and the optimal model θ^* at iteration $k + 1$, measured in terms of the expected squared norm $\mathbb{E}\|\theta^{k+1} - \theta^*\|^2$, satisfies the following inequality:*

$$\mathbb{E}\|\theta^{k+1} - \theta^*\|^2 = \underbrace{\mathbb{E}\|\bar{\theta}^{k,\tau} - \theta^*\|^2}_I + \underbrace{\mathbb{E}\|\tilde{\theta}^{k+1} - \bar{\theta}^{k,\tau}\|^2}_{II} + \underbrace{\mathbb{E}\|\theta^{k+1} - \tilde{\theta}^{k+1}\|^2}_{III} + \underbrace{2\mathbb{E}\langle \bar{\eta}^k, \bar{\theta}^{k,\tau} - \theta^* \rangle}_{IV}. \quad (6)$$

where $Q(\cdot)$ is denoted as quantizer. And we give the following definition: $\theta^{k+1} := \theta^k + \frac{1}{s} \sum_{n \in \mathcal{S}_k} Q(\theta_n^{k,\tau} - \theta^k)$, $\bar{\theta}^{k+1} := \theta^k + \frac{1}{N} \sum_{n \in [N]} Q(\theta_n^{k,\tau} - \theta^k)$, $\bar{\theta}^{k,t} := \frac{1}{N} \sum_{n \in [N]} \theta_n^{k,t}$, $\bar{\eta}^k := \frac{1}{N} \sum_{n \in [N]} \eta_n^k$.

Lemma 1 divides $\mathbb{E}\|\theta^{k+1} - \theta^*\|^2$ into I, II, III, IV . In Lemmas 2 to 5, we will analyze the relationship between I, II, III, IV and $\mathbb{E}\|\theta^k - \theta^*\|^2$, and then analyze the upper bound of $\mathbb{E}\|\theta^{k+1} - \theta^*\|^2$.

Proof. Since model parameters are updated based on historical information $\mathcal{F}^{k,t}$ during iteration, the probabilities of random events (such as customer sampling) involved in each update step are conditional probabilities. To simplify the presentation, we have omitted historical information.

$$\begin{aligned} \mathbb{E}_{\mathcal{S}_k} \theta^{k+1} &= \mathbb{E}[\theta^k + \frac{1}{s} \sum_{n \in \mathcal{S}_k} Q(\theta_n^{k,\tau} - \theta^k)] \\ &= \theta^k + \frac{1}{s} \sum_{\mathcal{A} \in [N], |\mathcal{A}|=s} \Pr[\mathcal{S}_k = \mathcal{A}] \sum_{n \in \mathcal{S}_k} Q(\theta_n^{k,\tau} - \theta^k). \end{aligned} \quad (7)$$

$\sum_{\mathcal{A} \in [N], |\mathcal{A}|=s} \sum_{n \in \mathcal{S}_k} Q(\theta_n^{k,\tau} - \theta^k)$ represents the sum of $Q(\theta_n^{k,\tau} - \theta^k)$ in all \mathcal{S}_k combinations. Under the condition that the n -th client is selected, the remaining $N - 1$ clients select $s - 1$ to form \mathcal{S}_k . Thus the number of combinations containing the n -th client is $\binom{N-1}{s-1}$. So $Q(\theta_n^{k,\tau} - \theta^k)$ computes $\binom{N-1}{s-1}$ times.

$$\begin{aligned} \mathbb{E}_{\mathcal{S}_k} \theta^{k+1} &= \theta^k + \frac{s!(N-s)!}{N!} \cdot \frac{1}{s} \cdot \frac{(N-1)!}{(s-1)!(N-s)!} \sum_{n \in [N]} Q(\theta_n^{k,\tau} - \theta^k) \\ &= \theta^k + \frac{1}{N} \sum_{n \in [N]} Q(\theta_n^{k,\tau} - \theta^k) \\ &= \bar{\theta}^{k+1}. \end{aligned} \quad (8)$$

490 According to Assumption 1, it holds that

$$\begin{aligned}
\mathbb{E}_Q \tilde{\theta}^{k+1} &= \theta^k + \frac{1}{N} \sum_{n \in [N]} \mathbb{E}_Q Q(\theta_n^{k,\tau} - \theta^k) \\
&= \theta^k + \frac{1}{N} \sum_{n \in [N]} (\theta_n^{k,\tau} - \theta^k + \eta_n^k) \\
&= \frac{1}{N} \sum_{n \in [N]} \theta_n^{k,\tau} + \frac{1}{N} \sum_{n \in [N]} \eta_n^k \\
&= \bar{\theta}^{k,\tau} + \bar{\eta}^k.
\end{aligned} \tag{9}$$

491 Since sampling randomness and quantization randomness are independent, we obtain

$$\begin{aligned}
\mathbb{E} \|\theta^{k+1} - \theta^*\|^2 &= \mathbb{E} \|\theta^{k+1} - \tilde{\theta}^{k+1} + \tilde{\theta}^{k+1} - \bar{\theta}^{k,\tau} + \bar{\theta}^{k,\tau} - \theta^*\|^2 \\
&= \mathbb{E} \|\theta^{k+1} - \tilde{\theta}^{k+1}\|^2 + \mathbb{E} \|\tilde{\theta}^{k+1} - \bar{\theta}^{k,\tau}\|^2 + \mathbb{E} \|\bar{\theta}^{k,\tau} - \theta^*\|^2 \\
&\quad + \mathbb{E}[2\langle \theta^{k+1} - \tilde{\theta}^{k+1}, \tilde{\theta}^{k+1} - \bar{\theta}^{k,\tau} \rangle] + 2\langle \tilde{\theta}^{k+1} - \bar{\theta}^{k,\tau}, \bar{\theta}^{k,\tau} - \theta^* \rangle \\
&= \mathbb{E} \|\theta^{k+1} - \tilde{\theta}^{k+1}\|^2 + \mathbb{E} \|\tilde{\theta}^{k+1} - \bar{\theta}^{k,\tau}\|^2 + \mathbb{E} \|\bar{\theta}^{k,\tau} - \theta^*\|^2 + 2\mathbb{E} \langle \bar{\eta}^k, \bar{\theta}^{k,\tau} - \theta^* \rangle.
\end{aligned} \tag{10}$$

492 \square

493 B.2 Proof of Lemma 2

494 **Lemma 2.** Under Assumptions 2 to 4 and the update method of Algorithm 2, select the learning rate
495 $\alpha_k = \min\{\frac{\mu}{L^2}, \frac{1}{L\sqrt{\tau(\tau-1)}}\}$, then I and $\mathbb{E} \|\theta^k - \theta^*\|^2$ satisfies the following inequality:

$$I \leq (1 + N\alpha_k^2)(1 - \mu\alpha_k)^\tau \mathbb{E} \|\theta^k - \theta^*\|^2 + \frac{(\tau-1)\sigma^2 e}{N} + \frac{\tau^2 \sigma^2 \alpha_k^2}{N} + \frac{\tau \sigma^2 \alpha_k^2 e}{N}. \tag{11}$$

496 Compared with Lemma 2 in (Reisizadeh et al. 2020), we obtain a tighter upper bound and relax the
497 condition of learning rate, which is beneficial to accelerating convergence.

498 *Proof.* We define the following sequence:

$$\begin{aligned}
\theta_n^{k,t+1} &= \theta_n^{k,t} - \alpha_k \hat{\nabla} f_i(\theta_n^{k,t}), k = 0, 1, \dots, K-1, t = 0, 1, \dots, \tau-1 \\
\gamma^{k,t+1} &:= \gamma^{k,t} - \alpha_k \nabla f(\gamma^{k,t}), k = 0, 1, \dots, K-1, t = 0, 1, \dots, \tau-1 \\
\Gamma^{k,t} &:= \frac{1}{N} \sum_{n \in [N]} \hat{\nabla} f_n(\theta_n^{k,t}) - \nabla f(\gamma^{k,t}), k = 0, 1, \dots, K-1, t = 0, 1, \dots, \tau-1 \\
w^{k,t} &:= \frac{1}{N} \sum_{n \in [N]} \mathbb{E} \|\theta_n^{k,t} - \gamma^{k,t}\|^2, k = 0, 1, \dots, K-1, t = 0, 1, \dots, \tau-1
\end{aligned}$$

502 where $\hat{\nabla} f_i(\theta_n^{k,t}) = \hat{\nabla} f_i(\theta_n^{k,t}, \hat{D}_n)$, $\theta_n^{k,0} = \gamma^{k,0} = \theta^k$. $\{\theta_n^{k,t}\}$ is the local model update pro-
503 cess. $\{\gamma^{k,t}\}$ represents the sequence of gradient descent update processes in centralized training.
504 $\{\Gamma^{k,t}\}, \{w^{k,t}\}$ are defined as the gradient error and the mean model error of two different update
505 methods, respectively. In addition, let us denote $\Gamma^k = \sum_{t=0}^{\tau-1} \Gamma^{k,t}$. From (36)(37)(38) in (Reisizadeh
506 et al. 2020), under the setting of $\alpha_k = \frac{\mu}{L^2}$, we yield

$$I \leq (1 + N\alpha_k^2)(1 - \mu\alpha_k)^\tau \mathbb{E} \|\theta^k - \theta^*\|^2 + \frac{1}{N} \|\mathbb{E} \Gamma^k\|^2 + \alpha_k^2 \mathbb{E} \|\Gamma^k\|^2. \tag{12}$$

507 Next, we analysis $\|\mathbb{E} \Gamma^k\|^2$. From the unbiasedness of random gradients in Assumption 3, it holds
508 that: $\mathbb{E} \Gamma^{k,0} = 0$. Using Jensen's inequality and (40) in (Reisizadeh et al. 2020), we obtain

$$\begin{aligned}
\|\mathbb{E} \Gamma^k\|^2 &\leq (\tau-1) \sum_{t=1}^{\tau-1} \|\mathbb{E} \Gamma^{k,t}\|^2 \\
&\leq (\tau-1) L^2 \sum_{t=1}^{\tau-1} w^{k,t}.
\end{aligned} \tag{13}$$

By analyzing (63) in (Reisizadeh et al. 2020), we obtain

$$\begin{aligned}\tau L^2 \alpha_k^2 \sum_{t=1}^{\tau-1} w^{k,t} &\leq \tau \alpha_k^2 \sigma^2 + \tau L^2 \alpha_k^2 \sum_{t=1}^{\tau-1} w^{k_1 t} \\ &\leq \tau \alpha_k^2 \sigma^2 (1 + \tau L^2 \alpha_k^2)^{\tau-1}.\end{aligned}$$

Therefore, we have

$$L^2 \sum_{t=1}^{\tau-1} w^{k,t} \leq \sigma^2 (1 + \tau L^2 \alpha_k^2)^{\tau-1}. \quad (14)$$

