

# High-Dimensional Stochastic Gradient Quantization for Communication-Efficient Edge Learning

Yuqing Du, Sheng Yang and Kaibin Huang

## Abstract

Edge machine learning involves the deployment of learning algorithms at the wireless network edge so as to leverage massive mobile data for enabling intelligent applications. The mainstream edge learning approach, federated learning, has been developed based on distributed gradient descent. Based on the approach, stochastic gradients are computed at edge devices and then transmitted to an edge server for updating a global AI model. Since each stochastic gradient is typically high-dimensional (with millions to billions of coefficients), communication overhead becomes a bottleneck for edge learning. To address this issue, we propose in this work a novel framework of hierarchical stochastic gradient quantization and study its effect on the learning performance. First, the framework features a practical hierarchical architecture for decomposing the stochastic gradient into its norm and normalized block gradients, and efficiently quantizes them using a uniform quantizer and a low-dimensional codebook on a Grassmann manifold, respectively. Subsequently, the quantized normalized block gradients are scaled and cascaded to yield the quantized normalized stochastic gradient using a so-called hinge vector designed under the criterion of minimum distortion. The hinge vector is also efficiently compressed using another low-dimensional Grassmannian quantizer. The other feature of the framework is a bit-allocation scheme for reducing the quantization error. The scheme divides the total bits from gradient quantization to determine the resolutions of the low-dimensional quantizers in the proposed framework. The framework is proved to guarantee model convergency by analyzing the convergence rate as a function of the quantization bits. Furthermore, by simulation, our design is shown to substantially reduce the communication overhead compared with the state-of-the-art signSGD scheme, while both achieve similar learning accuracies.

## I. INTRODUCTION

Recently, a large amount of data are generated in real-time and distributed at edge devices (e.g., smartphones and sensors). To allow rapid access to the enormous real-time data generated

Y. Du and K. Huang are with The University of Hong Kong, Hong Kong (Email: yqdu@eee.hku.hk, huangkb@eee.hku.hk). S. Yang is with Laboratory of Signals and Systems, CentraleSupélec, University of Paris-Saclay, 91190 Gif-sur-Yvette, France (e-mail: sheng.yang@centralesupelec.fr).

by edge devices for *artificial intelligence* (AI)-model training, several edge learning frameworks such as *federated edge learning* (FEEL) have been developed based on distributed stochastic gradient descent [1], [2]. Based on the approach, stochastic gradients are computed at edge devices and then transmitted to an edge server for aggregation and then updating a global AI-model. Typical stochastic gradients are of high dimensionality (each constitutes e.g., millions of parameters). Thus their transmission over communication networks can result in extensive overhead and a bottleneck for fast edge learning. To tackle this challenge, numerous schemes have been developed for compressing stochastic gradients to reduce the said communication overhead [3]–[5]. However, due to gradients’ high dimensionality, most of the existing schemes focus merely on scalar quantization. The area of high-dimensional *vector quantization* (VQ) targeting stochastic gradients is largely uncharted. This motivates us to make the first attempt on filling the void. Specifically, we propose a novel hierarchical vector quantization scheme using low-dimensional Grassmannian codebooks. By both simulation and theoretic analyses, this scheme is shown to be of low-complexity, communication-efficient, and guarantee learning convergence.

#### A. Stochastic Gradient Quantization

Recently, the topic of stochastic gradient quantization has been attracting growing interests for its being a key approach for improving the communication efficiency of edge learning [6]–[11]. In [6], a scheme called “Quantized SGD” (QSGD) is proposed, where a scalar quantizer is deployed and its efficiency is improved by Elias integer coding exploiting the distribution of quantized gradient values. Building on QSGD, the effect of the quantization error can be further alleviated using an error-compensation scheme presented in [8]. To be specific, the accumulated quantization error is exploited to accelerate the model convergence. A recent key advancement in the area is the finding that despite its supposed low resolution, the combination of one-bit scalar quantizer for gradient-coefficient quantization, named “signSGD” [7], and momentum in descent can be proved to achieve a convergence rate of the same order as its counterpart without quantization, namely the famous “ADAM” scheme [10]. The promising result has motivated a series of followup work. For example, the original signSGD can be improved by dividing the large-number of one-bit coefficients of a quantized gradient into blocks and scaling each block by the norm of the unquantized counterpart [9]. Such modifications are shown to accelerate

learning. All the above schemes are based on scalar quantization. There are few results on the VQ of gradients despite its being a well-developed area [11].

VQ, namely the joint quantization of the entries of a vector, is required to achieve the optimal rate-distortion trade-off [11]. Such asymptotic optimality, however, comes at the price of an exponentially growing complexity with the vector length. This makes it infeasible to directly apply the classic VQ algorithms to the quantization of high-dimensional stochastic gradients and explains the current popularity of scalar quantization for stochastic gradient quantization in edge learning. However, the effectiveness of VQ proven in conventional data compression suggests its potential for improving the communication efficiency for edge learning. This motivates the current work on designing a new VQ framework for stochastic gradient compression targeting SGD.

### *B. Grassmannian Quantization in Wireless Communication*

A Grassmann manifold refers to a space of lines or subspaces embedded in a higher-dimensional space. A quantizer for partitioning the manifold typically uses a Grassmannian codebook that comprises a set of lines or subspaces and a subspace distance as the distortion measure (e.g., the sine of two lines' separation angle). Consequently, the quantizer attempts to minimize the deviation in direction between a line and its quantized version, or the deviation in orientation for the case of a subspace. Thus, Grassmannian quantization is a suitable tool for compressing data containing information on vector direction or subspace orientation.

In wireless communication, Grassmannian quantization is widely adopted in one particular area, limited feedback, for efficient feedback of a quantized beamformer/precoder from a receiver to a transmitter to enable adaptive multi-antenna transmission [12]–[17]. In [13], for beamformer quantization and feedback, a randomly generated Grassmannian codebook, which comprises a set of unitary vectors uniformly distributed on the Grassmann manifold, is proved to be asymptotically optimal under the criterion of rate maximization as the codebook size grows. On the other hand, for a finite codebook size and a MIMO channel with rich scattering, an important finding is reported in [13], [14] that the optimal beamforming/precoding codebook design can be translated into the mathematical problem of Grassmannian line/subspace packing. The result, however, may not hold for correlated channels. This motivates a vein of research on developing systematic methods for Grassmannian codebook construction targeting

spatially/temporally correlated channels, where the correlation is exploited for reducing the required codebook size and hence the feedback overhead [15]–[18]. The latest research in the area is focused on next-generation massive MIMO communication with large-scale arrays. The resultant channels are high-dimensional and thus the limited feedback techniques developed previously for small-scale MIMO cannot be directly applied. The main approach for overcoming the challenge is to decompose a massive MIMO channel into the low-dimensional slow and fast time varying components, corresponding to channel spatial correlation and small-scale fading, respectively [19]–[21]. Then periodic feedback is needed only for a beamformer/precoder matched to the small-scale MIMO fading channels. The low-dimensional feedback can be then compressed using a traditional Grassmannian quantization technique, thereby reining in feedback overhead.

