# DQ-SGD: Dynamic Quantization in SGD for Communication-Efficient Distributed Learning

Guangfeng Yan[*], Shao-Lun Huang[†], Tian Lan[‡] and Linqi Song[*]
[*]*Department of Computer Science, City University of Hong Kong, Hong Kong SAR*
[*]*City University of Hong Kong Shenzhen Research Institute, Shenzhen, China*
[†]*Data Science and Information Technology Research Center, Tsinghua-Berkeley Shenzhen Institute, Shenzhen, China*
[‡]*Department of Electrical and Computer Engineering, George Washington University, Washington, DC, USA*
*Email: gfyan2-c@my.cityu.edu.hk, twn2gold@gmail.com, tlan@gwu.edu, linqi.song@cityu.edu.hk*

*Abstract*—Gradient quantization is an emerging technique in reducing communication costs in distributed learning. Existing gradient quantization algorithms often rely on engineering heuristics or empirical observations, lacking a systematic approach to dynamically quantize gradients. This paper addresses this issue by proposing a novel dynamically quantized SGD (DQ-SGD) framework, enabling us to dynamically adjust the quantization scheme for each gradient descent step by exploring the trade-off between communication cost and convergence error. We derive an upper bound, tight in some cases, of the convergence error for a restricted family of quantization schemes and loss functions. We design our DQ-SGD algorithm via minimizing the communication cost under the convergence error constraints. Finally, through extensive experiments on large-scale natural language processing and computer vision tasks on AG-News, CIFAR-10, and CIFAR-100 datasets, we demonstrate that our quantization scheme achieves better tradeoffs between the communication cost and learning performance than other state-of-the-art gradient quantization methods.

*Keywords*-Distributed Learning, Communication-efficient, Quantization

## I. Introduction

Distributed Stochastic Gradient Descent (SGD) is the core in a vast majority of distributed learning algorithms. Due to the limited bandwidth in practical networks, communication overhead for transferring gradients often becomes the performance bottleneck. Gradient quantization is an effective approach towards communication-efficient distributed learning, which uses fewer number of bits to approximate the original real value [2], [5], [6], [7], [8]. The lossy quantization inevitably brings in gradient noise, which hurts the convergence of the model. Hence, a key question is how to effectively select the number of quantization bits to balance the trade-off between the communication cost and the convergence performance.

Existing algorithms often quantize parameters into a fixed number of bits for all training iterations, which is inefficient in balancing the communication-convergence trade-off. To further reduce the communication overhead, some empirical studies began to dynamically adjust the number of quantization bits according to current model parameters in

the training process, such as the gradient's mean to standard deviation ratio [9], the training loss [10], gradient's root-mean-squared value [11]. Though these empirical heuristics of adaptive quantization methods show good performance in specific tasks, their imprecise conjectures and the lack of theoretical guidelines in the conjecture framework have limited their generalization to a broad range of machine learning models/tasks.

This paper proposes a novel dynamically quantized SGD (DQ-SGD) framework for minimizing communication overhead in distributed learning while maintaining the desired model performance. Under the assumption of smoothness and strong convexity, we first derive an upper bound on the gap between the loss after $T$ iterations and the optimal loss to characterize the convergence error caused by limited iteration steps, sampling, and quantization. Based on the above theoretical analysis, we design a dynamic quantization algorithm by minimizing the total communication cost under desired model performance constraints. Our dynamic quantization algorithm can adjust the number of quantization bits adaptively by taking into account the desired model performance, the remaining number of iterations, and the norm of gradients. We validate our theoretical analysis through extensive experiments on large-scale Natural Language Processing (NLP) and Computer Vision (CV) tasks, including text classification tasks on AG-News and image classification tasks on CIFAR-10 and CIFAR-100. Numerical results show that our proposed DQ-SGD significantly outperforms the baseline quantization methods.

To summarize, our key contributions are as follows:

• We propose a novel framework to characterize the trade-off between communication cost and convergence error by dynamically quantizing gradients in the distributed learning.

• We derive an upper bound on the convergence error for smooth strongly-convex objectives, and the upper bound is shown to be tight for a special case of quadratic functions with isotropic Hessian matrix.

• We develop a dynamically quantized SGD strategy, which is shown to achieve a fewer communication cost than fixed-bit quantization methods.

- We validate the proposed DQ-SGD on a variety of real-world datasets and machine learning models, demonstrating that our proposed DQ-SGD significantly outperforms state-of-the-art gradient quantization methods in terms of mitigating communication costs.

## II. RELATED WORK

To mitigate the communication bottleneck in distributed SGD, gradient quantization has been investigated. Different fixed number of bits quantization methods have been studied, such as 1BitSGD [5], [6], TernGrad (ternary levels) [12], QSGD (arbitrary fixed number of bits) [2].

However, these fixed-bit quantization methods may not be efficient in communication; and more efficient schemes that can dynamically adjust the number of quantization bits in different gradient descent step may have the potential to improve the communication-convergence tradeoff performance. Several studies try to construct adaptive quantization schemes through engineering heuristics or empirical observations. However, they do not come up with a solid theoretical analysis [9], [10], [11], which even results in contradicted conclusions. More specifically, MQGrad [10], and AdaQS [11] suggest using few quantization bits in early epochs and gradually increase the number of bits in later epochs; while the scheme proposed by Anders [9] states that more quantization bits should be used for the gradient with a larger root-mean-squared (RMS) value, choosing to use more bits in the early training stage and fewer bits in the later stage.

This paper's key contribution is to develop a systematic framework to crystallize the design trade-off in dynamic gradient quantization and settle this contradiction.

## III. PROBLEM FORMULATION

We consider a distributed learning system with $W$ workers and a parameter server. Data are distributed over $W$ workers, with a shared model to be jointly optimized. The local dataset of worker $i$ is $D_i$. We aim to minimize the objective function $F : \mathbb{R}^d \to \mathbb{R}$ with model parameter $\mathbf{x}$

$$\min_{\mathbf{x} \in \mathbb{R}^d} F(\mathbf{x}) = \frac{1}{W} \sum_{i=1}^{W} \mathbb{E}_{\xi \sim D_i}[l(\mathbf{x}; \xi)], \tag{1}$$

where $l(\mathbf{x}; \xi)$ is the loss of the model $\mathbf{x}$ at data point $\xi$. A standard approach to solve this problem is distributed SGD, where each worker $i$ computes its local stochastic gradient at iteration $t$ with model parameter $\mathbf{x}_t$: $\mathbf{g}_t^{(i)} = \nabla l(\mathbf{x}_t; \xi^{(i)})$. Then these local gradients are sent to the parameter server, and the server aggregates these gradients to update the model: $\mathbf{x}_{t+1} = \mathbf{x}_t - \frac{\eta}{W} \sum_{i=1}^{W} \mathbf{g}_t^{(i)}$, where $\eta$ is the learning rate. To reduce the communication cost, we consider to quantize the local stochastic gradients before sending them to the server:

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \frac{\eta}{W} \sum_{i=1}^{W} \mathcal{Q}_{b_t}[\mathbf{g}_t^{(i)}], \tag{2}$$

where $\mathcal{Q}_{b_t}[\cdot]$ is the quantization operator and $b_t$ is the number of quantization bits at iteration $t$ (in other words, we may allocate a different number of quantization bits at different iteration steps).