Plugging Equation (14) in Equation (13) implies that

$$\begin{aligned}\|\mathbb{E}\Gamma^k\|^2 &\leq (\tau-1)\sigma^2(1 + \tau L^2 \alpha_k^2)^{\tau-1} \\ &\leq (\tau-1)\sigma^2 e^{\tau(\tau-1)L^2 \alpha_k^2}.\end{aligned} \quad (15)$$

Let $\alpha_k = \min\{\frac{\mu}{L^2}, \frac{1}{L\sqrt{\tau(\tau-1)}}\}$, which means $\tau(\tau-1)L^2 \alpha_k^2 \leq 1$. It holds that

$$\|\mathbb{E}\Gamma^k\|^2 \leq (\tau-1)\sigma^2 e. \quad (16)$$

Finally, we analysis $\mathbb{E}\|\Gamma^k\|^2$. According to Jensen inequality, we have

$$\mathbb{E}\|\Gamma^k\|^2 \leq \tau \sum_{t=0}^{\tau-1} \mathbb{E}\|\Gamma^{k,t}\|^2. \quad (17)$$

When $t = 0$, we obtain the following inequality from (47) in (Reisizadeh et al. 2020):

$$\mathbb{E}\|\Gamma^k\|^2 \leq \frac{\sigma^2}{N}. \quad (18)$$

When $t \geq 1$, we yield

$$\begin{aligned}\mathbb{E}\|\Gamma^{k,t}\|^2 &= \mathbb{E}\left\|\frac{1}{N} \sum_{n \in [N]} \hat{\nabla} f_n(\theta_n^{k,t}) - \nabla f(\gamma^{k,t})\right\|^2 \\ &= \mathbb{E}\left\|\frac{1}{N} \sum_{n \in [N]} [\hat{\nabla} f_n(\theta_n^{k,t}) - \nabla f(\theta_n^{k,t})] + \frac{1}{N} \sum_{n \in [N]} [\nabla f(\theta_n^{k,t}) - \nabla f(\gamma^{k,t})]\right\|^2 \\ &= \mathbb{E}\left\|\frac{1}{N} \sum_{n \in [N]} [\hat{\nabla} f_n(\theta_n^{k,t}) - \nabla f(\theta_n^{k,t})]\right\|^2 + \mathbb{E}\left\|\frac{1}{N} \sum_{n \in [N]} [\nabla f(\theta_n^{k,t}) - \nabla f(\gamma^{k,t})]\right\|^2 \\ &\quad + 2\mathbb{E}\left\langle \frac{1}{N} \sum_{n \in [N]} [\hat{\nabla} f_n(\theta_n^{k,t}) - \nabla f(\theta_n^{k,t})], \frac{1}{N} \sum_{n \in [N]} [\nabla f(\theta_n^{k,t}) - \nabla f(\gamma^{k,t})] \right\rangle \\ &\leq \frac{1}{N^2} \mathbb{E}\left\|\sum_{n \in [N]} [\hat{\nabla} f_n(\theta_n^{k,t}) - \nabla f(\theta_n^{k,t})]\right\|^2 + \frac{1}{N^2} \mathbb{E}\left\|\sum_{n \in [N]} \nabla f(\theta_n^{k,t}) - \nabla f(\gamma^{k,t})\right\|^2 \\ &\quad + 2\mathbb{E}\left\langle \frac{1}{N} \sum_{n \in [N]} [\hat{\nabla} f_n(\theta_n^{k,t}) - \nabla f(\theta_n^{k,t})], \frac{1}{N} \sum_{n \in [N]} [\nabla f(\theta_n^{k,t}) - \nabla f(\gamma^{k,t})] \right\rangle.\end{aligned} \quad (19)$$

According to Assumption 3, we have

$$\frac{1}{N^2} \mathbb{E}\left\|\sum_{n \in [N]} [\hat{\nabla} f_n(\theta_n^{k,t}) - \nabla f(\theta_n^{k,t})]\right\|^2 = \frac{\sigma^2}{N}. \quad (20)$$

517

$$2\mathbb{E}\left\langle \frac{1}{N} \sum_{n \in [N]} [\hat{\nabla} f_n(\theta_n^{k,t}) - \nabla f(\theta_n^{k,t})], \frac{1}{N} \sum_{n \in [N]} [\nabla f(\theta_n^{k,t}) - \nabla f(\gamma^{k,t})] \right\rangle = 0. \quad (21)$$

From Assumption 2, we get

$$\frac{1}{N^2} \mathbb{E}\left\|\sum_{n \in [N]} [\nabla f(\theta_n^{k,t}) - \nabla f(\gamma^{k,t})]\right\|^2 = \frac{L^2}{N^2} \sum_{n \in [N]} \mathbb{E}\|\theta_n^{k,t} - \gamma^{k,t}\|^2. \quad (22)$$

519 Substituting Equations (20) to (22) into Equation (19), we get

$$\begin{aligned}\mathbb{E}\|\Gamma^{k,t}\|^2 &\leq \frac{\sigma^2}{N} + \frac{L^2}{N^2} \sum_{n \in [N]} \mathbb{E}\|\theta_n^{k,t} - \gamma^{k,t}\|^2 \\ &= \frac{\sigma^2}{N} + \frac{L^2}{N} w^{k,t}.\end{aligned}\tag{23}$$

520 Plugging Equations (18) and (23) in Equation (17) implies that

$$\begin{aligned}\mathbb{E}\|\Gamma^k\|^2 &\leq \tau \sum_{t=0}^{\tau-1} \mathbb{E}\|\Gamma^{k,t}\|^2 \\ &\leq \tau \frac{\sigma^2}{N} + \tau \sum_{t=1}^{\tau-1} \mathbb{E}\|\Gamma^{k,t}\|^2 \\ &\leq \tau \frac{\sigma^2}{N} + \tau \sum_{t=1}^{\tau-1} \left(\frac{\sigma^2}{N} + \frac{L^2}{N} w^{k,t} \right) \\ &\leq \tau \frac{\sigma^2}{N} + \tau(\tau-1) \frac{\sigma^2}{N} + \frac{\tau}{N} L^2 \sum_{t=1}^{\tau-1} w^{k,t} \\ &= \tau^2 \frac{\sigma^2}{N} + \frac{\tau}{N} L^2 \sum_{t=1}^{\tau-1} w^{k,t}.\end{aligned}\tag{24}$$

521 Finally, substituting Equation (14) into Equation (24), we get the upper bound of $\mathbb{E}\|\Gamma^k\|^2$:

$$\begin{aligned}\mathbb{E}\|\Gamma^k\|^2 &\leq \tau^2 \frac{\sigma^2}{N} + \frac{\tau}{N} \sigma^2 (1 + \tau L^2 \alpha_k^2)^{\tau-1} \\ &\leq \tau^2 \frac{\sigma^2}{N} + \frac{\tau}{N} \sigma^2 e^{\tau(\tau-1)L^2 \alpha_k^2} \\ &\leq \tau^2 \frac{\sigma^2}{N} + \frac{\tau}{N} \sigma^2 e.\end{aligned}\tag{25}$$

522 Substituting Equations (16) and (25) into Equation (12), we get the upper bound of I :

$$\begin{aligned}I &\leq (1 + N\alpha_k^2)(1 - \mu\alpha_k)^T \mathbb{E}\|\theta^k - \theta^*\|^2 + \frac{1}{N} \|\mathbb{E}\Gamma^k\|^2 + \alpha_k^2 \mathbb{E}\|\Gamma^k\|^2 \\ &\leq (1 + N\alpha_k^2)(1 - \mu\alpha_k)^\tau \mathbb{E}\|\theta^k - \theta^*\|^2 + \frac{(\tau-1)\sigma^2 e}{N} + \frac{\tau^2 \sigma^2 \alpha_k^2}{N} + \frac{\tau \sigma^2 \alpha_k^2 e}{N}.\end{aligned}\tag{26}$$

523 □

524 B.3 Proof of Lemma 3

525 **Lemma 3.** Under Assumptions 1 to 4 and the update method of Algorithm 2 selecting the learning
526 rate $\alpha_k = \min\{\frac{\mu}{L^2}, \frac{1}{L\sqrt{\tau(\tau-1)}}, \frac{1}{\mu}\}$, then II and $\mathbb{E}\|\theta^k - \theta^*\|^2$ satisfy the following inequality:

$$II \leq 4q \frac{\alpha_k \tau L^2}{\mu} \mathbb{E}\|\theta^k - \theta^*\|^2 + 4q \alpha_k^2 \tau \sigma^2 (\tau + e) + 2\mathcal{G}^k,\tag{27}$$

527 where $\|\eta_n^k\| \leq G_n^k$, $\mathcal{G}^k = \frac{1}{N} \sum_{n \in [N]} G_n^k{}^2$.

528 *Proof.* According to the definition of $\tilde{\theta}^{k+1}$ and $\tilde{\theta}^{k,\tau}$, we get

$$\begin{aligned}II &= \mathbb{E} \|\tilde{\theta}^{k+1} - \tilde{\theta}^{k,\tau}\|^2 \\ &= \mathbb{E} \left\| \theta^k + \frac{1}{N} \sum_{n \in [N]} Q(\theta_n^{k,\tau} - \theta^k) - \frac{1}{N} \sum_{n \in [N]} \theta_n^{k,\tau} \right\|^2 \\ &= \mathbb{E} \left\| \frac{1}{N} \sum_{n \in [N]} [Q(\theta_n^{k,\tau} - \theta^k) - (\theta_n^{k,\tau} - \theta^k)] \right\|^2.\end{aligned}\tag{28}$$

529 According to Jensen's inequality, we obtain

$$\begin{aligned}
II &\leq \frac{1}{N} \sum_{n \in [N]} \mathbb{E} \| Q(\theta_n^{k,r} - \theta^k) - (\theta_n^{k,r} - \theta^k) \|^2 \\
&= \frac{1}{N} \sum_{n \in [N]} \mathbb{E} \| Q(\theta_n^{k,r} - \theta^k) - (\theta_n^{k,r} + \eta_n^k - \theta^k) - \eta_n^k \|^2 \\
&= \frac{1}{N} \sum_{n \in [N]} [\mathbb{E} \| Q(\theta_n^{k,r} - \theta^k) - (\theta_n^{k,r} + \eta_n^k - \theta^k) \|^2 + \mathbb{E} \|\eta_n^k\|^2 \\
&\quad - 2\langle Q(\theta_n^{k,r} - \theta^k) - (\theta_n^{k,r} + \eta_n^k - \theta^k), \eta_n^k \rangle] \\
&\leq \frac{1}{N} \sum_{n \in [N]} [2\mathbb{E} \| Q(\theta_n^{k,r} - \theta^k) - (\theta_n^{k,r} + \eta_n^k - \theta^k) \|^2 + 2\mathbb{E} \|\eta_n^k\|^2].
\end{aligned} \tag{29}$$