Stochastic gradient quantization for FEEL is related to limited feedback in that they both aim at compressing a vector for efficient transmission to convey directional information. To be specific, one critical information conveyed by a stochastic gradient is the descending direction on a surface generated by a given loss function, which measures the learning accuracy. Though Grassmannian quantization seems to be a suitable tool for gradient compression, the direct application is impractical as the codebook size and computation complexity both increase exponentially with the gradient’s dimensionality. The techniques developed in the other area of limited feedback for massive MIMO with large-scale channels cannot be transferred to high-dimensional gradient quantization. The reason is that the former’s effectiveness hinges on the channel’s structural and multi-time-scale properties but no similar counterparts exist for stochastic gradients. This calls for the development of a new approach for high-dimensional gradient quantization.

### *C. Contributions and Organization*

This work addresses the issue of practical quantization of high-dimensional stochastic gradients to realize communication-efficient FEEL in a wireless system. To this end, we propose a novel framework of hierarchical gradient quantization based on gradient decomposition and low-dimensional Grassmannian quantization. The effect of the framework on the learning performance is characterized by analyzing the model convergence rate as a function of the number of quantization bits, which quantifies the communication overhead. To the best of the authors’ knowledge, this work represents the first attempt on applying Grassmannian quantization to the

compression of high-dimensional stochastic gradients.

The specific contributions of this work are summarized as follows.

- **Hierarchical Quantization Architecture:** The practicality of the proposed framework arises from a hierarchical architecture for intelligent gradient decomposition and low-dimensional component quantization. First, the stochastic gradient is decomposed into its norm and the normalized stochastic gradient. The norm is easily compressed using a scalar quantizer. Next, as a key feature of the framework, the high-dimensional normalized stochastic gradient is intelligently decomposed into 1) a set of equal-length unitary vectors called normalized block gradients and 2) a mentioned hinge vector, which is also unitary and integrates normalized block gradients to yield the normalized stochastic gradient. Such decomposition has two advantages. The hinge vector harnesses a certain level of high-dimensional VQ gain even though only practical low-dimensional quantizers are deployed. The other advantage is that the unitary nature of normalized block gradients and hinge vectors allow them to be efficiently compressed using two Grassmannian quantizers.
- **Bit-Allocation Scheme:** Under an overhead constraint, the total number of bits from quantizing a stochastic gradient is fixed. The allocation of the bits to control the resolutions of the quantizers in the proposed framework can be optimized under the criterion of minimum distortion, yielding a bit-allocation scheme. By average distortion analysis, as the length of block gradients increases, it can be proved that the randomness of the hinge vector diminishes and it converges to a known fixed point (a vector) on the Grassmann manifold, which is a vector and denoted as  $\mathbf{h}_{\text{ref}}$ . This suggests that in this asymptotic regime, all bits should be allocated to quantizing and transmitting normalized block gradients and the stochastic gradient norm with  $\mathbf{h}_{\text{ref}}$  being used at the server as a surrogate for the hinge vector. Otherwise, in the non-asymptotic regime, the optimal bit-allocation is derived in closed-form by minimizing the sum distortion.
- **Analysis of Learning Convergence Rate:** It is proved that the proposed hierarchical quantization scheme leads to the convergence of the FEEL algorithm even if the loss function is non-convex. The specific findings are three-fold: First, given a large number of edge devices, the convergence rate is asymptotically  $O\left(\frac{1}{\sqrt{N}}\right)$  with  $N$  denoting the total number of iterations; Second, the quantization error leads to a biased term on the upper bound of the expected gradient norm, where the increment of quantization bits reduces

the value of this biased term, giving rise to a faster convergence speed; Third, given the quantization bits, this bias vanishes at the rate of  $O(\frac{1}{K})$  with  $K$  denoting the total number of edge devices, which also accelerates model convergence.

*Organization:* The remainder of the paper is organized as follows. Section II introduces the FEEL system model and provides the problem formulation. Section III presents the hierarchical quantization scheme. The distortion analysis for the proposed scheme is given in Section IV, building on which, the bit-allocation strategy is derived in Section V. Section VI presents the convergence rate analysis of the learning algorithm with the proposed hierarchical quantization. Simulation results are provided in Section VII followed by concluding remarks in Section VIII.

## II. SYSTEM MODEL AND PROBLEM FORMULATION

Consider a FEEL system as illustrated in Fig. 1, where an edge server trains an AI model (e.g., a classifier), represented by the parameter set  $\theta$ , using training datasets distributed among  $K$  edge devices.

To facilitate the learning, the loss function measuring the model error is defined as follows. Let  $\{\mathcal{D}_k\}$  denote the local dataset collected at the  $k$ -th edge device. The local loss function of the model vector  $\theta$  on  $\{\mathcal{D}_k\}$  is given by

$$\text{(Local loss function)} \quad f_k(\theta) = \frac{1}{|\mathcal{D}_k|} \sum_{(\mathbf{x}_i, y_i) \in \mathcal{D}_k} f(\theta, \mathbf{x}_i, y_i), \quad (1)$$

where  $f(\theta, \mathbf{x}_i, y_i)$  is the sample-wise loss function quantifying the prediction error of the model  $\theta$  on the training sample  $\mathbf{x}_i$  w.r.t its ground-true label  $y_i$ . Without loss of generality, by assuming uniform sizes for local datasets:  $|\mathcal{D}_k| = D$ , the global loss function of the model vector  $\theta$  on all distributed local datasets can be written as

$$\text{(Global loss function)} \quad F(\theta) = \frac{1}{K} \sum_{k=1}^K f_k(\theta). \quad (2)$$

The learning process is to minimize the global loss function  $F(\theta)$ , which can be mathematically written as

$$\theta^* = \arg \min F(\theta). \quad (3)$$

In the context of FEEL, the gradient-averaging implementation is proposed in [2] to tackle the privacy issue by avoiding uploading all the local data. Specifically, in each iteration, say the  $n$ -th