It is clear that the lossy compression inevitably affects the convergence of model training and deteriorates the learning performance. Therefore, we use the gap between the loss after $T$ iterations and the optimal loss to characterize the learning performance. We say the algorithm achieves an $\epsilon$-suboptimal solution if

$$F(\mathbf{x}_T) - F(\mathbf{x}^*) \le \epsilon \tag{3}$$

where $\mathbf{x}^*$ is the optimal solution to minimize $F$. Note that this suboptimal gap $\epsilon$ depends not only on the constrained communications between workers and servers, but also on the limited number of iterations $T$, the stochastic sampling, and the initial model parameter.

In this work, given the total number of training iterations $T$, the number of workers $W$, and the desired model performance $\epsilon$, we aim to adaptively adjust the number of quantization bits $b_t$ for each step to minimize the total communication cost under the model performance constraints.

Formally, **our design of DQ-SGD is to solve the following *Dynamic Quantization Problem (DQP)*:**

$$\text{(DQP):} \quad \min_{\{b_t\}} f_{\mathcal{Q}}(T, W, \{b_t\}) \\ s.t. \quad F(\mathbf{x}_T) - F(\mathbf{x}^*) \le \epsilon, \tag{4}$$

where $f_{\mathcal{Q}}(T, W, \{b_t\})$ is the incurred total communication cost of $T$ iterations and our goal is to find appropriate dynamic quantization schemes $\{b_t\}$ for $T$ iterations.

## IV. DYNAMICALLY QUANTIZED SGD

In general, the DQP problem is not easy to solve and relaxations are needed to approach this problem. Therefore, we propose to solve a relaxed version of the DQP problem and design a DQ-SGD algorithm based on the solution, which we show performs sufficiently well in practice in the experiments.

More specifically, we relax the DQP problem from the following two perspectives.

- We restrict our quantization scheme to a family of Element-Wise Uniform (EWU) quantization schemes, which are unbias with bounded variance.
- We relax the constraint in Eq. (2) to upper bound the convergence rate for smooth strongly-convex loss functions.

### A. Element-Wise Uniform Quantization

There are several types of quantization operations – categorized from different perspectives, such as grid quantization, uniform and non-uniform quantization, biased and unbiased quantization. Here, we adopt a family of stochastic quantization –EWU, similar to [2], to quantize the gradients.

In this EWU scheme, The $j$-th component of the stochastic gradient vector $\mathbf{g}$ (for any worker $i$) is quantized as

$$\mathcal{Q}_b[g_j] = \|\mathbf{g}\|_p \cdot \mathrm{sgn}(g_j) \cdot \zeta(g_j, s), \tag{5}$$

where $\|\mathbf{g}\|_p$ is the $l_p$ norm of $\mathbf{g}$; $\mathrm{sgn}(g_j) = \{+1, -1\}$ is the sign of $g_j$; $s$ is the quantization level. Note that, the quantization level is roughly exponential to the number of quantized bits. If we use $b$ bits to quantize $g_j$, we will use one bit to represent its sign and the other $b-1$ bits to represent $\zeta(g_j, s)$, thus resulting in a quantization level $s = 2^{b-1} - 1$. And $\zeta(g_j, s)$ is an unbiased stochastic function that maps scalar $|g_j|/\|\mathbf{g}\|_p$ to one of the values in set $\{0, 1/s, 2/s, \ldots, s/s\}$: if $|g_j|/\|\mathbf{g}\|_p \in [l/s, (l+1)/s]$, we have

$$\zeta(g_j, s) = \begin{cases} l/s, & \text{with probability } 1 - p_r, \\ (l+1)/s, & \text{with probability } p_r = s\dfrac{|g_j|}{\|\mathbf{g}\|_p} - l. \end{cases} \tag{6}$$

Hence, the incurred total communication cost is:

$$f_{\mathcal{Q}}(T, W, \{b_t\}) = W \sum_{t=0}^{T-1} [db_t + B_{pre}], \tag{7}$$

where $B_{pre}$ is the number of bits of full-precision floating point (e.g., $B_{pre} = 32$ or $B_{pre} = 64$) to represent $\|\mathbf{g}\|_p$. If we make the commonly used assumption for stochastic gradients as follow:

**Assumption 1** (Unbiasness and Bounded Variance of Stochastic Gradient)**.** *The stochastic gradient oracle gives us an independent unbiased estimate $\mathbf{g}$ with a bounded variance:*

$$\mathbb{E}_{\xi \sim D_i}[\mathbf{g}_t^{(i)}] = \nabla F(\mathbf{x}_t), \tag{8}$$

$$\mathbb{E}_{\xi \sim D_i}[\|\mathbf{g}_t^{(i)} - \nabla F(\mathbf{x}_t)\|_2^2] \leq \sigma^2. \tag{9}$$

Then we have the following lemma to characterize the aggregated stochastic gradient $\hat{\mathbf{g}}_t \triangleq \frac{1}{W} \sum_{i=1}^{W} \mathcal{Q}_{b_t}[\mathbf{g}_t^{(i)}]$, and the proof is given in Appendix A.

**Lemma 1** (Unbiasness and Bounded Variance of EWU)**.** *For the local gradient $\mathbf{g}_t^{(i)}$, if the number of quantization bits for all $W$ workers are all $b_t$, then the aggregated gradient $\hat{\mathbf{g}}_t$ satisfies:*

$$\mathbb{E}[\hat{\mathbf{g}}_t] = \nabla F(\mathbf{x}_t) \tag{10}$$

*and*

$$\mathbb{E}\left[\|\hat{\mathbf{g}}_t\|_2^2\right] \leq \|\nabla F(\mathbf{x}_t)\|_2^2 + \underbrace{\frac{\sigma^2}{W}}_{\text{Sampling Noise}} + \underbrace{\frac{d}{4W(2^{b_t-1}-1)^2}\bar{G}_t^2}_{\text{Quantization Noise}}, \tag{11}$$

where $\bar{G}_t^2 = \frac{1}{W}\sum_{i=1}^{W} \|\mathbf{g}_t^{(i)}\|_p^2$, is the mean square of all local gradient $l_p$ norms at iteration $t$.

Eq. (10) means that the aggregated gradient $\hat{\mathbf{g}}_t$ is the unbiased estimate of $\nabla F(\mathbf{x})$. Eq. (11) implies that the difference between $\|\hat{\mathbf{g}}_t\|_2^2$ and $\|\nabla F(\mathbf{x}_t)\|_2^2$ consists of two parts: the first part is the sampling noise, which inversely proportional to $W$; the second part is the quantization noise, which is proportional to $\bar{G}_t^2$ and decays exponentially with the increase of the number quantization bits $b_t$.

### B. Upper Bounded Convergence Rate for Smooth and Strongly Convex Functions

We first state some assumptions as follows.

**Assumption 2** (Smoothness)**.** *The objective function $F(\mathbf{x})$ is $L$-smooth, if $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, $\|\nabla F(\mathbf{x}) - \nabla F(\mathbf{y})\|_2 \leq L\|\mathbf{x} - \mathbf{y}\|_2$.*

It implies that $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, we have

$$F(\mathbf{y}) \leq F(\mathbf{x}) + \nabla F(\mathbf{x})^{\mathrm{T}}(\mathbf{y} - \mathbf{x}) + \frac{L}{2}\|\mathbf{y} - \mathbf{x}\|_2^2 \tag{12}$$

$$\|\nabla F(\mathbf{x})\|_2^2 \leq 2L[F(\mathbf{x}) - F(\mathbf{x}^*)] \tag{13}$$

**Assumption 3** (Strong convexity)**.** *The objective function $F(\mathbf{x})$ is $\mu$-strongly convex, if $\exists \mu > 0$, $F(\mathbf{x}) - \dfrac{\mu}{2}\mathbf{x}^{\mathrm{T}}\mathbf{x}$ is a convex function.*

From Assumption 3, we have: $\forall \mathbf{x} \in \mathbb{R}^d$,

$$\|\nabla F(\mathbf{x})\|_2^2 \geq 2\mu[F(\mathbf{x}) - F(\mathbf{x}^*)] \tag{14}$$

Putting the quantized SGD (2) on smooth, strongly convex functions yield the following result with proof given in Appendix B.