530 From Assumption [1](#) we get

$$\begin{aligned}
\mathbb{E} \| Q(\theta_n^{k,r} - \theta^k) - (\theta_n^{k,r} + \eta_n^k - \theta^k) \|^2 &\leq q \mathbb{E} \|\theta_n^{k,r} - \theta^k\|^2. \\
\mathbb{E} \|\eta_n^k\|^2 &\leq G_n^{k^2}.
\end{aligned} \tag{30}$$

531 Next, we analysis $\mathbb{E} \|\theta_n^{k,\tau} - \theta^k\|^2$. We denote $h^k := \sum_{t=0}^{\tau-1} \nabla f(\gamma^{k,t})$, $\Gamma_n^k := \sum_{t=0}^{\tau-1} \hat{\nabla} f_n(\theta_n^{k,t}) - \nabla f(\gamma^{k,t})$. It holds that

$$\begin{aligned}
\theta_n^{k,\tau} &= \theta^k - \alpha_k \sum_{t=0}^{\tau-1} \hat{\nabla} f_n(\theta_n^{k,t}) \\
&= \theta^k - \alpha_k \left\{ \sum_{t=0}^{\tau-1} [\hat{\nabla} f_n(\theta_n^{k,t}) - \nabla f_n(\gamma^{k,t})] + \sum_{t=0}^{\tau-1} \nabla f_n(\gamma^{k,t}) \right\} \\
&= \theta^k - \alpha_k (\Gamma_n^k + h^k).
\end{aligned} \tag{31}$$

533 So

$$\begin{aligned}
\mathbb{E} \|\theta_n^{k,\tau} - \theta^k\|^2 &= \alpha_k^2 \mathbb{E} \|\Gamma_n^k + h^k\|^2 \\
&\leq \alpha_k^2 \mathbb{E} \|\Gamma_n^k\|^2 + \alpha_k^2 \mathbb{E} \|h^k\|^2 + 2\alpha_k^2 \langle \Gamma_n^k, h^k \rangle \\
&\leq 2\alpha_k^2 \mathbb{E} \|\Gamma_n^k\|^2 + 2\alpha_k^2 \mathbb{E} \|h^k\|^2.
\end{aligned} \tag{32}$$

534 Set learning rate $\alpha_k = \min\{\frac{\mu}{L^2}, \frac{1}{L\sqrt{\tau(\tau-1)}}, \frac{1}{\mu}\}$. We use Jensen inequality to deduce the following

535 result

$$\begin{aligned}
\|h^k\|^2 &\leq \tau \sum_{t=0}^{\tau-1} \|\nabla f(\gamma^{k,t})\|^2 \\
&\leq \tau L^2 \sum_{t=0}^{\tau-1} \|\nabla f(\gamma^{k,t}) - \nabla f(\theta^*)\|^2 \\
&\leq \tau L^2 \sum_{t=0}^{\tau-1} \|\gamma^{k,t} - \theta^*\|^2 \\
&\leq \tau L^2 \|\theta^k - \theta^*\|^2 \sum_{t=0}^{\tau-1} (1 - \mu\alpha_k)^t \\
&\leq \tau L^2 \|\theta^k - \theta^*\|^2 \frac{1 - (1 - \mu\alpha_k)^\tau}{\mu\alpha_k} \\
&\leq \frac{1}{\mu\alpha_k} \tau L^2 \|\theta^k - \theta^*\|^2.
\end{aligned} \tag{33}$$

536 What's more, substituting $N = 1$ into Equation [\(25\)](#), the following holds

$$\mathbb{E} \|\Gamma_n^k\|^2 \leq \tau^2 \sigma^2 + \tau \sigma^2 e. \tag{34}$$

537 We imply that

$$\begin{aligned}
\mathbb{E}\|\theta_n^{k,\tau} - \theta^k\|^2 &\leq 2\alpha_k^2 \mathbb{E}\|\Gamma_n^k\|^2 + 2\alpha_k^2 \mathbb{E}\|h^k\|^2 \\
&\leq 2\alpha_k^2 \tau^2 \sigma^2 + 2\alpha_k^2 \tau \sigma^2 e + 2\alpha_k^2 \frac{1}{\mu \alpha_k} \tau L^2 \|\theta^k - \theta^*\|^2 \\
&= 2\frac{\alpha_k \tau L^2}{\mu} \mathbb{E}\|\theta^k - \theta^*\|^2 + 2\alpha_k^2 \tau \sigma^2 (\tau + e).
\end{aligned} \tag{35}$$

538 Substituting Equation (35) into Equation (29) yields:

$$\begin{aligned}
II &\leq \frac{1}{N} \sum_{n \in [N]} [2\mathbb{E}\|Q(\theta_n^{k,\tau} - \theta^k) - (\theta_n^{k,\tau} + \eta_n^k - \theta^k)\|^2 + 2\mathbb{E}\|\eta_n^k\|^2] \\
&\leq \frac{1}{N} \sum_{n \in [N]} [2q\mathbb{E}\|\theta_n^{k,\tau} - \theta^k\|^2 + 2G_n^k] \\
&\leq 4q\frac{\alpha_k \tau L^2}{\mu} \mathbb{E}\|\theta^k - \theta^*\|^2 + 4q\alpha_k^2 \tau \sigma^2 (\tau + e) + 2G^k.
\end{aligned} \tag{36}$$

539

□

540 B.4 Proof of Lemma 4

541 **Lemma 4.** Under Assumptions 1 to 4 and the update method of Algorithm 2 select the learning rate
542 $\alpha_k = \min\{\frac{\mu}{L^2}, \frac{1}{L\sqrt{\tau(\tau-1)}}, \frac{1}{\mu}\}$, then III and $\mathbb{E}\|\theta^k - \theta^*\|^2$ satisfies the following inequality:

$$III \leq \frac{4(N-s)}{s(N-1)} \left[\frac{2(q+4)\alpha_k \tau L^2}{\mu} \mathbb{E}\|\theta^k - \theta^*\|^2 + 2(q+4)\alpha_k^2 \sigma^2 \tau (\tau + e) + 3G^k \right]. \tag{37}$$

543 *Proof.* On the one hand, we get the following result from $\mathbb{E}\|Q(x) - x - \eta\|^2 \leq q\|x\|^2$ in Assump-
544 tion 1

$$\begin{aligned}
\mathbb{E}\|Q(x) - x - \eta\|^2 &= \mathbb{E}[\|Q(x)\|^2 + \|x\|^2 + \|\eta\|^2 \\
&\quad - 2\langle Q(x), x \rangle - 2\langle Q(x), \eta \rangle - 2\langle x, \eta \rangle] \\
&= \mathbb{E}\|Q(x)\|^2 + \|x\|^2 + \|\eta\|^2 \\
&\quad - 2\langle x + \eta, x \rangle - 2\langle x + \eta, \eta \rangle - 2\langle x, \eta \rangle \\
&= \mathbb{E}\|Q(x)\|^2 + \|x\|^2 + \|\eta\|^2 - 2\|x\|^2 - 2\|\eta\|^2 - 6\langle \eta, x \rangle \\
&= \mathbb{E}\|Q(x)\|^2 - \|x\|^2 - \|\eta\|^2 - 6\langle \eta, x \rangle \\
&\leq q\|x\|^2.
\end{aligned} \tag{38}$$

545 So we have

$$\begin{aligned}
\mathbb{E}\|Q(x)\|^2 &\leq (q+1)\|x\|^2 + \|\eta\|^2 + 6\langle x, \eta \rangle \\
&= (q+4)\|x\|^2 + 3\|\eta\|^2.
\end{aligned} \tag{39}$$

546 On the other hand, from (59) in (Reisizadeh et al. 2020), we get:

$$\begin{aligned}
III &\leq \frac{N-s}{sN(N-1)} \sum_{n \in [N]} \mathbb{E}_Q \|Q(\theta_n^{k,\tau} - \theta^k) - \frac{1}{N} \sum_{n \in [N]} Q(\theta_n^{k,\tau} - \theta^k)\|^2 \\
&\leq \frac{N-s}{sN(N-1)} \sum_{n \in [N]} [\mathbb{E}_Q \|Q(\theta_n^{k,\tau} - \theta^k)\|^2 + \|\frac{1}{N} \sum_{n \in [N]} Q(\theta_n^{k,\tau} - \theta^k)\|^2 \\
&\quad - 2\langle Q(\theta_n^{k,\tau} - \theta^k), \frac{1}{N} \sum_{n \in [N]} Q(\theta_n^{k,\tau} - \theta^k) \rangle] \\
&\leq \frac{N-s}{sN(N-1)} \sum_{n \in [N]} [2\mathbb{E}_Q \|Q(\theta_n^{k,\tau} - \theta^k)\|^2 + 2N\|\frac{1}{N} \sum_{n \in [N]} Q(\theta_n^{k,\tau} - \theta^k)\|^2] \\
&\leq \frac{N-s}{sN(N-1)} \sum_{n \in [N]} [2\mathbb{E}_Q \|Q(\theta_n^{k,\tau} - \theta^k)\|^2 + \frac{2}{N} \sum_{n \in [N]} \|Q(\theta_n^{k,\tau} - \theta^k)\|^2].
\end{aligned} \tag{40}$$

547 According to $N \geq 1$, we yield

$$III \leq \frac{N-s}{sN(N-1)} \sum_{n \in [N]} [4\mathbb{E}_Q \|Q(\theta_n^{k,\tau} - \theta^k)\|^2]. \quad (41)$$

548 By substituting Equation (39) into Equation (41), it holds that

$$III \leq \frac{4(N-s)}{sN(N-1)} \sum_{n \in [N]} [(q+4)\|\theta_n^{k,\tau} - \theta^k\|^2 + 3\|\eta_n^k\|^2]. \quad (42)$$

549 Substituting Equation (35) into Equation (42), we get the upper bound of III :

$$\begin{aligned} III &\leq \frac{4(N-s)}{sN(N-1)} \sum_{n \in [N]} [(q+4) \cdot \frac{2\alpha_k \tau L^2}{\mu} \mathbb{E}\|\theta^k - \theta^*\|^2 + 2(q+4)\alpha_k^2 \tau \sigma^2 (\tau + e) + 3G_n^{k^2}] \\ &\leq \frac{4(N-s)}{s(N-1)} \left[\frac{2(q+4)\alpha_k \tau L^2}{\mu} \mathbb{E}\|\theta^k - \theta^*\|^2 + 2(q+4)\alpha_k^2 \sigma^2 \tau (\tau + e) + 3\mathcal{G}^k \right]. \end{aligned} \quad (43)$$

550 \square

551 B.5 Proof of Lemma 5

552 **Lemma 5.** Under Assumptions 1 to 4 and the update method of Algorithm 2, select the learning rate
553 $\alpha_k = \min\{\frac{\mu}{L^2}, \frac{1}{L\sqrt{\tau(\tau-1)}}\}$, then IV and $\mathbb{E}\|\theta^k - \theta^*\|^2$ satisfies the following inequality:

$$IV \leq (1 + N\alpha_k^2)(1 - \mu\alpha_k)^\tau \mathbb{E}\|\theta^k - \theta^*\|^2 + \frac{(\tau-1)\sigma^2 e}{N} + \frac{\tau^2 \sigma^2 \alpha_k^2}{N} + \frac{\tau \sigma^2 \alpha_k^2 e}{N} + \mathcal{G}^k. \quad (44)$$