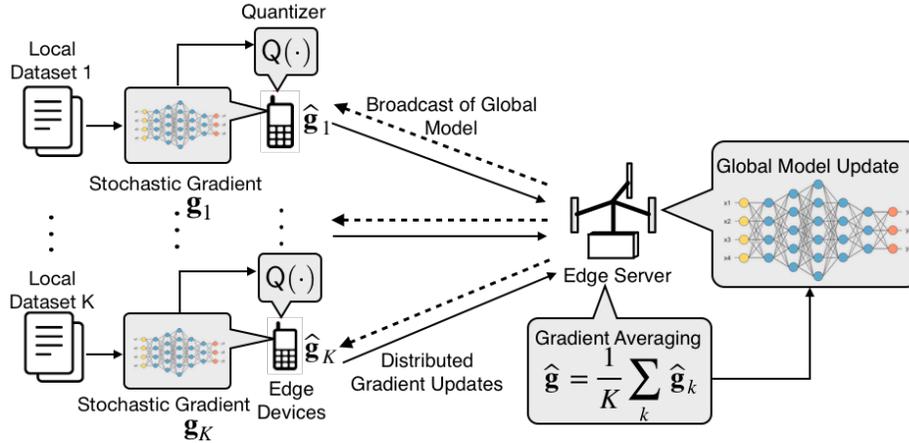


Figure 1. Federated edge learning system.

iteration, the edge server broadcasts the current model under training  $\theta[n]$  to all edge devices. Based on the received current model  $\theta[n]$ , each device computes the stochastic gradient by differentiating the local loss function defined in (1). Mathematically, for device  $k$ , the stochastic gradient of the  $n$ -th iteration can be computed as:

$$\text{(Stochastic gradient)} \quad \mathbf{g}_k[n] = \nabla f_k(\theta[n]), \quad (4)$$

where  $\nabla$  denotes the gradient operation. Conventionally, upon its completion, the local gradients are sent to the edge server for averaging. However, in practical applications, communicating the gradients in each iteration has been observed to be a significant performance bottleneck [6], which will become exacerbated especially when the gradient is dense. This motivates the lossy compression of the gradients before transmission, mathematically defined as follows:

$$\text{(Stochastic gradient quantization)} \quad \hat{\mathbf{g}}_k[n] = Q(\mathbf{g}_k[n]), \quad (5)$$

where  $Q(\cdot)$  maps any point  $\mathbf{g}$  in  $\mathbb{R}^{\text{Dim} \times 1}$  to one of the codewords in the codebook  $\mathcal{C}$ , i.e.  $\hat{\mathbf{g}}_k[n] \in \mathcal{C}$  with  $\hat{\mathbf{g}}_k[n]$  denoting the quantized version of the stochastic gradient  $\mathbf{g}_k[n]$ . Rather than conveying the quantized version, the edge devices communicate the codeword index to the edge server. It is further assumed that the edge server has the knowledge of the codebook  $\mathcal{C}$  and the codeword index is perfectly transmitted. It means that the quantized version of the stochastic gradients can be perfectly transmitted. Then, by averaging all the quantized stochastic gradients, the approximated

global gradient can be computed as:

$$\text{(Approximated Global gradient)} \quad \widehat{\mathbf{g}}[n] = \frac{1}{K} \sum_{k=1}^K \widehat{\mathbf{g}}_k[n], \quad (6)$$

where  $\widehat{\mathbf{g}}[n]$  denotes an estimate of the global gradient at the  $n$ -th iteration. Then, the global model  $\boldsymbol{\theta}$  is updated as follows:

$$\text{(Model updating)} \quad \boldsymbol{\theta}[n+1] = \boldsymbol{\theta}[n] - \eta \widehat{\mathbf{g}}[n], \quad (7)$$

where  $\eta$  is the step size. The learning process involves the iteration from (4) to (7) until the model converges.

Let us define the normalized stochastic gradient  $\mathbf{f} = \frac{\mathbf{g}}{\|\mathbf{g}\|}$  and the norm of the stochastic gradient  $\rho = \|\mathbf{g}\|$ . For analytical tractability, we make the following assumption.

**Assumption 1.** *The normalized stochastic gradient  $\mathbf{f}$  is uniformly distributed on the Grassmann manifold.*

In other word, we assume that  $\mathbf{g}$  is isotropic (i.e., statistically invariant under unitary transformation). To support our assumption, we have run a hypothesis test on a real stochastic-gradient-dataset. Specifically, we apply the *Kolmogorov-Smirnov test* (KS-test), which is a nonparametric hypothesis test quantifying a distance between the empirical distribution function and the referenced one. Let the null hypothesis  $\mathcal{H}_0$  be such that  $\mathbf{f}$  is uniformly distributed, and  $\mathcal{H}_1$  the alternative one. The obtained *p-value*, which indicates whether to reject or accept the null hypothesis  $\mathcal{H}_0$ , is close to 0.1. Given that the commonly used threshold for *p-value* is 0.05, this result suggests that the null hypothesis  $\mathcal{H}_0$  will be accepted, i.e., the normalized stochastic gradient is believed to be uniformly distributed on the Grassmann manifold.

We are interested in the *mean squared error* (MSE) of the stochastic gradient quantization problem. For a given *quantization codebook*  $\mathcal{C}$  of  $B$  bits (i.e., containing  $2^B$  codewords), the optimal quantization function in the MSE sense is such that  $\mathbf{Q}_{\mathcal{C}}^E(\mathbf{g}) \in \arg \min_{\hat{\mathbf{g}} \in \mathcal{C}} \|\mathbf{g} - \hat{\mathbf{g}}\|^2$ . And the distortion of the quantizer is denoted by  $D_{\mathcal{C}}^E = \mathbb{E} \{ \|\mathbf{g} - \mathbf{Q}_{\mathcal{C}}^E(\mathbf{g})\|^2 \}$ . Here, the superscript ‘E’ stands for Euclidean distance. The optimal codebook is therefore

$$\mathcal{C}^* \in \arg \min_{\mathcal{C}} D_{\mathcal{C}}^E. \quad (8)$$

Two main challenges of the above optimal quantization problem are: 1) codebook optimization which is NP hard; 2) VQ which is also NP hard with respect to the dimension for a general

codebook. Therefore, the goal of this work is to propose a hierarchical codebook design which enables VQ with low complexity.