**Theorem 1** (Convergence Error Bound of Strongly Convex Objectives)**.** *For the problem in Eq.* (1) *under Assumption 2, 3 and 1 with initial parameter $\mathbf{x}_0$, using quantized gradients in Eq.* (2) *for iteration, we can upper bound the convergence error by*

$$\mathbb{E}[F(\mathbf{x}_T) - F(\mathbf{x}^*)]$$

$$\leq \underbrace{\alpha(\eta)^T[F(\mathbf{x}_0) - F(\mathbf{x}^*)] + \frac{L\eta^2\sigma^2[1 - \alpha(\eta)^T]}{2W(1 - \alpha(\eta))}}_{\text{Error of Distributed SGD}}$$

$$+ \underbrace{\frac{Ld\eta^2}{8W}\sum_{t=0}^{T-1}\alpha(\eta)^{T-1-t}\frac{\bar{G}_t^2}{(2^{b_t-1}-1)^2}}_{\text{Quantization Error}}$$

$$\overset{T\to\infty}{\to} \frac{L\eta^2\sigma^2}{2W(1 - \alpha(\eta))} + \frac{Ld\eta^2}{8W}\sum_{t=0}^{T-1}\alpha(\eta)^{T-1-t}\frac{\bar{G}_t^2}{(2^{b_t-1}-1)^2} \tag{15}$$

*where $\alpha(\eta) =: 1 - 2\mu\eta + L\mu\eta^2$ (We abbreviate $\alpha(\eta)$ as $\alpha$ in the following section.).*

We can see that the convergence error consists of two parts: the first two terms are the error of the distributed SGD method, which is independent of the quantization algorithms. This part error can be reduced by increasing the number of iterations $T$ and also depends on the learning rate $\eta$ (from the expression of $\alpha$, we can see that when $\eta \leq 1/L$, with the increase of $\eta$, $\alpha$ decrease, and the convergence rate of the model is accelerated); The last term is **quantization error**, resulted from the lossy quantization of gradients. It is obtained by the weighted accumulation of quantization noise at each iteration and directly increases the convergence error floor. Note that $\alpha$ is less than 1. Thus the weight given to quantization noise decays exponentially as the number of intervening iterations increases. Accordingly, this is sometimes called an **exponential recency-weighted average**.

We can see that the quantization error decays exponentially in the number of quantization bits $b_t$. When the number of quantization bits at each iteration is large enough (e.g., $b_t = 32$), the quantization error tends to 0, but the communication cost is significantly high.

We aim to use as little communication cost as possible to ensure that the quantization error is below a given level $\epsilon_Q = [1 - \gamma]\epsilon$, where $1 - \gamma$ is a tradeoff factor representing the contribution of convergence error by quantization. Furthermore, we have $\lim_{T\to\infty} \epsilon_Q = \epsilon - \frac{L\eta^2\sigma^2}{2W(1-\alpha)}$.

### C. DQ-SGD Algorithm

Given the above two relaxations, we can rewrite the DQP as

$$
\min_{\{b_t\}} \quad W \sum_{t=0}^{T-1} (db_t + B_{pre}),
$$
$$
s.t. \quad \frac{Ld\eta^2}{8W} \sum_{t=0}^{T-1} \alpha^{T-1-t} \frac{\bar{G}_t^2}{(2^{b_t-1}-1)^2} = \epsilon_Q. \tag{16}
$$

By solving the above optimization problem, we can determine the $\{b_t\}$ at every iteration step:

$$
b_t = \log_2 \left[ \sqrt{\frac{T}{\hat{\epsilon}_Q}} \alpha^{(T-1-t)/2} \bar{G}_t + 1 \right] + 1 \tag{17}
$$

where $\hat{\epsilon}_Q \triangleq \frac{8W}{Ld\eta^2}\epsilon_Q$. We can see that the number of quantization bits is determined by three key factors: (i) the desired quantization error upper bound $\epsilon_Q$, the smaller the desired quantization error is, the more quantization bits are needed; (ii) the iteration step $t$, the number of bits is increasing as the training process goes on; (iii) the root-mean-square of local gradient norms $\bar{G}_t$, gradients with a larger norm should be quantized using more bits.

The pseudocode is given in Algorithm 1. We have a set of $W$ workers who proceed in synchronous steps, and each worker has a complete copy of the model. In each communication round, workers compute their local gradients and communicate quantized gradients with the parameter server (lines 3-5), while the server aggregates these gradients from workers and updates the model parameters (lines 8-10). If $Q_{b_t}[\mathbf{g}_t^{(i)}]$ is the quantized stochastic gradients in the $i$-th worker and $\mathbf{x}_t$ is the model parameter that the workers hold in iteration $t$, then the updated value of $\mathbf{x}$ by the end of this iteration is: $\mathbf{x}_{t+1} = \mathbf{x}_t + \eta\hat{\mathbf{g}}_t$, where $\hat{\mathbf{g}}_t = \frac{1}{W}\sum_{l=1}^{W} Q_{b_t}[\mathbf{g}_t^{(i)}]$. Note that we determine $b_{t+1}$ according to the gradients information at iteration $t$ (lines 11), so we update the compression bits every $\tau$ iterations in practice (In our experience, we take $\tau = 100$).

**DQ-SGD outperforms fixed-bit quantization based SGD in communication cost.** Compared with the fixed-bit algorithms, our proposed DQ-SGD can achieve the same performance with fewer communication costs. The communication cost of DQ-SGD and fixed-bit algorithms are shown as follows with proofs given in Appendix C.

**Theorem 2.** *For the problem in Eq.* (1) *under Assumptions 1, 2, 3, with initial parameter $\mathbf{x}_0$, using the dynamic quantizer in Eqs.* (17) *to quantize the gradients, then the total communication cost for DQ-SGD is upper bounded by*

$$
f_Q(T, W, b_t) \leq WdT \log_2 \sqrt{\frac{T\left[2L[F(\mathbf{x}_0) - \mathbf{x}^*] + \sigma^2\right]}{\hat{\epsilon}_Q}}
$$
$$
+ WTB_{pre} + WTd + \frac{WTd}{2}\log_2 GM(\alpha) \tag{18}
$$

*If we want to achieve the same model performance, the total communication cost of the fixed-bit algorithms is upper bounded by*

$$
f_Q(T, W, b_t) \leq WdT \log_2 \sqrt{\frac{T\left[2L[F(\mathbf{x}_0) - \mathbf{x}^*] + \sigma^2\right]}{\hat{\epsilon}_Q}}
$$
$$
+ WTB_{pre} + WTd + \frac{WTd}{2}\log_2 AM(\alpha) \tag{19}
$$

*where Arithmetic Mean $AM(\alpha) = \frac{1}{T}\sum_{t=0}^{T-1} \alpha^t = \frac{1}{T}\frac{1-\alpha^T}{1-\alpha}$ and Geometric Mean $GM(\alpha) = \left[\prod_{t=0}^{T-1} \alpha^t\right]^{\frac{1}{T}} = \alpha^{\frac{T-1}{2}}$.*

We can see that if we desire a lower quantization error, we need more communication costs. Note that $0 < \alpha < 1$, so $AM(\alpha) > GM(\alpha)$, which means our proposed DQ-SGD uses fewer communication cost compared with the fixed-bit algorithms.