Proof.

$$\begin{aligned} IV &= 2\mathbb{E}\langle \bar{\eta}^k, \bar{\theta}^{k,\tau} - \theta^* \rangle \\ &\leq \mathbb{E}\|\bar{\eta}^k\|^2 + \mathbb{E}\|\bar{\theta}^{k,\tau} - \theta^*\|^2. \end{aligned} \quad (45)$$

554 According to Jensen's inequality

$$\begin{aligned} \mathbb{E}\|\bar{\eta}^k\|^2 &= \mathbb{E}\left\| \frac{1}{N} \sum_{n \in [N]} \eta_n^k \right\|^2 \\ &\leq \frac{1}{N} \sum_{n \in [N]} \mathbb{E}\|\eta_n^k\|^2. \end{aligned} \quad (46)$$

555 From noise bounded in Assumption 1, we obtain

$$\begin{aligned} \mathbb{E}\|\bar{\eta}^k\|^2 &\leq \frac{1}{N} \sum_{n \in [N]} G_n^{k^2} \\ &= \mathcal{G}^{k^2}. \end{aligned} \quad (47)$$

556 Substituting Equation (47) and Equation (11) into Equation (45), we get the upper bound of IV :

$$IV \leq (1 + N\alpha_k^2)(1 - \mu\alpha_k)^\tau \mathbb{E}\|\theta^k - \theta^*\|^2 + \frac{(\tau-1)\sigma^2 e}{N} + \frac{\tau^2 \sigma^2 \alpha_k^2}{N} + \frac{\tau \sigma^2 \alpha_k^2 e}{N} + \mathcal{G}^k. \quad (48)$$

557 \square

558 B.6 Proof of Theorem 1

559 In this section we prove the **Theorem 1** of our work. Here, we first restate Theorem 1.

560 **Theorem B.1** (Convergence of Strongly Convex Settings). Suppose Assumptions 1 to 4
561 hold. Define $A_0 := \frac{2qL}{\mu} + \frac{4L^2(N-s)(q+4)}{s(N-1)\mu}$. Set a constant k_0 that satisfies $k_0 \geq$

562 $\max\{\frac{48\tau-1}{\tau}, \frac{48N-\mu^2\tau}{\mu^2\tau^2}, \frac{48L^2A_0\tau-\mu^3\tau+6L^2\mu}{\mu^2\tau^2-L^2\mu\tau}, \frac{4L^2-\mu^2}{\mu^2\tau}, \frac{4L\sqrt{\tau(\tau-1)}-\mu}{\mu\tau}\}$. Then for any $k \geq k_0$, if the
 563 learning rate satisfies $\alpha_k = \min\{\frac{1}{L^2}, \frac{1}{L\sqrt{\tau(\tau-1)}}, \frac{1}{\mu\tau}\}$, Algorithm 2 has

$$\mathbb{E}\|\theta^k - \theta^*\|^2 \leq (\frac{k_0\tau+1}{k\tau+1})^2 \mathbb{E}\|\theta^{k_0} - \theta^*\|^2 + A_2 \frac{1}{\tau} \cdot \frac{\tau}{k\tau+1} + A_3 \cdot \frac{\tau}{k\tau+1} + A_4 \frac{g^{k^2}}{\tau^2} \cdot \frac{\tau}{k\tau+1}. \quad (49)$$

564 where $\{g^k\}$ is a bounded positive sequence. The constants in **Theorem 1** are defined as

$$A_2 := \frac{16}{\mu^2} \left[2\frac{e\sigma^2 L^2}{N} + 2\frac{e\sigma^2}{N} + 4q\sigma^2 e + \frac{8e(N-s)(q+4)\sigma^2}{s(N-1)} \right]$$

$$A_3 := \frac{16}{\mu^2} \left[2\frac{\sigma^2}{N} + 4q\sigma^2 + \frac{4(N-s)2(q+4)\sigma^2}{s(N-1)} \right], A_4 := \frac{3sN+12N-15s}{s(N-1)}.$$

566 *Proof.* Substituting Equations (11), (27), (37) and (44) into Equation (6), we get

$$\begin{aligned} \mathbb{E}\|\theta^{k+1} - \theta^*\|^2 &= I + II + III + IV \\ &\leq [2(1+N\alpha_k^2)(1-\mu\alpha_k)^\tau + 4q\frac{\alpha_k\tau L^2}{\mu} + \frac{4(N-s)}{s(N-1)} \cdot \frac{2(q+4)\alpha_k\tau L^2}{\mu}] \\ &\quad \cdot \mathbb{E}\|\theta^k - \theta^*\|^2 + 2\frac{(\tau-1)\sigma^2 e}{N} + 2\frac{\tau^2\sigma^2 a_k^2}{N} + 2\frac{\tau\sigma^2 a_k^2 e}{N} + 4q\alpha_k^2\tau\sigma^2(\tau+e) \\ &\quad + \frac{4(N-s)2(q+4)}{s(N-1)}\alpha_k^2\sigma^2\tau(\tau+e) + \frac{3sN+12N-15s}{s(N-1)}\mathcal{G}^k \\ &\leq [2(1+N\alpha_k^2)(1-\mu\alpha_k)^\tau + 4q\frac{\alpha_k\tau L^2}{\mu} + \frac{4(N-s)}{s(N-1)} \cdot \frac{2(q+4)\alpha_k\tau L^2}{\mu}] \\ &\quad \cdot \mathbb{E}\|\theta^k - \theta^*\|^2 + [2\frac{e\sigma^2 L^2}{N} + 2\frac{e\sigma^2}{N} + 4q\sigma^2 e + \frac{8e(N-s)(q+4)\sigma^2}{s(N-1)}]\tau\alpha_k^2 \\ &\quad + [2\frac{\sigma^2}{N} + 4q\sigma^2 + \frac{4(N-s)2(q+4)\sigma^2}{s(N-1)}]\tau^2 a_k^2 + \frac{3sN+12N-15s}{s(N-1)}\mathcal{G}^k. \end{aligned} \quad (50)$$

567 We define the following constants

$$A_1 := 2(1+N\alpha_k^2)(1-\mu\alpha_k)^\tau + 4q\frac{\alpha_k\tau L^2}{\mu} + \frac{4(N-s)}{s(N-1)} \cdot \frac{2(q+4)\alpha_k\tau L^2}{\mu},$$

$$A_2 := \frac{16}{\mu^2} \left[2\frac{e\sigma^2 L^2}{N} + 2\frac{e\sigma^2}{N} + 4q\sigma^2 e + \frac{8e(N-s)(q+4)\sigma^2}{s(N-1)} \right],$$

$$A_3 := \frac{16}{\mu^2} \left[2\frac{\sigma^2}{N} + 4q\sigma^2 + \frac{4(N-s)2(q+4)\sigma^2}{s(N-1)} \right], A_4 := \frac{3sN+12N-15s}{s(N-1)}.$$

570 Substituting A_1, A_2, A_3, A_4 into Equation (50), we have

$$\mathbb{E}\|\theta^{k+1} - \theta^*\|^2 \leq A_1 \mathbb{E}\|\theta^k - \theta^*\|^2 + \frac{A_2}{16}\mu^2\tau a_k^2 + \frac{A_3}{16}\mu^2\tau^2 a_k^2 + A_4\mathcal{G}^k \quad (51)$$

571 Let us find $g^k \geq \max_{n \in [N]} \{(k\tau+1)G_n^k\}$, then $\mathcal{G}^k \leq \frac{g^{k^2}}{(k\tau+1)^2}$. According to Lemma 6, when $k_0 \geq$

572 $\max\{\frac{48\tau-1}{\tau}, \frac{48N-\mu^2\tau}{\mu^2\tau^2}, \frac{48L^2A_0\tau-\mu^3\tau+6L^2\mu}{\mu^2\tau^2-L^2\mu\tau}, \frac{4L^2-\mu^2}{\mu^2\tau}, \frac{4L\sqrt{\tau(\tau-1)}-\mu}{\mu\tau}\}$, it holds that $A_1 \leq 1 - \frac{1}{2}\mu\tau\alpha_k$.

573 We set $\alpha_k = \frac{4}{\mu(k\tau+1)}$. Then yield

$$\begin{aligned} \mathbb{E}\|\theta^{k+1} - \theta^*\|^2 &\leq (1 - \frac{1}{2}\mu\tau\alpha_k)\mathbb{E}\|\theta^k - \theta^*\|^2 + \frac{A_2}{16}\mu^2\tau\alpha_k^2 + \frac{A_3}{16}\mu^2\tau^2\alpha_k^2 + A_4 \frac{g^{k^2}}{(k\tau+1)^2} \\ &\leq (1 - \frac{2\tau}{k\tau+1})\mathbb{E}\|\theta^k - \theta^*\|^2 + A_2 \frac{1}{\tau} \frac{\tau^2}{(k\tau+1)^2} + A_3 \frac{\tau^2}{(k\tau+1)^2} \\ &\quad + A_4 \frac{g^{k^2}}{\tau^2} \cdot \frac{\tau^2}{(k\tau+1)^2} \end{aligned} \quad (52)$$

574 We take $k_1 = \frac{1}{\tau}$, $a = \frac{A_2}{\tau} + A_3 + A_4 \frac{g^{k_2}}{\tau^2}$ in Lemma 7 to get the final result:

$$\begin{aligned} \mathbb{E}\|\theta^k - \theta^*\|^2 &\leq \left(\frac{k_0 + 1/\tau}{k + 1/\tau}\right)^2 \mathbb{E}\|\theta^{k_0} - \theta^*\|^2 + A_2 \frac{1}{\tau} \cdot \frac{1}{k + 1/\tau} + A_3 \frac{1}{k + 1/\tau} + A_4 \frac{g^{k_2}}{\tau^2} \cdot \frac{1}{k + 1/\tau} \\ &= \left(\frac{k_0\tau + 1}{k\tau + 1}\right)^2 \mathbb{E}\|\theta^{k_0} - \theta^*\|^2 + A_2 \frac{1}{\tau} \cdot \frac{\tau}{k\tau + 1} + A_3 \frac{\tau}{k\tau + 1} + A_4 \frac{g^{k_2}}{\tau^2} \cdot \frac{\tau}{k\tau + 1}. \end{aligned} \quad (53)$$

575 \square

576 B.7 Proof of Lemma 6

577 **Lemma 6.** Denote $A_1 := 2(1 + N\alpha_k^2)(1 - \mu\alpha_k)^\tau + 4q \frac{\alpha_k \tau L^2}{\mu} + \frac{4(N-s)}{s(N-1)} \cdot \frac{2(q+4)\alpha_k \tau L^2}{\mu}$. For
 578 any $k_0 \geq \max\{\frac{48\tau-1}{\tau}, \frac{48N-\mu^2\tau}{\mu^2\tau^2}, \frac{48L^2A_0\tau-\mu^3\tau+6L^2\mu}{\mu^2\tau^2-L^2\mu\tau}, \frac{4L^2-\mu^2}{\mu^2\tau}, \frac{4L\sqrt{\tau(\tau-1)}-\mu}{\mu\tau}\}$, it holds that $A_1 \leq$
 579 $1 - \frac{1}{2}\mu\tau\alpha_k$.