### III. HIERARCHICAL GRADIENT QUANTIZATION

In this paper, we propose to quantize the gradient norm  $\rho$  with a scalar codebook  $\mathcal{C}_\rho$  with  $B_\rho$  bits and the normalized stochastic gradient  $\mathbf{f}$  with a Grassmannian codebook<sup>1</sup>  $\mathcal{C}_\mathbf{f}$  with  $B_\mathbf{f}$  bits. This is motivated by the suitability of such a codebook for quantizing a vector that contains directional information and the tractability of relevant designs [22].

Nevertheless, directly designing the codebook for  $\mathbf{f}$  is impractical due to its high-dimensionality. To further reduce the complexity, we propose to decompose  $\mathbf{f}$  prior to quantization as follows. Assuming that  $\text{Dim} = LM$  for some integers<sup>2</sup>  $M$  and  $L$ , we partition the vector  $\mathbf{f}$  into  $M$  blocks of length  $L$ , i.e.,  $\mathbf{f}^T = [\mathbf{v}_1^T, \dots, \mathbf{v}_M^T]$ . We call  $\mathbf{v}_i$  the  $i$ -th *block gradient*. Also, let us define the *normalized block gradient*  $\mathbf{s}_i = \frac{\mathbf{v}_i}{\|\mathbf{v}_i\|}$  and the *hinge vector*  $\mathbf{h} = [h_1, \dots, h_M]^T$  where  $h_i = \|\mathbf{v}_i\|$ ,  $\forall i = 1, \dots, M$ . It follows that both the normalized block gradients and the hinge vector have unit norm and can be quantized with Grassmannian quantizers. In addition, one can show that the normalized block gradients are also isotropic, which is formally established in the following lemma.

**Lemma 1.** (*Uniformity of normalized block gradients*). *If  $\mathbf{f} = [f_1, f_2, \dots, f_{\text{Dim}}]^T$  is a uniformly distributed unitary random vector and given  $\mathbf{v} = [f_m, f_{m+1}, \dots, f_n]^T$  with  $m < n$  an arbitrary block gradient picked from  $\mathbf{f}$ , one can have that the normalized block gradient  $\mathbf{s} = \frac{\mathbf{v}}{\|\mathbf{v}\|}$  is uniformly distributed on the Grassmann manifold.*

*Proof:* See Appendix A. ■

Given the above decomposition and properties, we propose a quantization scheme with the following main ingredients, as shown in Fig. 2.

- For the stochastic gradient norm:  $B_\rho$ -bit scalar quantizer  $\mathcal{C}_\rho$ ;
- For the normalized block gradients:  $B_s$ -bit uniform and even<sup>3</sup> Grassmannian quantizer  $\mathcal{C}_s$ ;

<sup>1</sup>A Grassmannian codebook is a set of unit norm codewords.

<sup>2</sup>We can apply zero padding if  $\text{Dim} \neq LM$ .

<sup>3</sup>We call  $\mathcal{C}$  an even codebook if it can be partitioned as  $\mathcal{C} = \mathcal{C}^+ \cup \mathcal{C}^-$  with  $\mathcal{C}^+ \cap \mathcal{C}^- = \emptyset$  such that  $-\mathbf{c} \in \mathcal{C}^-, \forall \mathbf{c} \in \mathcal{C}^+$ .

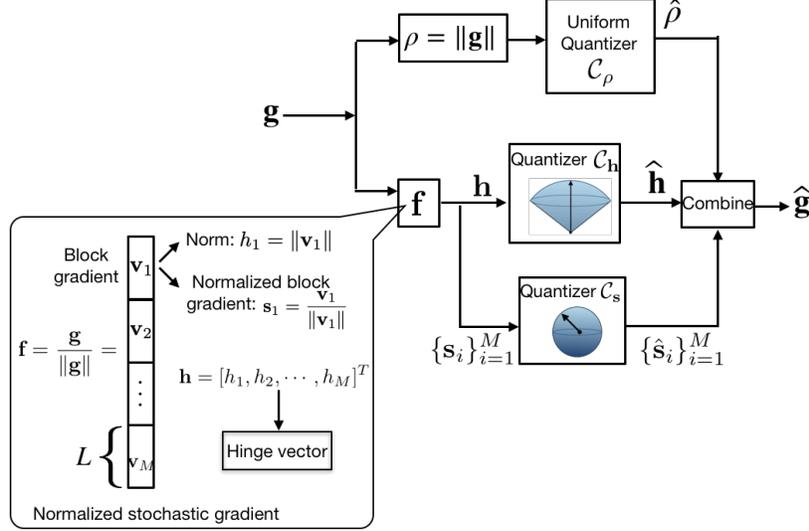


Figure 2. Hierarchical quantization scheme.

- For the hinge vector:  $B_h$ -bit positive<sup>4</sup> Grassmannian quantizer  $\mathcal{C}_h$ .

In total, we need  $B = B_\rho + MB_s + B_h$  bits. The quantized version of  $\mathbf{g}$  is

$$\hat{\mathbf{g}}^T = \hat{\rho}[\hat{\mathbf{f}}_1^T, \dots, \hat{\mathbf{f}}_M^T] = \hat{\rho}[\hat{h}_1\hat{\mathbf{s}}_1^T, \dots, \hat{h}_M\hat{\mathbf{s}}_M^T], \quad (9)$$

where  $\hat{\rho}$ ,  $\hat{h}_i$ , and  $\hat{\mathbf{s}}_i$  denote the quantized versions of  $\rho$ ,  $h_i$ , and  $\mathbf{s}_i, \forall i$ .

For the above quantizers, we focus on quantization functions that minimize the Euclidean distance (MSE) between  $\mathbf{g}$  and  $\hat{\mathbf{g}}$ . Let  $\mathbf{x}$  and  $\hat{\mathbf{x}}$  be two unit norm vectors, and let  $d_c(\mathbf{x}, \hat{\mathbf{x}}) = \sqrt{1 - |\hat{\mathbf{x}}^T \mathbf{x}|^2}$  be the chordal distance that measures angular deviation between  $\mathbf{x}$  and  $\hat{\mathbf{x}}$ . The two following lemmas are straightforward.