**Algorithm 1:** DQ-SGD in Distributed Learning

---
**Input:** Iterations number $T$, desired quantization error upper bound $\epsilon_Q$, learning rate $\eta$, initial point $\mathbf{x}_0 \in \mathbb{R}^d$, initial number of quantization bits $b_0$, hyper-parameters $\alpha$

**Output:** $\mathbf{x}_T$

1 **for** $t = 0, 1, ..., T-1$ **do**
2    **On each worker** $l = 1, ..., W$**:**
3    Compute local gradient $\mathbf{g}_t^{(i)}$;
4    Quantize the gradient $\mathcal{Q}_{b_t}[\mathbf{g}_t^{(i)}]$ according to Eq. (5);
5    Send $\mathcal{Q}_{b_t}[\mathbf{g}_t^{(i)}]$ to server;
6    Receive $\mathbf{x}_{t+1}$ and $b_{t+1}$ from server;
7    **On server:**
8    Collect all $W$ quantized gradients $\mathcal{Q}_{b_t}[\mathbf{g}_t^{(i)}]$ from workers;
9    Average: $\hat{\mathbf{g}}_t = \frac{1}{W}\sum_{l=1}^{W} \mathcal{Q}_{b_t}[\mathbf{g}_t^{(i)}]$;
10    Update the global parameters $\mathbf{x}_{t+1} = \mathbf{x}_t - \eta\hat{\mathbf{g}}_t$;
11    Update the quantization bits $b_{t+1}$ according to Eq. (17);
12    Send $\mathbf{x}_{t+1}$ and $b_{t+1}$ to all workers;
13 **end**

---

## V. DISCUSSIONS

### A. Convergence Error for Quadratic Objectives.

In previous sections, we use the upper bound of convergence error to measure the model performance. In this subsection, we will prove that there exist strongly convex functions $F(\mathbf{x})$ where the convergence error bound in Theorem 1 is tight (i.e., The '=' in Eq. (15) can be achieved.).

For general quadratic functions, we can employ gradient flow[1] to calculate an exact convergence error. We have the relationship between the aggregated stochastic gradients and full gradients: $\hat{\mathbf{g}}_t = \nabla F(\mathbf{x}_t) + \boldsymbol{\epsilon}_t$. Based on the central limit theorem, it is assumed that $\boldsymbol{\epsilon}_t$ follows the Gaussian distribution, that is $\boldsymbol{\epsilon}_t \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}(\mathbf{x}_t))$. Then using analysis within the gradient flow framework, we can get the following theorem.

**Theorem 3** (Exact Convergence Error for Quadratic Objectives). *For a quadratic optimization objective function $F(\mathbf{x}) = 1/2\mathbf{x}^{\mathrm{T}}\mathbf{H}\mathbf{x} + \mathbf{A}^{\mathrm{T}}\mathbf{x} + B$, consider the perturbed gradient descent dynamics*

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta\nabla F(\mathbf{x}_t) - \eta\boldsymbol{\epsilon}_t, \boldsymbol{\epsilon}_t \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}(\mathbf{x}_t)) \quad (20)$$

---

[1]when the learning rate is infinitesimal, the stochastic gradient descent process can be regarded as a stochastic dynamic system.

*We can achieve*

$$\mathbb{E}[F(\mathbf{x}_T) - F(\mathbf{x}^*)]$$
$$= \frac{1}{2}(\mathbf{x}_0 - \mathbf{x}^*)^{\mathrm{T}}(\boldsymbol{\rho}(\eta)^T)^{\mathrm{T}}\mathbf{H}\boldsymbol{\rho}(\eta)^T(\mathbf{x}_0 - \mathbf{x}^*)$$
$$+ \frac{\eta^2}{2}\sum_{t=0}^{T-1}\mathrm{Tr}\left[\boldsymbol{\rho}(\eta)^{T-1-t}\boldsymbol{\Sigma}(\mathbf{x}_t)\mathbf{H}\left(\boldsymbol{\rho}(\eta)^{T-1-t}\right)^{\mathrm{T}}\right] \quad (21)$$

*where $\boldsymbol{\rho}(\eta) := \mathbf{I} - \eta\mathbf{H}$, and $\mathbf{H}$ is the Hessian matrix.*

Detailed proof is in Appendix D. We can see that the convergence error consists of two parts: the error of the gradient descent method, which is linearly convergent; the error due to gradient estimation error (data sampling noise, gradient quantization error).

Consider the case where the Hessian matrix is isotropic $\mathbf{H} = \lambda\mathbf{I}$, and let $\beta(\eta) := 1 - 2\eta\lambda + \eta^2\lambda^2$, then Eq.(21) can be rewrite as

$$\mathbb{E}[F(\mathbf{x}_T) - F(\mathbf{x}^*)] = \beta(\eta)^T[F(\mathbf{x}_0) - F(\mathbf{x}^*)]$$
$$+ \frac{\lambda\eta^2}{2}\sum_{t=0}^{T-1}\beta(\eta)^{T-1-t}\mathrm{Tr}[\boldsymbol{\Sigma}(\mathbf{x}_t)] \quad (22)$$

In Lemma 1, if we let the gradient noise always reach the upper bound value, then we can get

$$\mathrm{Tr}[\boldsymbol{\Sigma}(\mathbf{x}_t)] = \mathbb{E}\left[\|\hat{\mathbf{g}}_t - \nabla F(\mathbf{x}_t)\|_2^2\right]$$
$$= \frac{\sigma^2}{W} + \frac{d}{4W(2^{b_t-1}-1)^2}\bar{G}_t \quad (23)$$

Plugging Eq. (23) into Eq. (22), then the '=' in Eq. (15) can be achieved, and proves that the upper bound for strongly convex objectives in Theorem 1 is tight in this case.

### B. Algorithm Implementation Details

Although Eq. (17) provide valuable insights about how to adjust $b_t$ over time, it is still challenging to use it in practice due to the convergence rate $\alpha$ being known. Inspired by [16], we propose a straightforward rule where we approximate $F(\mathbf{x}^*)$ to 0 and the learning rate $\eta$ is small enough. We estimate $\alpha$ as follows according to Theorem 1:

$$\alpha_{est} = \left[\frac{F(\mathbf{x}_t)}{F(\mathbf{x}_0)}\right]^{1/t} \quad (24)$$

where $F(\mathbf{x}_0)$ and $F(\mathbf{x}_t)$ can be easily obtained in the training.

### C. Dynamic Adjustment in the Number of Bits

From Eq. (17), we can see that two factors may affect how to adjust the number of quantization bits: the increased weight $\alpha^{(T-1-t)/2}$ and the root-mean-square of local gradient norm $\bar{G}_t$. $\alpha^{(T-1-t)/2}$ increases with the iterations $t$, and $\bar{G}_t$ gets smaller as the training process goes on. Therefore,

• **Decreasing in Communication.** If the decreasing rate of the root-mean-square of local gradient norm (i.e., $\frac{\bar{G}_{t+1}}{\bar{G}_t}$)

smaller than $\sqrt{\alpha}$, then $b_{t+1} < b_t$, which means the number of quantization bits decreases with the iteration step;

• **Increasing in Communication.** On the contrary, if the decreasing rate $\frac{\bar{G}_{t+1}}{\bar{G}_t}$ is bigger than $\sqrt{\alpha}$, then $b_{t+1} > b_t$, meaning that the number of quantization bits increases with the iteration step.