580 *Proof.* According to $(1+x) \leq e^x$, we have $(1 - \mu\alpha_k)^\tau \leq e^{-\tau\mu\alpha_k}$. Due to $-\tau\mu\alpha_k \leq 0$, we get
 581 $(1 - \mu\alpha_k)^\tau \leq 1 - \tau\mu\alpha_k + \tau^2\mu^2\alpha_k^2$. Let us define $A_0 := \frac{2qL}{\mu} + \frac{4L^2(N-s)(q+4)}{s(N-1)\mu}$, then

$$\begin{aligned} A_1 &\leq 2(1 + N\alpha_k^2)(1 - \tau\mu\alpha_k + \tau^2\mu^2\alpha_k^2) + 2A_0\alpha_k\tau \\ &= 2(1 - \tau\mu\alpha_k + \tau^2\mu^2\alpha_k^2) + 2N\alpha_k^2(1 - \tau\mu\alpha_k + \tau^2\mu^2\alpha_k^2) + 2A_0\alpha_k\tau \\ &= 2 + 2(A_0 - \mu)\alpha_k\tau + 2\mu^2\alpha_k^2\tau^2 + 2N\alpha_k^2(1 - \tau\mu\alpha_k + \tau^2\mu^2\alpha_k^2). \end{aligned} \quad (54)$$

582 Since $\tau\mu\alpha_k \leq 1$, the following holds

$$A_1 \leq 2 + 2(A_0 - \mu)\alpha_k\tau + 2\mu^2\alpha_k^2\tau^2 + 2N\alpha_k^2. \quad (55)$$

583 If $a_k\tau \leq \frac{1}{12\mu}$, $a_k\tau \leq \frac{\mu\tau}{12N}$ and $\alpha_k\tau \leq \frac{\mu^2\tau-6L^2}{12L^2A_0}$ are satisfied, then $2\mu^2\alpha_k^2\tau^2 \leq \frac{1}{6}\mu\alpha_k\tau$, $2N\alpha_k^2 \leq$
 584 $\frac{1}{6}\mu\alpha_k\tau$ and $1 + 2A_0\alpha_k\tau \leq \frac{1}{6} \frac{\mu^2\tau}{L^2} \leq \frac{1}{6}\mu\alpha_k\tau$ are satisfied respectively.

585 What's more, $\alpha_k = \min\{\frac{\mu}{L^2}, \frac{1}{L\sqrt{\tau(\tau-1)}}, \frac{1}{\mu\tau}\}$. Let $\alpha_k = \frac{4}{\mu(k\mu+1)}$. So for any $k_0 \geq$
 586 $\max\{\frac{48\tau-1}{\tau}, \frac{48N-\mu^2\tau}{\mu^2\tau^2}, \frac{48L^2A_0\tau-\mu^3\tau+6L^2\mu}{\mu^2\tau^2-L^2\mu\tau}, \frac{4L^2-\mu^2}{\mu^2\tau}, \frac{4L\sqrt{\tau(\tau-1)}-\mu}{\mu\tau}\}$, it holds that $A_1 \leq 1 -$
 587 $\frac{1}{2}\mu\tau\alpha_k, \forall k \geq k_0$. \square

588 B.8 Proof of Lemma 7

589 **Lemma 7.** Let $\{v^k\}$ be a nonnegative sequence and $v^{k+1} \leq (1 - \frac{2}{k+k_1})v^k + \frac{a}{(k+k_1)^2}, \forall k \geq k_0$,
 590 where $a, k_0, k_1 > 0, k + k_1 \geq 1$. Then $\forall k \geq k_0$,

$$v^k \leq \frac{(k_0 + k_1)^2}{(k + k_1)^2} v^{k_0} + \frac{a}{k + k_1}. \quad (56)$$

591 *Proof.* The conclusion is trivial for $k = k_0$. If $v^k \leq \frac{(k_0+k_1)^2}{(k+k_1)^2} v^{k_0} + \frac{a}{k+k_1}$ holds when $k > k_0$. We
 592 imply the following inequality:

$$\begin{aligned} v^{k+1} &\leq \left(1 - \frac{2}{k+k_1}\right)v^k + \frac{a}{(k+k_1)^2} \\ &\leq \left(1 - \frac{2}{k+k_1}\right)\left[\frac{(k_0+k_1)^2}{(k+k_1)^2} v^{k_0} + \frac{a}{k+k_1}\right] + \frac{a}{(k+k_1)^2} \\ &= \frac{k+k_1-2}{k+k_1} \frac{(k_0+k_1)^2}{(k+k_1)^2} v^{k_0} + \frac{k+k_1-1}{(k+k_1)^2} a. \end{aligned} \quad (57)$$

593 Due to

$$\begin{aligned} \frac{1}{(k+k_1)^3 - 3(k+k_1) - 2} &= \frac{1}{[(k+k_1)^2 + 2(k+k_1) + 1](k+k_1 - 2)} \\ &\geq \frac{1}{(k+k_1)^3}, \quad \forall k+k_1 \geq 1, \\ \frac{1}{(k+k_1)^2 - 1} &\geq \frac{1}{(k+k_1)^2}, \quad \forall k+k_1 \geq 1. \end{aligned} \quad (58)$$

594 We get

$$\begin{aligned} \frac{k+k_1-2}{(k+k_1)^3} &\leq \frac{1}{(k+k_1+1)^2}, \\ \frac{k+k_1-1}{(k+k_1)^2} &\leq \frac{1}{k+k_1+1}. \end{aligned} \quad (59)$$

595 Plugging Equation (59) in Equation (58) implies that

$$\begin{aligned} v^{k+1} &\leq \frac{(k_0+k_1)^2}{(k+k_1+1)^2} v^{k_0} + \frac{a}{k+k_1+1} \\ &\leq \frac{(k_0+k_1)^2}{(k+k_1)^2} v^{k_0} + \frac{a}{k+k_1}. \end{aligned} \quad (60)$$

596

□

597 C Proof of Theorem 2

598 In this section, we analyze the convergence of non-convex problems.

599 C.1 Proof of Lemma 8

600 **Lemma 8.** Under Assumptions 1 to 3 and the update method of Algorithm 2 it holds that

$$\mathbb{E}f(\theta^{k+1}) \leq \mathbb{E}f(\bar{\theta}^{k,\tau}) + \underbrace{\frac{L}{2} \mathbb{E}\|\tilde{\theta}^{k+1} - \bar{\theta}^{k,\tau}\|^2}_V + \underbrace{\frac{L}{2} \mathbb{E}\|\theta^{k+1} - \tilde{\theta}^{k+1}\|^2}_{VI} + \underbrace{\langle \nabla f(\bar{\theta}^{k,\tau}), \bar{\eta}^k \rangle}_{VII}. \quad (61)$$

601 In Lemmas 9 to 11, we will analyze the upper bound of V, VI, VII .

602 *Proof.* According to Assumptions 2 and 3 it hold that

$$\begin{aligned} \mathbb{E}f(\theta^{k+1}) &= \mathbb{E}f(\theta^{k+1} - \tilde{\theta}^{k+1} + \tilde{\theta}^{k+1}) \\ &\leq \mathbb{E}f(\tilde{\theta}^{k+1}) + \mathbb{E}\langle \nabla f(\tilde{\theta}^{k+1}), \theta^{k+1} - \tilde{\theta}^{k+1} \rangle + \frac{L}{2} \mathbb{E}\|\theta^{k+1} - \tilde{\theta}^{k+1}\|^2 \\ &= \mathbb{E}f(\tilde{\theta}^{k+1}) + \frac{L}{2} \mathbb{E}\|\theta^{k+1} - \tilde{\theta}^{k+1}\|^2. \end{aligned} \quad (62)$$

603 According to Assumptions 1 and 2 the following holds

$$\begin{aligned} \mathbb{E}f(\tilde{\theta}^{k+1}) &= \mathbb{E}f(\tilde{\theta}^{k+1} - \bar{\theta}^{k,\tau} + \bar{\theta}^{k,\tau}) \\ &\leq \mathbb{E}f(\bar{\theta}^{k,\tau}) + \mathbb{E}\langle \nabla f(\bar{\theta}^{k,\tau}), \tilde{\theta}^{k+1} - \bar{\theta}^{k,\tau} \rangle + \frac{L}{2} \mathbb{E}\|\tilde{\theta}^{k+1} - \bar{\theta}^{k,\tau}\|^2 \\ &= \mathbb{E}f(\bar{\theta}^{k,\tau}) + \langle \nabla f(\bar{\theta}^{k,\tau}), \bar{\theta}^{k,\tau} + \bar{\eta}^k - \bar{\theta}^{k,\tau} \rangle + \frac{L}{2} \mathbb{E}\|\tilde{\theta}^{k+1} - \bar{\theta}^{k,\tau}\|^2 \\ &= \mathbb{E}f(\bar{\theta}^{k,\tau}) + \frac{L}{2} \mathbb{E}\|\tilde{\theta}^{k+1} - \bar{\theta}^{k,\tau}\|^2 + \langle \nabla f(\bar{\theta}^{k,\tau}), \bar{\eta}^k \rangle. \end{aligned} \quad (63)$$

604 Adding Equation (62) to Equation (63) then we yield the conclusion:

$$\mathbb{E}f(\theta^{k+1}) \leq \mathbb{E}f(\bar{\theta}^{k,\tau}) + \frac{L}{2} \mathbb{E}\|\theta^{k+1} - \tilde{\theta}^{k+1}\|^2 + \frac{L}{2} \mathbb{E}\|\tilde{\theta}^{k+1} - \bar{\theta}^{k,\tau}\|^2 + \langle \nabla f(\bar{\theta}^{k,\tau}), \bar{\eta}^k \rangle. \quad (64)$$

605

□

606 **C.2 Proof of Lemma 9**

607 **Lemma 9.** Under Assumptions 1 and 3 and the update method of Algorithm 2 it holds that

$$V \leq \frac{4(N-s)}{s(N-1)} \left[\frac{2(q+4)\alpha_k \tau L^2}{\mu} \mathbb{E} \|\theta^k - \theta^*\|^2 + 2(q+4)\alpha_k^2 \sigma^2 \tau (\tau + e) + 3\mathcal{G}^k \right]. \quad (65)$$