**Lemma 2.** *For a positive Grassmannian codebook  $\mathcal{C}$  and a given unitary vector  $\mathbf{x}$  with positive entries, if  $\hat{\mathbf{x}} \in \arg \min_{\hat{\mathbf{x}} \in \mathcal{C}} d_c(\mathbf{x}, \hat{\mathbf{x}})$ , then  $\hat{\mathbf{x}} \in \arg \min_{\hat{\mathbf{x}} \in \mathcal{C}} \|\mathbf{x} - \hat{\mathbf{x}}\|^2$ .*

**Lemma 3.** *For an even Grassmannian codebook  $\mathcal{C}$  and a given unit norm vector  $\mathbf{x}$ , if  $\hat{\mathbf{x}} \in \arg \min_{\hat{\mathbf{x}} \in \mathcal{C}^+} d_c(\mathbf{x}, \hat{\mathbf{x}})$ , then either  $\hat{\mathbf{x}}$  or  $-\hat{\mathbf{x}}$  belongs to  $\arg \min_{\hat{\mathbf{x}} \in \mathcal{C}} \|\mathbf{x} - \hat{\mathbf{x}}\|^2$ .*

Hence, to find a codeword in  $\mathcal{C}$  with the shortest Euclidean distance to  $\mathbf{x}$ , we first find a codeword  $\hat{\mathbf{x}}$  in  $\mathcal{C}^+$  with the shortest chordal distance to  $\mathbf{x}$ . If the inner product between the pair

<sup>4</sup>We call  $\mathcal{C}$  a positive codebook if all codewords have only positive entries.

is negative, then we flip the sign of the codeword, which is a codeword in  $\mathcal{C}^-$  and, therefore, is still inside  $\mathcal{C}$ . In addition, if  $\hat{\mathbf{x}}$  is the quantized version of  $\mathbf{x}$  in an even codebook, then

$$\frac{1}{2}\|\mathbf{x} - \hat{\mathbf{x}}\|^2 = 1 - \sqrt{1 - d_c^2(\mathbf{x}, \hat{\mathbf{x}})}. \quad (10)$$

**Proposition 1.** *The mean square error of the proposed quantizer is*

$$\mathbb{E} \{ \|\mathbf{g} - \hat{\mathbf{g}}\|^2 \} = D_{\mathcal{C}_\rho}^E + E_g D_{\mathcal{C}_f}^E, \quad (11)$$

where  $E_g = \mathbb{E} \{ \|\mathbf{g}\|^2 \}$ , and

$$D_{\mathcal{C}_f}^E = D_{\mathcal{C}_s}^E + D_{\mathcal{C}_h}^E - \frac{1}{2} D_{\mathcal{C}_s}^E D_{\mathcal{C}_h}^E, \quad (12)$$

with  $D_{\mathcal{C}_s}^E = \mathbb{E} \{ \|\mathbf{s} - \mathbf{Q}_{\mathcal{C}_s}^E(\mathbf{s})\|^2 \}$ ,  $D_{\mathcal{C}_h}^E = \mathbb{E} \{ \|\mathbf{h} - \mathbf{Q}_{\mathcal{C}_h}^E(\mathbf{h})\|^2 \}$ .

**Corollary 1.** *When  $\max\{D_{\mathcal{C}_s}^E, D_{\mathcal{C}_h}^E\}$  is small, we have*

$$\mathbb{E} \{ \|\mathbf{g} - \hat{\mathbf{g}}\|^2 \} \lesssim D_{\mathcal{C}_\rho}^E + E_g (D_{\mathcal{C}_s}^E + D_{\mathcal{C}_h}^E), \quad (13)$$

where  $\lesssim$  means the upper bound is asymptotically tight.

There are two sub-problems for the codebook design problem. First, for given bit allocation  $(B_\rho, B_s, B_h)$ , we need to jointly design the codebooks  $(\mathcal{C}_\rho, \mathcal{C}_s, \mathcal{C}_h)$ . From Corollary 1, it is asymptotically optimal, in the sense of scaling law, to design the codebooks  $\mathcal{C}_\rho, \mathcal{C}_s, \mathcal{C}_h$  separately. Second, we should optimize the bit allocation such that the MSE (13) is minimized. Given that it is hard to obtain tractable MSE expression as a function of the codebook size, we investigate the more tractable expected squared chordal distance instead. Such choice is justified by the closeness between the two measures in the following sense

$$d_c^2(\mathbf{x}, \hat{\mathbf{x}}) \leq \|\mathbf{x} - \hat{\mathbf{x}}\|^2 \leq 2d_c^2(\mathbf{x}, \hat{\mathbf{x}}). \quad (14)$$

Mathematically, we have the following optimization problem

$$\min_{B_\rho, B_s, B_h} \left\{ \min_{\mathcal{C}_\rho} D_{\mathcal{C}_\rho}^E + E_g \left( \min_{\mathcal{C}_s} D_{\mathcal{C}_s}^C + \min_{\mathcal{C}_h} D_{\mathcal{C}_h}^C \right) \right\}, \quad (15)$$

$$\text{s.t. } B_\rho + MB_s + B_h \leq B, \quad (16)$$

where  $D_{\mathcal{C}}^C = \mathbb{E} \{ d_c^2(\mathbf{x}, \mathbf{Q}_{\mathcal{C}}^C(\hat{\mathbf{x}})) \}$  with the superscript ‘C’ for ‘chordal’ distance.

1) *The design of Grassmannian codebook  $\mathcal{C}_s$* : Recall that the codebook  $\mathcal{C}_s$  is an even codebook. Hence, it is enough to first construct a codebook  $\mathcal{C}_s^+$  of size  $2^{B_s-1}$ , then construct  $\mathcal{C}_s^- = \{\mathbf{c} : -\mathbf{c} \in \mathcal{C}_s^+\}$ , and finally let  $\mathcal{C}_s = \mathcal{C}_s^+ \cup \mathcal{C}_s^-$ . The optimal construction of  $\mathcal{C}_s^+$  for isotropic sources is known as *Grassmannian line packing* [14], formulated as

$$(\text{codebook design for } \mathcal{C}_s^+) \quad \max_{\mathcal{C}_s^+} \min_{\hat{\mathbf{s}} \neq \hat{\mathbf{s}}' \in \mathcal{C}_s^+} d_c(\hat{\mathbf{s}}, \hat{\mathbf{s}}'). \quad (17)$$

2) *The design of Grassmannian codebook  $\mathcal{C}_h$* : An optimal codebook  $\mathcal{C}_h$  should satisfy

$$\mathcal{C}_h \in \arg \min_{\mathcal{C}_h} \mathbb{E}[d_c^2(\mathbf{h}, \mathbf{Q}_{\mathcal{C}_h}^C(\mathbf{h}))]. \quad (18)$$

Since the hinge vector is not isotropic, uniform quantization is not optimal in general. A practical suboptimal solution is the *Lloyd algorithm* [23] on the Grassmann manifold that can be implemented iteratively.