## VI. EXPERIMENTS

In this section, we conduct experiments on NLP and CV tasks on three datasets: AG-News [13], CIFAR-10 [14], and CIFAR-100 [14], to validate the effectiveness of our proposed DQ-SGD method. We conduct experiments for $W = 8$ workers and use canonical networks to evaluate the performance of different algorithms: BiLSTM on the text classification task on the AG-News dataset, Resnet18 on the image classification task on the CIFAR-10 dataset, and Resnet34 on the image classification task on the CIFAR-100 dataset. Other parameters information is shown in Table I. We use test accuracy to measure the learning performance. We compare our proposed DQ-SGD with the following baselines: SignSGD [5], TernGrad [12], QSGD [2], Adaptive [11] and AdaQS [9].

**Test Accuracy vs Communication Cost**. Figure 1 and table II compare the test accuracy and communication cost of different algorithms under different tasks. The communication cost of Vanilla SGD are 1313.29 GB, 1998.06 GB, and 3805.54 GB, and the test accuracy can be achieved are 0.9016, 0.8815, and 0.6969 on AG-News, CIFAR-10, CIFAR-100, respectively. We set a communication budget as $15\%$ of the communication cost incurred by the SGD and a performance threshold as $99.7\%$ of the test accuracy achieved by the SGD for all three tasks. We can see that our proposed algorithm is the only one that satisfies a high-performance and low communication cost (the upper left region). Other baselines cannot achieve the performance threshold given the communication budget.

**Fixed Quantization vs. Dynamic Quantization**. Figure 2 shows the comparison results of the Fixed Bits algorithm and our proposed DQ-SGD on CIFAR-10. Figure 2 (a) and Figure 2 (b) show the test accuracy curves and the training loss curves. Figure 2 (c) shows the bits allocation of each iteration of DQ-SGD. Fixed Bits (6 bits) and DQ-SGD can get almost the same accuracy as SGD. However, the communication cost of DQ-SGD is reduced up to $25\%$ compared with that of Fixed Bits (6 bits). It can be seen that our dynamic quantization strategy can effectively reduce the communication cost compared with the fixed quantization scheme. Figure 3 shows the accuracy of Fixed Bits and DQ-SGD under different communication costs. It can be seen that DQ-SGD can achieve higher test accuracy than Fixed Bits under the same communication cost.

## VII. CONCLUSION

This paper proposes a novel adaptive gradient quantization strategy called DQ-SGD to reduce the communication cost of distributed computing based on theoretical analysis. DQ-SGD adjusts quantization bits automatically by considering the norm of gradient and current iteration number. The experimental results of image classification and text classification show that DQ-SGD is superior to state-of-the-art gradient quantization methods in reducing communication costs.

## VIII. ACKNOWLEDGMENT

## REFERENCES

[1] Sebastian U Stich, Jean-Baptiste Cordonnier, and Martin Jaggi. Sparsified SGD with memory. In Advances in Neural Information Processing Systems, pp. 4447–4458, 2018.

[2] Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. QSGD: Communication-efficient SGD via gradient quantization and encoding. In Advances in Neural Information Processing Systems, pp. 1709–1720, 2017.

[3] McMahan, B., Moore, E., Ramage, D., Hampson, S., and y Arcas, B. A. Communication-efficient learning of deep networks from decentralized data. In Artificial Intelligence and Statistics, pp. 1273–1282. 2017.

[4] Stich, S. U. Local sgd converges fast and communicates little. In International Conference on Learning Representations, 2018.

[5] Frank Seide, Hao Fu, Jasha Droppo, Gang Li, and Dong Yu. 1-bit stochastic gradient descent and its application to data-parallel distributed training of speech DNNs. Conference of the International Speech Communication Association, pp. 1058–1062, 2014.

[6] Bernstein, J., Wang, Y. X., Azizzadenesheli, K., and Anandkumar, A. SignSGD: Compressed optimisation for non-convex problems. In International Conference on Machine Learning (pp. 560-569), 2018.

[7] Wu, J., Huang, W., Huang, J., and Zhang, T. Error compensated quantized SGD and its applications to large-scale distributed optimization. In International Conference on Machine Learning (pp. 5325-5333), 2018.

[8] Ananda Theertha Suresh, Felix X Yu, Sanjiv Kumar, and H Brendan Mcmahan. Distributed mean estimation with limited communication. International Conference on Machine Learning, pp.3329–3337, 2017.

[9] Jinrong Guo, Wantao Liu, Wang Wang, Jizhong Han, Ruixuan Li, Yijun Lu, and Songlin Hu. Accelerating distributed deep learning by adaptive gradient quantization. In IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 1603–1607, 2020.

(a) AG-News (SGD:1313.29GB, 0.9016)     (b) CIFAR10 (SGD:1998.06GB, 0.8815)     (c) CIFAR100 (SGD:3805.54GB, 0.6969)

Figure 1.   Performance comparison with SOTA on AG-News, CIFAR-10, and CIFAR-100.



(a) Test accuracy     (b) Training loss     (c) Bits allocation

Figure 2.   The learning process of Fixed Bits and DQ-SGD on CIFAR-10.

Table I

PARAMETERS

| Dataset | Net | Learning rate | Batchsize | Interations | Hyperparameters in DQSGD |
|---------|-----|---------------|-----------|-------------|--------------------------|
| AG-News | BiLSTM | 0.005 | 32 | 1000 | $k = 5, \alpha = 0.994$ |
| CIFAR-10 | Resnet18 | 0.1 | 32 | 6000 | $k = 20, \alpha = 0.999$ |
| CIFAR-100 | Resnet34 | 0.01 | 64 | 6000 | $k = 10, \alpha = 0.999$ |

Table II

TEST ACCURACY VS. COMMUNICATION COST.

| | AG-News | | CIFAR-10 | | CIFAR-100 | |
|---|---|---|---|---|---|---|
| | Accuracy | Communication Cost | Accuracy | Communication Cost | Accuracy | Communication Cost |
| Vanilla SGD | 0.9016 | 1313.29 | 0.8815 | 1998.06 | 0.6969 | 3805.54 |
| SignSGD | 0.8663 | 41.04 | 0.5191 | 62.44 | 0.3955 | 118.92 |
| TernGrad | 0.8480 | 82.08 | 0.7418 | 124.88 | 0.6174 | 237.85 |
| QSGD (4 bits) | 0.8894 | 164.16 | 0.8545 | 249.76 | 0.6837 | 475.69 |
| QSGD (6 bits) | 0.9006 | 246.24 | 0.8803 | 374.64 | 0.6969 | 713.54 |
| Adaptive | 0.8991 | 201.12 | 0.8787 | 361.97 | 0.6943 | 641.74 |
| AdaQS | 0.9001 | 201.13 | 0.8809 | 373.47 | 0.6960 | 744.72 |
| DQSGD (Ours) | 0.8997 | **192.85** | 0.8793 | **280.98** | 0.6959 | **565.46** |

[10] Guoxin Cui, Jun Xu, Wei Zeng, Yanyan Lan, Jiafeng Guo, and Xueqi Cheng. MQGrad: Reinforcement learning of gradient quantization in parameter server. In Proceedings of the 2018 ACM SIGIR International Conference on Theory of Information Retrieval, pp. 83–90, 2018.

[11] Anders Oland and Bhiksha Raj. Reducing communication overhead in distributed learning by an order of magnitude (almost). In IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 2219–2223, 2015.

[12] Wei Wen, Cong Xu, Feng Yan, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. TernGrad: Ternary gradients to reduce communication in distributed deep learning. Advances in Neural Information Processing Systems, pp. 1508–1518, 2017.

Figure 3. Test accuracy of Fixed Bits and DQ-SGD under different communication cost on CIFAR-10.

[13] Bottou, L., Curtis, F. E., and Nocedal, J. Optimization methods for large-scale machine learning. Siam Review, 60(2):223–311, 2018.