608 *Proof.* According to the definition of $\tilde{\theta}^{k+1}$ and $\tilde{\theta}^{k,\tau}$, we have

$$\begin{aligned} V &= \mathbb{E} \|\theta^k + \frac{1}{N} \sum_{n \in [N]} [Q(\theta_n^{k,\tau} - \theta^k) - \theta_n^{k,\tau}]\|^2 \\ &= \frac{1}{N^2} \mathbb{E} \left\| \sum_{n \in [N]} [(\theta^k - \theta_n^{k,\tau}) + Q(\theta_n^{k,\tau} - \theta^k)] \right\|^2 \\ &= \frac{1}{N^2} \mathbb{E} \left\langle \sum_{i \in [N]} (\theta^k - \theta_i^{k,\tau}) + \sum_{j \in [N]} Q(\theta_j^{k,\tau} - \theta^k), \sum_{i \in [N]} (\theta^k - \theta_i^{k,\tau}) + \sum_{j \in [N]} Q(\theta_j^{k,\tau} - \theta^k) \right\rangle. \end{aligned} \quad (66)$$

609 Let us denote $a_i = \theta^k - \theta_i^{k,\tau}$, $b_j = Q(\theta_j^{k,\tau} - \theta^k)$, then

$$\begin{aligned} \mathbb{E} b_i &= \theta_i^{k,\tau} - \theta^k + \eta_i^k \\ &= -a_i + \eta_i^k. \end{aligned} \quad (67)$$

610 We yield

$$\begin{aligned} V &= \frac{1}{N^2} \mathbb{E} \left\langle \sum_{i \in [N]} a_i + \sum_{j \in [N]} b_j, \sum_{i \in [N]} a_i + \sum_{j \in [N]} b_j \right\rangle \\ &= \frac{1}{N^2} \mathbb{E} \left[\left\langle \sum_{i \in [N]} a_i, \sum_{j \in [N]} a_j \right\rangle + 2 \left\langle \sum_{i \in [N]} a_i, \sum_{j \in [N]} b_j \right\rangle \right. \\ &\quad \left. + \left\langle \sum_{i \in [N]} b_i, \sum_{j \in [N]} b_j \right\rangle \right] \\ &= \frac{1}{N^2} \mathbb{E} \left[\left\langle \sum_{i \in [N]} a_i, \sum_{j \in [N]} a_j \right\rangle + 2 \left\langle \sum_{i \in [N]} a_i, \sum_{j \in [N]} b_j \right\rangle + \sum_{i \in [N]} \|b_i\|^2 \right. \\ &\quad \left. + 2 \sum_{i \in [N]} \sum_{j \in [N] \setminus \{i\}} \langle b_i, b_j \rangle \right] \\ &= \frac{1}{N^2} \left[\left\langle \sum_{i \in [N]} a_i, \sum_{j \in [N]} a_j \right\rangle - 2 \left\langle \sum_{i \in [N]} a_i, \sum_{j \in [N]} a_j \right\rangle + 2 \left\langle \sum_{i \in [N]} a_i, \sum_{j \in [N]} \eta_j^k \right\rangle \right. \\ &\quad \left. + \sum_{i \in [N]} \mathbb{E} \|b_i\|^2 + 2 \sum_{i \in [N]} \sum_{j \in [N] \setminus \{i\}} \langle \mathbb{E} b_i, \mathbb{E} b_j \rangle \right] \\ &= \frac{1}{N^2} \left[- \left\langle \sum_{i \in [N]} a_i, \sum_{j \in [N]} a_j \right\rangle + 2 \left\langle \sum_{i \in [N]} a_i, \sum_{j \in [N]} \eta_j^k \right\rangle + \sum_{i \in [N]} \mathbb{E} \|b_i\|^2 \right. \\ &\quad \left. + 2 \sum_{i \in [N]} \sum_{j \in [N] \setminus \{i\}} \langle -a_i + \eta_i^k, -a_j + \eta_j^k \rangle \right] \\ &= \frac{1}{N^2} \left[- \sum_{i \in [N]} \|a_i\|^2 - 2 \sum_{i \in [N]} \sum_{j \in [N] \setminus \{i\}} \langle a_i, a_j \rangle + 2 \left\langle \sum_{i \in [N]} a_i, \sum_{j \in [N]} \eta_j^k \right\rangle + \sum_{i \in [N]} \mathbb{E} \|b_i\|^2 \right. \\ &\quad \left. + 2 \sum_{i \in [N]} \sum_{j \in [N] \setminus \{i\}} \langle a_i, a_j \rangle - 4 \sum_{i \in [N]} \sum_{j \in [N] \setminus \{i\}} \langle a_i, \eta_j^k \rangle + 2 \sum_{i \in [N]} \sum_{j \in [N] \setminus \{i\}} \langle \eta_i^k, \eta_j^k \rangle \right]. \end{aligned} \quad (68)$$

611 According to Equation (39), it holds that

$$\mathbb{E} \|b_i\|^2 \leq (q+4) \|a_i\|^2 + 3 \|\eta_i\|^2. \quad (69)$$

612 So

$$\begin{aligned}
V &= \frac{1}{N^2} \left[- \sum_{i \in [N]} \|a_i\|^2 + \sum_{i \in [N]} \mathbb{E} \|b_i\|^2 + 2 \left\langle \sum_{i \in [N]} a_i, \sum_{j \in [N]} \eta_j^k \right\rangle - 4 \sum_{i \in [N]} \sum_{j \in [N] \setminus \{i\}} \langle a_i, \eta_j^k \rangle \right. \\
&\quad \left. + 2 \sum_{i \in [N]} \sum_{j \in [N] \setminus \{i\}} \langle \eta_i^k, \eta_j^k \rangle \right] \\
&\leq \frac{1}{N^2} \left[(q+3) \sum_{i \in [N]} \|a_i\|^2 + 3 \sum_{i \in [N]} \|\eta_i^k\|^2 + 2 \sum_{i \in [N]} \langle a_i, \eta_i^k \rangle + 2 \sum_{i \in [N]} \sum_{j \in [N] \setminus \{i\}} \langle a_i, \eta_j^k \rangle \right. \\
&\quad \left. - 4 \sum_{i \in [N]} \sum_{j \in [N] \setminus \{i\}} \langle a_i, \eta_j^k \rangle + 2 \sum_{i \in [N]} \sum_{j \in [N] \setminus \{i\}} \langle \eta_i^k, \eta_j^k \rangle \right] \\
&= \frac{1}{N^2} \left[(q+3) \sum_{i \in [N]} \|a_i\|^2 + 3 \sum_{i \in [N]} \|\eta_i^k\|^2 + 2 \sum_{i \in [N]} \langle a_i, \eta_i^k \rangle - 2 \sum_{i \in [N]} \sum_{j \in [N] \setminus \{i\}} \langle a_i, \eta_j^k \rangle \right. \\
&\quad \left. + 2 \sum_{i \in [N]} \sum_{j \in [N] \setminus \{i\}} \langle \eta_i^k, \eta_j^k \rangle \right] \\
&\leq \frac{1}{N^2} \left[(q+3) \sum_{i \in [N]} \|a_i\|^2 + 3 \sum_{i \in [N]} \|\eta_i^k\|^2 + \sum_{i \in [N]} \|a_i\|^2 + \sum_{i \in [N]} \|\eta_i^k\|^2 + \sum_{i \in [N]} (N-1) \|a_i\|^2 \right. \\
&\quad \left. + \sum_{i \in [N]} \sum_{j \in [N] \setminus \{i\}} \|\eta_j^k\|^2 + \sum_{i \in [N]} \sum_{j \in [N] \setminus \{i\}} (\|\eta_i^k\|^2 + \|\eta_j^k\|^2) \right] \\
&\leq \frac{1}{N^2} \left[(q+3) \sum_{i \in [N]} \|a_i\|^2 + 3 \sum_{i \in [N]} \|\eta_i^k\|^2 + \sum_{i \in [N]} \|a_i\|^2 + \sum_{i \in [N]} \|\eta_i^k\|^2 + \sum_{i \in [N]} (N-1) \|a_i\|^2 \right. \\
&\quad \left. + 2 \sum_{i \in [N]} \sum_{j \in [N] \setminus \{i\}} \|\eta_j^k\|^2 + (N-1) \sum_{i \in [N]} \|\eta_i^k\|^2 \right] \\
&\leq \frac{1}{N^2} \left[(N+q+3) \sum_{i \in [N]} \|a_i\|^2 + (N+3) \sum_{i \in [N]} \|\eta_i^k\|^2 + 2 \sum_{i \in [N]} \sum_{j \in [N] \setminus \{i\}} \|\eta_j^k\|^2 \right] \\
&\leq \frac{1}{N^2} \left[(N+q+3) \sum_{i \in [N]} \|a_i\|^2 + (N+3) \sum_{i \in [N]} \|\eta_i^k\|^2 + 2(N-1) \sum_{i \in [N]} \|\eta_i^k\|^2 \right] \\
&\leq \frac{1}{N^2} \left[(N+q+3) \sum_{i \in [N]} \|\theta^k - \theta_n^{k,\tau}\|^2 + (3N+1) \sum_{i \in [N]} \|\eta_i^k\|^2 \right].
\end{aligned} \tag{70}$$

613 On the other hand,

$$\begin{aligned}
\mathbb{E} \|\theta^k - \theta_n^{k,\tau}\|^2 &= \mathbb{E} \|\theta^k - \theta^k + \alpha_k \sum_{t=0}^{\tau-1} \nabla f_n(\theta_n^{k,t})\|^2 \\
&= \alpha_k^2 \mathbb{E} \left\| \sum_{t=0}^{\tau-1} [\hat{\nabla} f_n(\theta_n^{k,t}) - \nabla f_n(\theta_n^{k,t})] \right\|^2 + \alpha_k^2 \mathbb{E} \left\| \sum_{t=0}^{\tau-1} \nabla f_n(\theta_n^{k,t}) \right\|^2 \\
&\quad + 2\alpha_k^2 \mathbb{E} \left\langle \sum_{t=0}^{\tau-1} [\hat{\nabla} f_n(\theta_n^{k,t}) - \nabla f_n(\theta_n^{k,t})], \sum_{t=0}^{\tau-1} \nabla f_n(\theta_n^{k,t}) \right\rangle.
\end{aligned} \tag{71}$$