3) *The design of scalar codebook  $\mathcal{C}_\rho$* : Similarly, an optimal codebook  $\mathcal{C}_\rho$  should satisfy

$$\mathcal{C}_\rho \in \arg \min_{\mathcal{C}_\rho} \mathbb{E}[\|\rho - \mathbf{Q}_{\mathcal{C}_\rho}^E(\rho)\|^2]. \quad (19)$$

For simplicity, we adopt a uniform quantizer for  $\rho = \|\mathbf{g}\|$ . So far, we have developed the hierarchical quantization scheme for stochastic gradients using three codebooks (see Fig. 2).

#### IV. DISTORTION ANALYSIS

In this section, the distortion from quantizing the stochastic gradients under the proposed hierarchical scheme will be analyzed, which involves the derivations of  $D_{\mathcal{C}_s}^C$ ,  $D_{\mathcal{C}_h}^C$  and  $D_{\mathcal{C}_\rho}^E$  in (15), respectively.

1) *Distortion analysis on the normalized block gradient*: According to [24], [25], the codebook designed by line packing is asymptotically optimal, and thus the resulting distortion for quantizing uniformly distributed unitary random vectors is asymptotically identical to that using the random codebook. Mathematically, the average distortion for the normalized block gradient is characterized as follows.

**Lemma 4.** ( [25], Theorem 2). *Let  $\mathcal{C}_s$  be a codebook designed by line packing, with resolution  $B_s - 1$  and dimensionality  $L$ . The average distortion, denoted as  $D_{\mathcal{C}_s}^C$ , incurred by quantizing the uniformly distributed unitary random vector under the codebook  $\mathcal{C}_s$  can be bounded as*

$$\frac{L-1}{L+1} 2^{-\frac{2(B_s-1)}{L-1}} + o(1) \leq D_{\mathcal{C}_s}^C \leq 2^{-\frac{2(B_s-1)}{L-1}} + o(1), \quad (20)$$

where  $o(1)$  indicates a vanishing function of  $L$  as  $L \rightarrow \infty$ .

One can observe from the above lemma that  $D_{\mathcal{C}_s}^{\mathcal{C}}$  decreases exponentially as the codebook resolution  $B_s$  increases. Moreover, fix  $B_s$ , as the length of the block increases, the quantization performance degrades accordingly in that pairwise distances between codewords enlarges according to the well-known results in line packing [24].

2) *Geometric properties of the hinge vector*: To facilitate the derivation of  $D_{\mathcal{C}_h}^{\mathcal{C}}$ , the geometric properties of the hinge vector will be analyzed in this sub-section.

To begin with, we investigate the statistic distribution of the hinge vector by introducing the following lemma.

**Lemma 5.** *If  $X \sim \mathcal{X}^2(a)$  and  $Y \sim \mathcal{X}^2(b)$  are independent Chi-squared random variables, then  $\frac{X}{X+Y}$  follows the Beta  $(\frac{a}{2}, \frac{b}{2})$  distribution.*

Since the normalized stochastic gradient  $\mathbf{f}$  is isotropic, thereby, it can be further generated as  $\mathbf{f} = \frac{\mathbf{x}}{\|\mathbf{x}\|}$ , where the elements of  $\mathbf{x}$  are i.i.d. Gaussian random variables with zero mean and unit variance. It follows that the square of each element of the hinge vector  $\mathbf{h} = [h_1, h_2, \dots, h_M]^T$  are beta distributed, i.e.,

$$h_i^2 = \frac{Z_1}{Z_1 + Z_2} \sim \text{Beta}\left(\frac{L}{2}, \frac{\text{Dim} - L}{2}\right), \quad \forall i \in [1, M], \quad (21)$$

where  $Z_1 \sim \mathcal{X}^2(L)$  and  $Z_2 \sim \mathcal{X}^2(\text{Dim} - L)$  are independent.

With this distribution at hand, the geometric properties of the hinge vector will be analyzed. One specific result is characterized as below.

**Proposition 2** (Convergence of hinge vector). *For sufficiently large  $L$ ,  $\mathbf{h}$  converges to a constant vector  $\mathbf{h}_{\text{ref}} = \frac{1}{\sqrt{M}}\mathbf{1}_{M \times 1}$ . In particular, given  $r = \sqrt{\frac{2}{1+2L^{\frac{1}{4}}}}$ , we have*

$$\Pr(d_c(\mathbf{h}_{\text{ref}}, \mathbf{h}) > r) < L^{-\frac{1}{2}} + O(L^{-\frac{3}{2}}), \quad (22)$$

*Proof: The convergence is implied by (22) that is proved in Appendix B. ■*

**Remark 1.** (Geometric Interpretation). According to Proposition 2, the hinge vector converges to  $\mathbf{h}_{\text{ref}}$  at least geometrically fast in block dimensionality  $L$ . To be specific, the hinge vector locates with high probability within a ball of radius  $r$  (with respect to the chordal distance) on

the Grassmann manifold centered at  $\mathbf{h}_{\text{ref}}$ , namely,  $\mathcal{B}_{\mathbf{h}_{\text{ref}}}(r) = \{\mathbf{h} \mid \mathbf{h}^T \mathbf{h} = 1, d_c(\mathbf{h}_{\text{ref}}, \mathbf{h}) \leq r\}$ , where  $r$  converges to zero as  $L$  grows.

3) *Distortion analysis on the hinge vector:* Based on the above analyses, we are now ready to derive the average distortion for quantizing the hinge vector. To this end, we first introduce the following lemma.