[14] Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level convolutional networks for text clas-sification. InAdvances in Neural Information Processing Systems, pp. 649–657, 2015.

[15] Krizhevsky Alex, and Geoffrey Hinton. Learning multiple layers of features from tiny images, 2009.

[16] Wang Jianyu and Gauri Joshi. "Adaptive Communication Strategies to Achieve the Best Error-Runtime Trade-off in Local-Update SGD." MLSys. 2019.

## APPENDIX A.
### PROOF OF LEMMA 1

According to Eq. (6), we have

$$\mathbb{E}[\zeta(g_j, s)] = \frac{l}{s}[1 - s\frac{|g_j|}{\|\mathbf{g}\|_p} + l] + \frac{l+1}{s}[s\frac{|g_j|}{\|\mathbf{g}\|_p} - l] = \frac{|g_j|}{\|\mathbf{g}\|_p}$$

Then, we have

$$\mathbb{E}[\zeta(g_j, s)^2] = \mathbb{E}[\zeta(g_j, s)]^2 + \mathbb{V}[\zeta(g_j, s)]$$
$$= \frac{|g_j|^2}{\|\mathbf{g}\|_p^2} + \frac{1}{s^2}p_r(1 - p_r) \leq \frac{|g_j|^2}{\|\mathbf{g}\|_p^2} + \frac{1}{4s^2}$$

Considering that $\mathcal{Q}_b(g_j) = \|\mathbf{g}\|_p \cdot \text{sgn}(g_j) \cdot \zeta(g_j, s)$, we have

$$\mathbb{E}[\|\mathcal{Q}_b[\mathbf{g}]\|_2^2] = \sum_{j=0}^{d} \mathbb{E}[\|\mathbf{g}\|_p^2 \zeta(g_j, s)^2]$$
$$\leq \sum_{j=0}^{d} \|\mathbf{g}\|_p^2(\frac{|g_j|^2}{\|\mathbf{g}\|_p^2} + \frac{1}{4s^2}) = \|\mathbf{g}\|_2^2 + \frac{d}{4s^2}\|\mathbf{g}\|_p^2$$

So we can get $\mathbb{E}[\mathcal{Q}_b[\mathbf{g}]] = \mathbf{g}$ and

$$\text{Tr}\left\{\mathbb{V}\left[\mathcal{Q}_{b_t}[\mathbf{g}^{(i)})]\right]\right\} = \mathbb{E}[\|\mathcal{Q}_b[\mathbf{g}]\|_2^2 - \|\mathbf{g}\|_2^2 \leq \frac{d\|\mathbf{g}\|_p^2}{4(2^{b-1} - 1)^2}.$$

Hence, for the aggregated gradient $\hat{\mathbf{g}}_t \triangleq \frac{1}{W}\sum_{i=1}^{W} \mathcal{Q}_{b_t}[\mathbf{g}_t^{(i)}]$:

$$\mathbb{E}[\hat{\mathbf{g}}_t]] = \frac{1}{W}\sum_{i=1}^{W}\mathbb{E}[\mathbf{g}_t^{(i)}] \overset{(a)}{=} \nabla F(\mathbf{x}_t)$$

$$\mathbb{E}\left[\|\hat{\mathbf{g}}_t\|_2^2\right] = \text{Tr}\left\{\mathbb{V}[\hat{\mathbf{g}}_t]\right\} + \mathbb{E}[\hat{\mathbf{g}}_t]^{\text{T}}\mathbb{E}[\hat{\mathbf{g}}_t]$$
$$= \frac{1}{W^2}\sum_{i=1}^{W}\text{Tr}\left\{\mathbb{V}\left[\mathcal{Q}_{b_t}[\mathbf{g}_t^{(i)})]\right]\right\} + \|\frac{1}{W}\sum_{i=1}^{W}\mathbf{g}_t^{(i)}\|_2^2$$
$$\leq \frac{d}{4W^2(2^{b_t-1} - 1)^2}\sum_{i=1}^{W}\|\mathbf{g}_t^{(i)}\|_p^2 + \|\frac{1}{W}\sum_{i=1}^{W}\mathbf{g}_t^{(i)}\|_2^2$$
$$\leq \frac{d}{4W^2(2^{b_t-1} - 1)^2}\sum_{i=1}^{W}\|\mathbf{g}_t^{(i)}\|_p^2 + \frac{1}{W^2}\sum_{i=1}^{W}\|\mathbf{g}_t^{(i)}\|_2^2$$
$$\overset{(a)}{\leq} \frac{d}{4W(2^{b_t-1} - 1)^2}\bar{G}_t^2 + \|\nabla F(\mathbf{x}_t)\|_2^2 + \frac{\sigma^2}{W}$$

where $(a)$ uses the Assumption 1, and $\bar{G}_t^2 = \frac{1}{W}\sum_{i=1}^{W}\|\mathbf{g}_t^{(i)}\|_p^2$.

## APPENDIX B.
### PROOF OF THEOREM 1

Considering function $F$ is $L-$smooth, and using Assumption 2:

$$F(\mathbf{x}_{t+1}) \leq F(\mathbf{x}_t) + \nabla F(\mathbf{x}_t)^{\text{T}}(\mathbf{x}_{t+1} - \mathbf{x}_t) + \frac{L}{2}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|_2^2$$

Due to $\mathbf{x}_{t+1} = \mathbf{x}_t - \frac{\eta}{W}\sum_{i=1}^{W}\mathcal{Q}_{b_t}[\mathbf{g}_t^{(i)}]$, then:

$$F(\mathbf{x}_{t+1}) \leq F(\mathbf{x}_t) + \nabla F(\mathbf{x}_t)^{\text{T}}(-\frac{\eta}{W}\sum_{i=1}^{W}\mathcal{Q}_{b_t}[\mathbf{g}_t^{(i)}])$$
$$+ \frac{L}{2}\| - \frac{\eta}{W}\sum_{i=1}^{W}\mathcal{Q}_{b_t}[\mathbf{g}_t^{(i)}]\|_2^2$$

Taking total expectations, and using Lemma 1, this yields:

$$\mathbb{E}[F(\mathbf{x}_{t+1}) - F(\mathbf{x}_t)] \leq (-\eta + \frac{L\eta^2}{2})\|\nabla F(\mathbf{x}_t)\|_2^2$$
$$+ \frac{L\eta^2\sigma^2}{2W} + \frac{L\eta^2 d}{8W(2^{b_t-1} - 1)^2}\bar{G}_t^2$$

Considering that function $F$ is $\mu-$strongly convex, and using Assumption 3, so:

$$\mathbb{E}[F(\mathbf{x}_{t+1}) - F(\mathbf{x}_t)] \leq -(2\mu\eta - L\mu\eta^2)[F(\mathbf{x}_t) - F(\mathbf{x}^*)]$$
$$+ \frac{L\eta^2\sigma^2}{2W} + \frac{L\eta^2 d}{8W(2^{b_t-1} - 1)^2}\bar{G}_t^2$$

Subtracting $F(\mathbf{x}^*)$ from both sides, and let $\alpha(\eta) := 1 - $

$2\mu\eta + L\mu\eta^2$, so:

$$\mathbb{E}[F(\mathbf{x}_{t+1}) - F(\mathbf{x}^*)] \leq \alpha[F(\mathbf{x}_t) - F(\mathbf{x}^*)] + \frac{L\eta^2\sigma^2}{2W}$$
$$+ \frac{L\eta^2 d}{8W(2^{b_t-1}-1)^2}\bar{G}_t^2$$

Applying this recursively, we conclude the proof.