614 According to Assumption 3, we have $\mathbb{E}_\zeta [\hat{\nabla} f_i(\mathbf{x})] = \nabla f_i(\mathbf{x})$. It implies that $\mathbb{E} \langle \sum_{t=0}^{\tau-1} [\hat{\nabla} f_n(\theta_n^{k,t}) - \nabla f_n(\theta_n^{k,t})], \sum_{t=0}^{\tau-1} \nabla f_n(\theta_n^{k,t}) \rangle = 0$. We use Jensen's inequality to obtain the following result:

$$\begin{aligned}
\mathbb{E} \|\theta^k - \theta_n^{k,\tau}\|^2 &\leq \alpha_k^2 \tau \sum_{t=0}^{\tau-1} \mathbb{E} \|\hat{\nabla} f_n(\theta_n^{k,t}) - \nabla f_n(\theta_n^{k,t})\|^2 + \alpha_k^2 \tau \sum_{t=0}^{\tau-1} \mathbb{E} \|\nabla f_n(\theta_n^{k,t})\|^2 \\
&\leq \alpha_k^2 \tau^2 \sigma^2 + \alpha_k^2 \tau \sum_{t=0}^{\tau-1} \mathbb{E} \|\nabla f_n(\theta_n^{k,t})\|^2.
\end{aligned} \tag{72}$$

616 Substituting Equation (72) into Equation (70), we get

$$\begin{aligned}
V &\leq \frac{1}{N^2} [(N+q+3) \sum_{n \in [N]} \mathbb{E} \|\theta^k - \theta_n^{k,\tau}\|^2 + (3N+1) \sum_{n \in [N]} \|\eta_n^k\|^2] \\
&\leq \frac{1}{N^2} [(N+q+3)N\alpha_k^2\tau^2\sigma^2 + (N+q+3)\alpha_k^2\tau \sum_{n \in [N]} \sum_{t=0}^{\tau-1} \mathbb{E} \|\nabla f_n(\theta_n^{k,t})\|^2 \\
&\quad + (3N+1) \sum_{n \in [N]} \|\eta_n^k\|^2] \\
&\leq \frac{N+q+3}{N} \alpha_k^2\tau^2\sigma^2 + \frac{(N+q+3)\alpha_k^2\tau}{N^2} \sum_{n \in [N]} \sum_{t=0}^{\tau-1} \mathbb{E} \|\nabla f_n(\theta_n^{k,t})\|^2 + \frac{3N+1}{N} \mathcal{G}^k.
\end{aligned} \tag{73}$$

617

□

618 C.3 Proof of Lemma 10

619 **Lemma 10.** Under Assumptions 1 and 3 and the update method of Algorithm 2 it holds that

$$VI \leq \frac{4(N-s)}{sN(N-1)} [N(q+4)\alpha_k^2\tau^2\sigma^2 + (q+4)\alpha_k^2\tau \sum_{n \in [N]} \sum_{t=0}^{\tau-1} \mathbb{E} \|\nabla f_n(\theta_n^{k,t})\|^2 + 3N\mathcal{G}^k]. \tag{74}$$

620 *Proof.* Substituting Equation (72) into Equation (42), we yield

$$\begin{aligned}
VI &\leq \frac{4(N-s)}{sN(N-1)} \sum_{n \in [N]} [(q+4)\mathbb{E} \|\theta_n^{k,\tau} - \theta^k\|^2 + 3\|\eta_n^k\|^2] \\
&\leq \frac{4(N-s)}{sN(N-1)} \sum_{n \in [N]} [(q+4)\alpha_k^2\tau^2\sigma^2 + (q+4)\alpha_k^2\tau \sum_{t=0}^{\tau-1} \mathbb{E} \|\nabla f_n(\theta_n^{k,t})\|^2 + 3\|\eta_n^k\|^2] \\
&\leq \frac{4(N-s)}{sN(N-1)} [N(q+4)\alpha_k^2\tau^2\sigma^2 + (q+4)\alpha_k^2\tau \sum_{n \in [N]} \sum_{t=0}^{\tau-1} \mathbb{E} \|\nabla f_n(\theta_n^{k,t})\|^2 + 3N\mathcal{G}^k].
\end{aligned} \tag{75}$$

621

□

622 C.4 Proof of Lemma 11

623 **Lemma 11.** Under Assumption 1 and the update method of Algorithm 2 it holds that

$$VII = \|\nabla f(\bar{\theta}^{k,\tau})\| \sqrt{\mathcal{G}^k}. \tag{76}$$

624 *Proof.* According to Cauchy-Schwarz inequality, we have

$$\begin{aligned}
VII &= \langle \nabla f(\bar{\theta}^{k,\tau}), \bar{\eta}^k \rangle \\
&\leq \|\nabla f(\bar{\theta}^{k,\tau})\| \|\bar{\eta}^k\|^{\frac{1}{2}} \\
&= \|\nabla f(\bar{\theta}^{k,\tau})\| \left[\frac{1}{N} \sum_{n \in [N]} \|\eta_n^k\|^2 \right]^{\frac{1}{2}} \\
&\leq \|\nabla f(\bar{\theta}^{k,\tau})\| \left[\frac{1}{N} \sum_{n \in [N]} \|\eta_n^k\|^2 \right]^{\frac{1}{2}} \\
&\leq \|\nabla f(\bar{\theta}^{k,\tau})\| \left[\frac{1}{N} \sum_{n \in [N]} \mathcal{G}_n^k \right]^{\frac{1}{2}} \\
&= \|\nabla f(\bar{\theta}^{k,\tau})\| \sqrt{\mathcal{G}^k}.
\end{aligned} \tag{77}$$

625

□

626 C.5 Proof of Theorem 2

627 In this section we prove the **Theorem 2** of this study. Here, we first restate Theorem 2.

628 **Theorem C.1** (Convergence of Non-convex Settings). *Under Assumptions [1](#) to [3](#) and the update*
 629 *method of Algorithm [2](#) we define the following two constants: $F_1 = 2L^2\tau(\tau - 1)$, $F_2 = L +$*
 630 *$\frac{N+q+3}{N}\tau L + \frac{4L(N-s)}{s(N-1)}(q+4)\tau$. When the constraint $T \geq \frac{F_1^2}{(\sqrt{F_2^2+2F_1}-F_2)^2}$ is satisfied and the*
 631 *learning rate is $\alpha_k = \frac{1}{\sqrt{T}}$, the following first-order stability condition holds*

$$\frac{1}{T} \sum_{k=0}^{K-1} \sum_{t=0}^{\tau-1} \mathbb{E} \|\nabla f(\bar{\theta}^{k,t})\|^2 \leq \frac{2[f(\theta^0) - f^*]}{\sqrt{T}} + H_1 \frac{\tau-1}{T} + H_2 \frac{\tau}{\sqrt{T}} + H_3 \frac{1}{\sqrt{T}}. \quad (78)$$

632 where $g^k \leq \max_{n \in [N]} \{G_n^k\}$, $\|\eta_n^k\| \leq G_n^k$,

$$H_1 := \frac{L\sigma^2(N+1)}{N}, H_2 := \sigma^2 \left[\frac{L(N+q+3)}{N} + \frac{4L(N-s)}{sN(N-1)}N(q+4) + \frac{L}{N\tau} \right]$$

633

$$H_3 := \sum_{k=0}^{K-1} \left[\frac{12L(N-s)}{sN(N-1)}Ng^k + \frac{(3N+1)L}{N}g^k + \|\nabla f(\bar{\theta}^{k,\tau})\| \right] g^k.$$

634 *Proof.* We get the following result from Lemma 7 in (Reisizadeh et al. 2020):

$$\begin{aligned} \mathbb{E}f(\bar{\theta}^{k,\tau}) &\leq \mathbb{E}f(\theta^k) - \frac{1}{2}\alpha_k \sum_{t=0}^{\tau-1} \mathbb{E} \|\nabla f(\bar{\theta}^{k,t})\|^2 \\ &\quad - \alpha_k \left[\frac{1}{2N} - \frac{L\alpha_k}{2N} - \frac{L^2}{N}\tau(\tau-1)\alpha_k^2 \right] \sum_{t=0}^{\tau-1} \sum_{n \in [N]} \mathbb{E} \|\nabla f(\theta_n^{k,t})\|^2 + \alpha_k^2 \frac{L}{2} \frac{\sigma^2}{N} \tau \\ &\quad + \alpha_k^3 \frac{\sigma^2}{N} (N+1) \frac{\tau(\tau-1)}{2} L. \end{aligned} \quad (79)$$

635 Substituting Equations [\(65\)](#), [\(74\)](#), [\(76\)](#) and [\(79\)](#) into Equation [\(61\)](#), we get:

$$\begin{aligned} \mathbb{E}f(\theta^{k+1}) &\leq \mathbb{E}f(\theta^k) - \frac{1}{2}\alpha_k \sum_{t=0}^{\tau-1} \mathbb{E} \|\nabla f(\bar{\theta}^{k,t})\|^2 \\ &\quad - \alpha_k \left(\frac{1}{2N} - \frac{L\alpha_k}{2N} - \frac{L^2}{N}\tau(\tau-1)\alpha_k^2 \right) \sum_{t=0}^{\tau-1} \sum_{n \in [N]} \mathbb{E} \|\nabla f(\theta_n^{k,t})\|^2 + \alpha_k^2 \frac{L}{2} \frac{\sigma^2}{N} \tau \\ &\quad + \alpha_k^3 \frac{\sigma^2}{N} (N+1) \frac{\tau(\tau-1)}{2} L \\ &\quad + \frac{L}{2} \left[\frac{N+q+3}{N} \alpha_k^2 \tau^2 \sigma^2 + \frac{(N+q+3)\alpha_k^2 \tau}{N^2} \sum_{n \in [N]} \sum_{t=0}^{\tau-1} \mathbb{E} \|\nabla f_n(\theta_n^{k,t})\|^2 + \frac{3N+1}{N} \mathcal{G}^k \right] \\ &\quad + \frac{2L(N-s)}{sN(N-1)} [N(q+4)\alpha_k^2 \tau^2 \sigma^2 + (q+4)\alpha_k^2 \tau \sum_{n \in [N]} \sum_{t=0}^{\tau-1} \mathbb{E} \|\nabla f_n(\theta_n^{k,t})\|^2 + 3N\mathcal{G}^k] \\ &\quad + \|\nabla f(\bar{\theta}^{k,\tau})\| \sqrt{\mathcal{G}^k}. \end{aligned} \quad (80)$$

636 After rearranging Equation (80), we get:

$$\begin{aligned}
\mathbb{E}f(\theta^{k+1}) &\leq \mathbb{E}f(\theta^k) - \frac{1}{2}\alpha_k \sum_{t=0}^{\tau-1} \mathbb{E}\|\nabla f(\bar{\theta}^{k,t})\|^2 \\
&\quad - \alpha_k \left[\frac{1}{2N} - \frac{L\alpha_k}{2N} - \frac{L^2}{N}\tau(\tau-1)\alpha_k^2 + \frac{L(N+q+3)\alpha_k^2\tau}{2N^2} + \frac{2L(N-s)}{sN(N-1)}(q+4)\alpha_k^2\tau \right] \\
&\quad \cdot \sum_{t=0}^{\tau-1} \sum_{n \in [N]} \mathbb{E}\|\nabla f(\theta_n^{k,t})\|^2 \\
&\quad + \alpha_k^3 \frac{\sigma^2}{N}(N+1) \frac{\tau(\tau-1)}{2} L + \left[\frac{L(N+q+3)}{2N} + \frac{2L(N-s)}{sN(N-1)}N(q+4) + \frac{L}{2N\tau} \right] \alpha_k^2 \tau^2 \sigma^2 \\
&\quad + \left[\frac{6L(N-s)}{sN(N-1)}N + \frac{(3N+1)L}{2N} \right] \mathcal{G}^k + \|\nabla f(\bar{\theta}^{k,\tau})\| \sqrt{\mathcal{G}^k}.
\end{aligned} \tag{81}$$