**Lemma 6.** *If the set  $\mathcal{C}' = \mathcal{C} \cap \mathcal{B}_{\mathbf{h}_{\text{ref}}}((1 + \alpha)r)$  is not empty for some  $\alpha > 0$ , then, for any  $\mathbf{x} \in \mathcal{B}_{\mathbf{h}_{\text{ref}}}(r)$ , we have*

$$\min_{\mathbf{c} \in \mathcal{C}'} d_c(\mathbf{c}, \mathbf{x}) \leq \left(1 + \frac{2}{\alpha}\right) \min_{\mathbf{c} \in \mathcal{C}} d_c(\mathbf{c}, \mathbf{x}) \quad (23)$$

*Proof:* Let  $\mathbf{c} \in \mathcal{C}'$  and  $\mathbf{c}' \in \mathcal{C} \setminus \mathcal{C}'$ . Thus, we have  $d_c(\mathbf{c}, \mathbf{h}_{\text{ref}}) \leq (1 + \alpha)r$  and  $d_c(\mathbf{c}', \mathbf{h}_{\text{ref}}) \geq (1 + \alpha)r$ . By the triangle inequality, we have

$$d_c(\mathbf{c}, \mathbf{x}) \leq d_c(\mathbf{h}_{\text{ref}}, \mathbf{x}) + d_c(\mathbf{c}, \mathbf{h}_{\text{ref}}) \leq (2 + \alpha)r, \quad (24)$$

$$d_c(\mathbf{c}', \mathbf{x}) \geq d_c(\mathbf{c}', \mathbf{h}_{\text{ref}}) - d_c(\mathbf{x}, \mathbf{h}_{\text{ref}}) \geq \alpha r, \quad (25)$$

from which we have

$$d_c(\mathbf{c}, \mathbf{x}) \leq \left(1 + \frac{2}{\alpha}\right) d_c(\mathbf{c}', \mathbf{x}). \quad (26)$$

Taking the minimum on both sides, (23) is straightforward. ■

Let us now construct a codebook of  $N = 2^{B_h}$  codewords as follows. First, we draw  $N' \geq N$  points uniformly from the Grassmann manifold as for the normalized block gradient codebook. Then, we choose the  $N$  codewords that are closest to  $\mathbf{h}_{\text{ref}}$  to form the codebook  $\mathcal{C}_h$ . To analyze the average distortion, we introduce two balls  $\mathcal{B}_{\mathbf{h}_{\text{ref}}}(r)$  and  $\mathcal{B}_{\mathbf{h}_{\text{ref}}}((1 + \alpha)r)$  for some  $\alpha > 0$  and  $r > 0$  that can be optimized later on. Let us consider the following encoding rules.

- If  $\mathbf{h}$  lies outside of  $\mathcal{B}_{\mathbf{h}_{\text{ref}}}(r)$ , an encoding error is declared. This event has probability  $P_1(r)$ .
- If no codeword lies inside  $\mathcal{B}_{\mathbf{h}_{\text{ref}}}((1 + \alpha)r)$ , an encoding error is declared. This event has probability  $P_2(\alpha, r)$ .
- If there are more than  $N$  codewords inside  $\mathcal{B}_{\mathbf{h}_{\text{ref}}}((1 + \alpha)r)$ , an error is declared. This event has probability  $P_3(\alpha, r)$ .

We can upper-bound the distortion by 1 whenever an error is declared, then we have the following

upper bound on the average distortion from the union bound:

$$P_1(r) + P_2(\alpha, r) + P_3(\alpha, r) + (1 - P_1(r)) \left(1 + \frac{2}{\alpha}\right) D(N'), \quad (27)$$

where the last term is from Lemma 6 and  $D(N')$  is the average distortion for a uniform random quantizer with  $N'$  codewords.

In particular,  $P_3(\alpha, r)$  is the complementary cumulative distribution function of a binomial distribution with parameter  $N'$  and  $p$  where  $p$  is the probability that a uniformly distributed point falls inside the ball.

**Lemma 7.** *The distortion for quantizing the hinge vector can be upper-bounded as*

$$D_{C_h}^C \leq L^{-\frac{1}{2}} (\beta_L 2^{-\frac{2B_h}{M-1}} + 1) + O(L^{-\frac{3}{2}}), \quad (28)$$

where  $\beta_L = \frac{1+r^{\frac{1}{2}}}{1-r^{\frac{1}{2}}}$  with  $r = \sqrt{\frac{2}{1+2L^{\frac{1}{4}}}} \approx L^{-\frac{1}{8}}$ .

*Proof:* See Appendix C. ■

Several observations can be made from the above lemma. First, as the codebook resolution  $B_h$  increases, the upper bound of the quantization error reduces accordingly due to the reduction of pairwise codewords distance. Second, given  $B_h$ , the upper bound is a decreasing function of the block length  $L$ . This is because the hinge vector tends to converge to  $\mathbf{h}_{\text{ref}}$  given larger  $L$  and the resulting quantization error reduces.

4) *Distortion analysis on the stochastic-gradient norm:* given a uniform quantizer for the norm of stochastic gradients, the average distortion can be upper-bounded as

$$D_{C_\rho}^E \leq \left(\frac{\Delta^{\text{quant}}}{2}\right)^2, \quad (29)$$

where  $\Delta^{\text{quant}}$  denotes the quantization interval of a uniform quantizer.

## V. QUANTIZATION BIT ALLOCATION

In this section, a practical *bit-allocation scheme* will be developed. Specifically, given a fixed number  $B$  of bits, for quantizing the stochastic gradient vector, we aim to determine the scheme on how to allocate these bits to the three codebooks derived in the preceding section.

$$(\mathbf{P}') \quad \min_{B_\rho, B_s, B_h} \left\{ \frac{\rho_{\max}^2}{4} 2^{-2B_\rho} + E_g 2^{-\frac{2(B_s-1)}{L-1}} + E_g \beta_L 2^{-\frac{2B_h}{M-1}} L^{-\frac{1}{2}} \right\}, \quad (32)$$

$$\text{s.t. } B_\rho + MB_s + B_h = B, \quad (33)$$

Given the proposed hierarchical quantization scheme, the original bit-allocation problem is formulated as

$$\min_{B_\rho, B_s, B_h} \left\{ D_{C_\rho}^E + E_g (D_{C_s}^C + D_{C_h}^C) \right\}, \quad (30)$$

$$\text{s.t. } B_\rho + MB_s + B_h = B. \quad (31)$$

Here, the distortions are replaced by the upper bounds derived previously.

For the *Bit-allocation problem*, we assume for tractability that the elements of stochastic gradient  $\mathbf{g}$  are i.i.d. Gaussian with zero mean and unit variance. Note that this is also the worst case in the sense that for a given variance the differential entropy is maximized with Gaussian distributions. Then, it follows that  $E_g = ML$  and  $\rho = \|\mathbf{g}\| \sim \mathcal{X}(ML)$ . Given  $\rho \in [0, \rho_{\max}]^5$ , it follows from (29) that  $D_{C_\rho}^E \leq \frac{\rho_{\max}^2}{(2^{B_\rho+1}+2)^2} \leq \frac{\rho_{\max}^2}{4} 2^{-2B_\rho}$ . Then, the original bit-allocation problem can be relaxed as  $(\mathbf{P}')$ , which is given at the top of this page.