## APPENDIX C.
## PROOF OF THEOREM 2

The bits allocation is: $b_t \approx \log_2\left[\sqrt{\frac{T}{\hat{\epsilon}_Q}}\alpha^{(T-1-t)/2}\bar{G}_t + 1\right] + 1$, so

$$f_{\mathcal{Q}}(T, W, b_t) = W\sum_{t=0}^{T-1}[db_t + B_{pre}]$$

$$\approx Wd\sum_{t=0}^{T-1}\log_2\left[\sqrt{\frac{T}{\hat{\epsilon}_Q}}\alpha^{(T-1-t)/2}\bar{G}_t\right] + WTB_{pre} + WTd$$

$$< WdT\log_2\sqrt{\frac{T\bar{G}_0^2}{\hat{\epsilon}_Q}} + \frac{WdT(T-1)}{4}\log_2\alpha + WTB_{pre} + WTd$$

$$< WdT\log_2\sqrt{\frac{T[2L[F(\mathbf{x}_0)-\mathbf{x}^*]+\sigma^2]}{\hat{\epsilon}_Q}} + WTB_{pre}$$

$$+ WTd + \frac{WdT(T-1)}{4}\log_2\alpha$$

Accordingly, if we fix the number of quantization bits (i.e., $b_t = b$), then we have

$$\hat{\epsilon}_Q = \sum_{t=0}^{T-1}\alpha^{T-1-t}\frac{\bar{G}_t^2}{(2^{b_t-1}-1)^2}$$

$$\leq \frac{\bar{G}_0^2}{(2^{b-1}-1)^2}\sum_{t=0}^{T-1}\alpha^{T-1-t}$$

$$\leq \frac{2L[F(\mathbf{x}_0)-\mathbf{x}^*]+\sigma^2}{(2^{b-1}-1)^2}\frac{1-\alpha^T}{1-\alpha}$$

So,

$$b \leq \log_2\left[\sqrt{\frac{2L[F(\mathbf{x}_0)-\mathbf{x}^*]+\sigma^2}{\hat{\epsilon}_Q}\frac{1-\alpha^T}{1-\alpha}} + 1\right] + 1$$

So, the total communication cost for the fixed bits

algorithm is

$$f_{\mathcal{Q}}(T, W, b_t) = W\sum_{t=0}^{T-1}[db_t + B_{pre}] = WTdb + WTB_{pre}$$

$$= WTd\log_2\left[\sqrt{\frac{2L[F(\mathbf{x}_0)-\mathbf{x}^*]+\sigma^2}{\hat{\epsilon}_Q}\frac{1-\alpha^T}{1-\alpha}} + 1\right]$$

$$+ WTB_{pre} + WTd$$

$$\approx WTd\log_2\sqrt{\frac{2L[F(\mathbf{x}_0)-\mathbf{x}^*]+\sigma^2}{\hat{\epsilon}_Q}} + WTB_{pre}$$

$$+ WTd + WTd\log_2\sqrt{\frac{1-\alpha^T}{1-\alpha}}$$

If we let $AM(\alpha) = \frac{1}{T}\sum_{t=0}^{T-1}\alpha^t = \frac{1}{T}\frac{1-\alpha^T}{1-\alpha}$ and $GM(\alpha) = (\prod_{t=0}^{T-1}\alpha^t)^{(1/T)} = \alpha^{(T-1)/2}$.

Then the total communication cost for DQ-SGD is

$$f_{\mathcal{Q}}(T, W, \{b_t\}) \leq WdT\log_2\sqrt{\frac{T[2L[F(\mathbf{x}_0)-\mathbf{x}^*]+\sigma^2]}{\hat{\epsilon}_Q}}$$

$$+ WTB_{pre} + WTd + \frac{WTd}{2}\log_2 GM(\alpha)$$

and the total communication cost for fixed bits is

$$f_{\mathcal{Q}}(T, W, \{b_t\}) \leq WdT\log_2\sqrt{\frac{T[2L[F(\mathbf{x}_0)-\mathbf{x}^*]+\sigma^2]}{\hat{\epsilon}_Q}}$$

$$+ WTB_{pre} + WTd + \frac{WTd}{2}\log_2 AM(\alpha)$$

## APPENDIX D.
## PROOF OF THEOREM 3

For a quadratic optimization problem $F(\mathbf{x}) = 1/2\mathbf{x}^T\mathbf{H}\mathbf{x} + \mathbf{A}^T\mathbf{x} + B$, we consider a Gaussian noise case

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta\nabla F(\mathbf{x}_t) - \eta\boldsymbol{\epsilon}_t, \boldsymbol{\epsilon}_t \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}(\mathbf{x}_t))$$

Then we have

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta\nabla F(\mathbf{x}_t) - \eta\boldsymbol{\epsilon}_t$$
$$= \mathbf{x}_t - \eta[\mathbf{H}\mathbf{x}_t + \mathbf{A}] - \eta\boldsymbol{\epsilon}_t$$
$$= (\mathbf{I} - \eta\mathbf{H})\mathbf{x}_t - \eta\mathbf{A} - \eta\boldsymbol{\epsilon}_t$$

Considering $\nabla F(\mathbf{x}^*) = \eta\mathbf{A} + \eta\mathbf{H}\mathbf{x}^* = 0$, subtracting $\mathbf{x}^*$ from both sides, and rearranging, this yields:

$$\mathbf{x}_{t+1} - \mathbf{x}^* = (\mathbf{I} - \eta\mathbf{H})\mathbf{x}_t - \eta\mathbf{A} - \mathbf{x}^* - \eta\boldsymbol{\epsilon}_t$$
$$= (\mathbf{I} - \eta\mathbf{H})(\mathbf{x}_t - \mathbf{x}^*) - \eta\mathbf{A} - \eta\mathbf{H}\mathbf{x}^* - \eta\boldsymbol{\epsilon}_t$$
$$= (\mathbf{I} - \eta\mathbf{H})(\mathbf{x}_t - \mathbf{x}^*) - \eta\boldsymbol{\epsilon}_t$$

Applying this recursively, let $\boldsymbol{\rho} = \mathbf{I} - \eta\mathbf{H}$, we have:

$$\mathbf{x}_T - \mathbf{x}^* = \boldsymbol{\rho}^T(\mathbf{x}_0 - \mathbf{x}^*) - \sum_{t=0}^{T-1}[\eta\boldsymbol{\rho}^{T-1-t}\boldsymbol{\epsilon}_t]$$

Considering that $\epsilon_t \sim \mathcal{N}(0, \mathbf{\Sigma}(\mathbf{x}_t))$, then:

$$\sum_{t=0}^{T-1}[\eta\boldsymbol{\rho}^{T-1-t}\epsilon_t] = \sum_{t=0}^{T-1}[\eta\boldsymbol{\rho}^{T-1-t}\mathbf{\Sigma}(\mathbf{x}_t)^{\frac{1}{2}}\mathcal{N}(\mathbf{0},\mathbf{I})]$$

$$= \sum_{t=0}^{T-1}[\eta\boldsymbol{\rho}^{T-1-t}\mathbf{\Sigma}(\mathbf{x}_t)^{\frac{1}{2}}[\mathbf{W}(t+1)-\mathbf{W}(t)]\} \equiv I(T)$$

where, $\mathbf{W}$ is a standard $d$-dimensional Wiener process, and $I(T)$ is an Ito integral. Hence $\mathbf{x}_T = \mathbf{x}^* + \boldsymbol{\rho}^T(\mathbf{x}_0-\mathbf{x}^*) - I(T)$, then:

$$F(\mathbf{x}_T) = \frac{1}{2}\mathbf{x}_T^{\mathrm{T}}\mathbf{H}\mathbf{x}_T + \mathbf{A}^{\mathrm{T}}\mathbf{x}_T + B$$

$$= \frac{1}{2}(\mathbf{x}^{(0)}-\mathbf{x}^*)^{\mathrm{T}}(\boldsymbol{\rho}^T)^{\mathrm{T}}\mathbf{H}\boldsymbol{\rho}^T(\mathbf{x}_0-\mathbf{x}^*) + \frac{1}{2}I(T)^{\mathrm{T}}\mathbf{H}I(T)$$
$$+ F(\mathbf{x}^*) - [\boldsymbol{\rho}^T(\mathbf{x}_0-\mathbf{x}^*)+\mathbf{x}^*+\mathbf{A}]^{\mathrm{T}}\mathbf{H}I(T)$$