637 Let $F_1 = 2L^2\tau(\tau-1)$, $F_2 = L + \frac{N+q+3}{N}\tau L + \frac{4L(N-s)}{s(N-1)}(q+4)\tau$. Then by Lemma 12 we conclude
638 that there exists $\alpha_k = \frac{1}{\sqrt{T}}$ such that $1 - F_2\alpha_k - F_1\alpha_k^2 \geq 0$ when constraint $T \geq \frac{F_1^2}{(\sqrt{F_2^2+2F_1-F_2})^2}$
639 holds. So

$$\begin{aligned}
\mathbb{E}f(\theta^{k+1}) &\leq \mathbb{E}f(\theta^k) - \frac{1}{2}\alpha_k \sum_{t=0}^{\tau-1} \mathbb{E}\|\nabla f(\bar{\theta}^{k,t})\|^2 \\
&\quad + \alpha_k^3 \frac{\sigma^2}{N}(N+1) \frac{\tau(\tau-1)}{2} L + \left[\frac{L(N+q+3)}{2N} + \frac{2L(N-s)}{sN(N-1)}N(q+4) + \frac{L}{2N\tau} \right] \alpha_k^2 \tau^2 \sigma^2 \\
&\quad + \left[\frac{6L(N-s)}{sN(N-1)}N + \frac{(3N+1)L}{2N} \right] \mathcal{G}^k + \|\nabla f(\bar{\theta}^{k,\tau})\| \sqrt{\mathcal{G}^k}.
\end{aligned} \tag{82}$$

640 Set $\alpha_k = \frac{1}{\sqrt{T}}$. Let us find $g^k \geq \max_{n \in [N]} \{G_n^k\}$. Then $\mathcal{G}^k \leq g^{k^2}$. Summing Equation (82) over
641 $k = 0, 1, \dots, K-1$, we obtain the following result:

$$\begin{aligned}
f^* &\leq f(\theta^0) - \frac{1}{2\sqrt{T}} \sum_{k=0}^{K-1} \sum_{t=0}^{\tau-1} \mathbb{E}\|\nabla f(\bar{\theta}^{k,t})\|^2 \\
&\quad + \frac{1}{\sqrt{T}} \frac{\sigma^2}{N}(N+1) \frac{(\tau-1)}{2} L + \left[\frac{L(N+q+3)}{2N} + \frac{2L(N-s)}{sN(N-1)}N(q+4) + \frac{L}{2N\tau} \right] \tau^2 \sigma^2 \\
&\quad + \sum_{k=0}^{K-1} \left[\frac{6L(N-s)}{sN(N-1)}N + \frac{(3N+1)L}{2N} \right] g^{k^2} + \sum_{k=0}^{K-1} \|\nabla f(\bar{\theta}^{k,\tau})\| g^k.
\end{aligned} \tag{83}$$

642 Rearranging Equation (83), we yield the following result:

$$\begin{aligned}
\frac{1}{T} \sum_{k=0}^{K-1} \sum_{t=0}^{\tau-1} \mathbb{E}\|\nabla f(\bar{\theta}^{k,t})\|^2 &\leq \frac{2}{\sqrt{T}} [f(\theta^0) - f^*] + \frac{1}{T} \frac{\sigma^2}{N}(N+1)(\tau-1)L \\
&\quad + \frac{1}{\sqrt{T}} \left[\frac{L(N+q+3)}{N} + \frac{4L(N-s)}{sN(N-1)}N(q+4) + \frac{L}{N\tau} \right] \tau \sigma^2 \\
&\quad + \frac{1}{\sqrt{T}} \sum_{k=0}^{K-1} \left[\frac{12L(N-s)}{sN(N-1)}N g^k + \frac{(3N+1)L}{N} g^k + \|\nabla f(\bar{\theta}^{k,\tau})\| \right] g^k,
\end{aligned} \tag{84}$$

643 where $\{g^k\}$ denotes a sequence of positive bounded real numbers. \square

644 C.6 Proof of Lemma 11

645 **Lemma 12.** When $T \geq \frac{F_1^2}{(\sqrt{F_2^2+2F_1-F_2})^2}$, there exists a choice of step size $\alpha_k = \frac{1}{\sqrt{T}}$, s.t. $1 - F_2\alpha_k -$
646 $F_1\alpha_k^2 \geq 0$.

647 *Proof.* We now analyze the second term of the right-hand side in Equation (81).

$$\begin{aligned}
& -\alpha_k \left[\frac{1}{2N} - \frac{L\alpha_k}{2N} - \frac{L^2}{N} \tau(\tau-1)\alpha_k^2 \right] + \frac{L(N+q+3)\alpha_k^2\tau}{2N^2} + \frac{2L(N-s)}{sN(N-1)}(q+4)\alpha_k^2\tau \\
& = \frac{-\alpha_k}{2N} \left[1 - L\alpha_k - 2L^2\tau(\tau-1)\alpha_k^2 - \frac{N+q+3}{N}\tau L\alpha_k - \frac{4L(N-s)}{s(N-1)}(q+4)\tau\alpha_k \right] \quad (85) \\
& = -\frac{\alpha_k}{2N} \left[1 - \left(L + \frac{N+q+3}{N}\tau L + \frac{4L(N-s)}{s(N-1)}(q+4)\tau \right) \alpha_k - 2L^2\tau(\tau-1)\alpha_k^2 \right]
\end{aligned}$$

648 Let $F_1 = 2L^2\tau(\tau-1)$, $F_2 = L + \frac{N+q+3}{N}\tau L + \frac{4L(N-s)}{s(N-1)}(q+4)\tau$. Then the expression $1 - F_2\alpha_k - F_1\alpha_k^2$
649 defines a downward-opening parabola in α_k , passing through the point $(0, 1)$, with its axis of
650 symmetry located at $-\frac{F_2}{2F_1} < 0$. Therefore, there exists $\exists \alpha_k > 0$ such that $1 - F_2\alpha_k - F_1\alpha_k^2 \geq 0$.
651 Specifically, this inequality holds as long as $\alpha_k \leq \frac{\sqrt{F_2^2 + 2F_1} - F_2}{F_1}$. Hence, there exists $\exists T > 0$ such
652 that choosing $\alpha_k = \frac{1}{\sqrt{T}}$ satisfies the above condition, provided that $\frac{1}{\sqrt{T}} \leq \frac{\sqrt{F_2^2 + 2F_1} - F_2}{F_1}$.

653 Therefore, when the constraint $T \geq \frac{F_1^2}{(\sqrt{F_2^2 + 2F_1} - F_2)^2}$ is satisfied, we have $1 - F_2\alpha_k - F_1\alpha_k^2 \geq 0$. \square

654 D Other Experiments

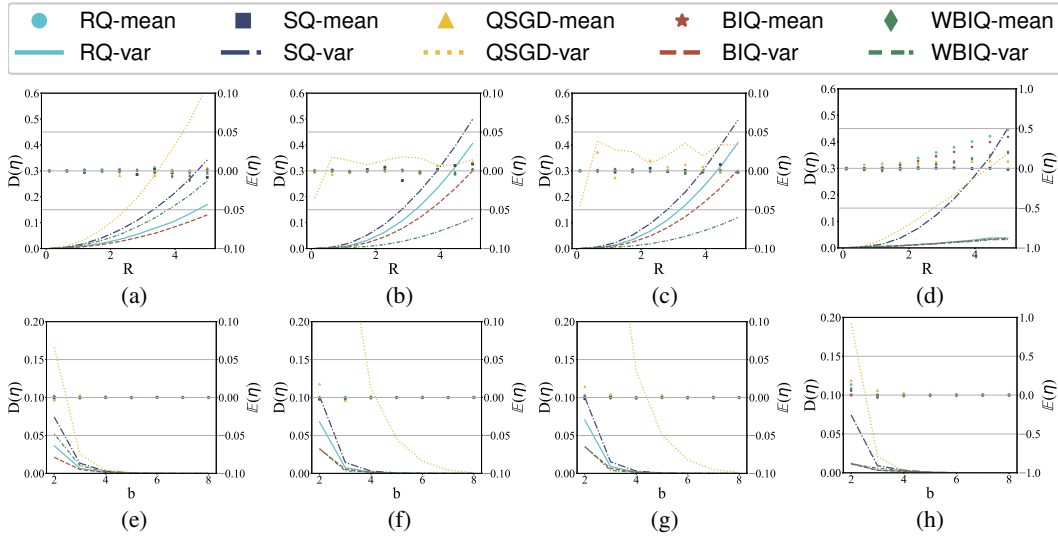


Figure 3: Mean and variance of quantization errors under different distributions, where points and lines represent the mean and variance of the quantization error, respectively: (a) (e) Uniform, (b) (f) Gaussian, (c) (g) Laplace, (d) (h) Power law distribution.

655 **Sampling.** We evaluate the performance of various quantization algorithms under uniform, Gaus-
656 sian, Laplace and Power law distributions through 10,000 sampling experiments, varying the quan-
657 tization range $[-R, R]$ from $[-0.1, 0.1]$ to $[-5, 5]$ and the quantization resolution b from 2 to 8 bits, as
658 shown in Figure 3. Specifically, we measure both the mean $\mathbb{E}(\eta)$ and variance $D(\eta)$ of the quan-
659 tization error, where $\mathbb{E}(\eta)$ reflects quantization accuracy and $D(\eta)$ indicates stability. Experimental
660 results show that $D(\eta)$ increases as R increases, consistent with the expectation that larger intervals
661 lead to greater uncertainty under low-resolution quantization. Notably, BIQ and WBIQ achieve
662 the smallest variance under Gaussian, Laplace and Power law distributions, demonstrating superior
663 stability. Under uniform distribution, BIQ also yields the lowest variance among all methods. While
664 all quantizers exhibit near-zero bias under uniform, Gaussian and Laplace settings. QSGD and SQ
665 show relatively large variances under Power law distribution despite small mean errors, indicating
666 reduced robustness. In contrast, RQ, BIQ and WBIQ, as biased quantizers, behave differently: BIQ

667 and WBIQ exhibit significantly lower variance and smaller deviation compared to RQ under Power
668 law distribution, highlighting their improved stability and accuracy. In the ablation with respect
669 to quantization resolution b , both BIQ and WBIQ achieve the lowest variance, particularly under
670 low-bit quantization. These results underscore the importance and effectiveness of BIQ and WBIQ in
671 enhancing quantization performance across diverse distributions.