Subtracting $F(\mathbf{x}^*)$ from both sides, taking total expectations, and rearranging, this yields:

$$\mathbb{E}[F(\mathbf{x}_T)-F(\mathbf{x}^*)]$$
$$= \frac{1}{2}(\mathbf{x}_0-\mathbf{x}^*)^{\mathrm{T}}(\boldsymbol{\rho}^T)^{\mathrm{T}}\mathbf{H}\boldsymbol{\rho}^T(\mathbf{x}_0-\mathbf{x}^*) + \frac{1}{2}\mathbb{E}[I(T)^{\mathrm{T}}\mathbf{H}I(T)]$$
$$- [\boldsymbol{\rho}^T(\mathbf{x}_0-\mathbf{x}^*)+\mathbf{x}^*+\mathbf{A}]^{\mathrm{T}}\mathbf{H}\mathbb{E}[I(T)]$$

The property of Ito integral $I(T)$ is:

$$\mathbb{E}[I(T)] = 0$$

$$\mathbb{E}[I(T)^{\mathrm{T}}\mathbf{H}I(T)] = \sum_{t=0}^{T-1}\eta^2\mathrm{Tr}[\boldsymbol{\rho}^{T-1-t}\mathbf{\Sigma}(\mathbf{x}_t)\mathbf{H}(\boldsymbol{\rho}^{T-1-t})^{\mathrm{T}}]$$

Using this property, we have:

$$\mathbb{E}[F(\mathbf{x}_T)-F(\mathbf{x}^*)] = \frac{1}{2}(\mathbf{x}_0-\mathbf{x}^*)^{\mathrm{T}}(\boldsymbol{\rho}^T)^{\mathrm{T}}\mathbf{H}\boldsymbol{\rho}^T(\mathbf{x}_0-\mathbf{x}^*)$$
$$+ \frac{\eta^2}{2}\sum_{t=0}^{T-1}\mathrm{Tr}[\boldsymbol{\rho}^{T-1-t}\mathbf{\Sigma}(\mathbf{x}_t)\mathbf{H}(\boldsymbol{\rho}^{T-1-t})^{\mathrm{T}}]$$

If we consider a simple example: the Hessian matrix is isotropic $\mathbf{H} = \lambda\mathbf{I}$, let $\alpha(\eta) := 1 - 2\eta\lambda + \eta^2\lambda^2$, so

$$first = \frac{1}{2}(\mathbf{x}_0-\mathbf{x}^*)^{\mathrm{T}}(\boldsymbol{\rho}^T)^{\mathrm{T}}\mathbf{H}\boldsymbol{\rho}^T(\mathbf{x}_0-\mathbf{x}^*)$$

$$= \alpha(\eta)^T\frac{1}{2}(\mathbf{x}_0-\mathbf{x}^*)^{\mathrm{T}}\mathbf{H}(\mathbf{x}_0-\mathbf{x}^*)$$

$$= \alpha(\eta)^T[\frac{1}{2}\mathbf{x}_0^{\mathrm{T}}\mathbf{H}\mathbf{x}_0 + \frac{1}{2}\mathbf{x}^{*\mathrm{T}}\mathbf{H}\mathbf{x}^* - \mathbf{x}_0^{\mathrm{T}}\mathbf{H}\mathbf{x}^*]$$

$$= \alpha(\eta)^T[\frac{1}{2}\mathbf{x}_0^{\mathrm{T}}\mathbf{H}\mathbf{x}_0 + \frac{1}{2}\mathbf{x}^{*\mathrm{T}}\mathbf{H}\mathbf{x}^* - \mathbf{x}_0^{\mathrm{T}}\mathbf{H}\mathbf{x}^*$$
$$+ \mathbf{x}_0^{\mathrm{T}}(\mathbf{H}\mathbf{x}^* + A) - \mathbf{x}^{*\mathrm{T}}(\mathbf{H}\mathbf{x}^* + A)]$$

$$= \alpha(\eta)^T[\frac{1}{2}\mathbf{x}_0^{\mathrm{T}}\mathbf{H}\mathbf{x}_0 - \frac{1}{2}\mathbf{x}^{*\mathrm{T}}\mathbf{H}\mathbf{x}^* + \mathbf{x}_0^{\mathrm{T}}A - \mathbf{x}^{*\mathrm{T}}A]$$
$$= \alpha(\eta)^T[F(\mathbf{x}_0)-F(\mathbf{x}^*)]$$

$$second = \frac{\lambda\eta^2}{2}\sum_{t=0}^{T-1}\alpha(\eta)^{T-1-t}\mathrm{Tr}[\mathbf{\Sigma}(\mathbf{x}_t)]$$

Thus,

$$\mathbb{E}[F(\mathbf{x}_T)-F(\mathbf{x}^*)] = \alpha(\eta)^T[F(\mathbf{x}_0)-F(\mathbf{x}^*)]$$
$$+ \frac{\lambda\eta^2}{2}\sum_{t=0}^{T-1}\alpha(\eta)^{T-1-t}\mathrm{Tr}[\mathbf{\Sigma}(\mathbf{x}_t)]$$

## APPENDIX E.
### PROOF OF EQ. (17)

For the constrain function of Eq. (16), we have

$$\frac{\partial^2[\sum_{t=0}^{T-1}\alpha^{T-1-t}\frac{\bar{G}_t^2}{(2^{b_t-1}-1)^2}]}{\partial b_t^2}$$

$$= \frac{(\ln 2)^2 * (2^{b_t}+1) * 2^{b_t}}{(2^{b_t-1}-1)^4}\alpha^{T-1-t}\bar{G}_t^2 > 0$$

Therefore, this optimization problem is a convex optimization problem, then we have the Lagrange function

$$\mathcal{L}(b_t, \lambda) = W\sum_{t=0}^{T-1}(db_t + B_{pre})$$

$$+ \lambda\left[\sum_{t=0}^{T-1}\alpha^{T-1-t}\frac{\bar{G}_t^2}{(2^{b_t-1}-1)^2} - \epsilon_Q\right]$$

where $\lambda$ is Lagrange multiplier. Then we can get

$$\frac{\partial\mathcal{L}(b_t,\lambda)}{\partial b_t} = Wd - \lambda\frac{2\ln 2 * 2^{b_t-1}}{(2^{b_t-1}-1)^3}\alpha^{T-1-t}\bar{G}_t^2 = 0$$

By solve this equation, we can get

$$b_t \approx \log_2\left[\sqrt{\frac{2\lambda\ln 2}{Wd}}\alpha^{(T-1-t)/2}\bar{G}_t + 1\right] + 1$$

Setting $\sum_{t=0}^{T-1}\alpha^{T-1-t}\frac{\bar{G}_t^2}{(2^{b_t-1}-1)^2} = \hat{\epsilon}_Q$, we can solve $\lambda$ and get the result.