The above problem  $(\mathbf{P}')$  is convex, and the optimal solutions can be derived by leveraging *Karush-Kuhn-Tucker* (KKT) conditions as follows

$$(\text{KKT conditions}) \quad \begin{cases} \lambda^* + \frac{\partial f(B_\rho^*, B_s^*, B_h^*)}{\partial B_\rho^*} = 0 \\ M\lambda^* + \frac{\partial f(B_\rho^*, B_s^*, B_h^*)}{\partial B_s^*} = 0 \\ \lambda^* + \frac{\partial f(B_\rho^*, B_s^*, B_h^*)}{\partial B_h^*} = 0, \end{cases} \quad (34)$$

where  $f(B_\rho^*, B_s^*, B_h^*)$  is the objective of the optimization problem, i.e. (32), and  $\lambda^*$  is the Lagrange multiplier. Solving the above equations, we obtain the following bit-allocation scheme.

<sup>5</sup>Due to the fact that for  $\rho \sim \mathcal{X}(ML)$ ,  $\mathbb{E}[\rho] + \sqrt{\text{Var}[\rho]} \leq \mathbb{E}[\rho^2] + \sqrt{\text{Var}[\rho^2]}$  with  $\sqrt{\text{Var}[\cdot]}$  denoting the standard deviation, for tractability, we take  $\rho_{\max} = \mathbb{E}[\rho^2] + \sqrt{\text{Var}[\rho^2]} = ML + \sqrt{2ML}$ .

**Scheme 1.** (Quantization Bit Allocation). To minimize the distortion, the bits  $B$  can be allocated to the three quantizers as follows:

$$B_\rho^* = \lfloor \log_2 \frac{ML + \sqrt{2ML}}{2} + \frac{1}{2} \log_2 \ln 2 + \frac{1}{2} - \frac{1}{2} \log_2 \lambda^* \rfloor, \quad (35)$$

$$B_s^* = \lfloor \frac{L-1}{2} \log_2 \frac{2L}{L-1} + 1 + \frac{L-1}{2} \log_2 \ln 2 - \frac{L-1}{2} \log_2 \lambda^* \rfloor, \quad (36)$$

$$B_h^* = \lfloor \frac{M-1}{2} \log_2 \frac{2}{M-1} + \frac{M-1}{2} \log_2 \beta_L M \sqrt{L} + \frac{M-1}{2} \log_2 \ln 2 - \frac{M-1}{2} \log_2 \lambda^* \rfloor, \quad (37)$$

where  $\beta_L = \frac{1+L^{-\frac{1}{16}}}{1-L^{-\frac{1}{16}}}$  and  $\lambda^*$  can be obtained by substituting the above equations into (33) as

$$\log_2 \lambda^* = \frac{2}{ML} \log_2 \frac{ML + \sqrt{2ML}}{2} + \frac{L-1}{L} \log_2 \frac{2L}{L-1} + \frac{2}{L} + \frac{1}{ML} + \frac{2(M-1)}{ML} \log_2 \frac{2}{M-1} + \frac{M-1}{ML} \log_2 \beta_L M \sqrt{L} + \log_2 \ln 2 - \frac{2B}{ML}. \quad (38)$$

Next, it is necessary to show that the above bit-allocation scheme is optimal in the sense of scaling law. To this end, we bound  $\mathbb{E}[\|\mathbf{g} - \hat{\mathbf{g}}\|^2]$  in the following theorem.

**Theorem 1.** (Optimality of Bit Allocation). *For sufficiently large  $L$  and at the low resolution regime<sup>6</sup>, i.e.  $B \leq ML$ , the distortion under MSE metric per dimension, i.e.  $\frac{\mathbb{E}[\|\mathbf{g} - \hat{\mathbf{g}}\|^2]}{ML}$  can be bounded as*

$$-\frac{2 \ln 2}{ML} B \leq \ln \frac{\mathbb{E}[\|\mathbf{g} - \hat{\mathbf{g}}\|^2]}{ML} \leq c_{\text{gap}} - \frac{2 \ln 2}{ML} B + O\left(L^{-\frac{3}{2}}\right), \quad (39)$$

where  $c_{\text{gap}} = \ln 2 - \ln \frac{2L}{L-1} + \frac{2}{ML} \ln \frac{ML + \sqrt{2ML}}{2} + \frac{L-1}{L} \ln \frac{2L}{L-1} + \frac{2}{L} \ln 2 + \frac{2(M-1)}{ML} \ln \frac{2}{M-1} + \frac{\ln 2}{ML} + \frac{M-1}{ML} \ln \beta_L M \sqrt{L} + 2(\beta_L + 1)L^{-\frac{1}{2}} - 2(\beta_L + 1)^2 L^{-1}$  with  $\beta_L = \frac{1+L^{-\frac{1}{16}}}{1-L^{-\frac{1}{16}}}$ ;  $B$  denotes the number of bits used for quantizing the stochastic gradient  $\mathbf{g}$ .

It can be observed from (39) that the scaling law of the upper bound is the same as that of the lower bound with respect to  $B$ . It is further noted that the upper bound in the above theorem is derived by setting  $B_h^* = 0$ . It means that the derived scaling law is independent of

<sup>6</sup>The term ‘low resolution’ is declared in the sense that only fewer than one bit is exploited for quantizing each coefficient of the stochastic gradient. This is a popular regime being explored in the area of edge learning [7].

the number of bits allocated for quantizing the hinge vector. This is aligned with the intuition that as block length  $L$  increases, the hinge vector tends to converge to the constant vector and the corresponding distortion is close to zero. From more theoretic point of view, it follows from (28) that, at the low resolution regime, no matter how many bits are allocated to  $\mathcal{C}_h$ , the decaying-rate of the average distortion is asymptotically bounded by  $L^{-\frac{1}{2}}$ . Thereby, this motivates a practical bit-allocation scheme as given below.

**Scheme 2.** (Practical Bit Allocation). For the high-dimensional stochastic gradient  $\mathbf{g}$ ,  $B_\rho^* = \left\lceil \log_2 \frac{ML + \sqrt{2ML}}{2} + \frac{1}{2} \log_2 \ln 2 + \frac{1}{2} - \frac{1}{2} \log_2 \lambda^* \right\rceil$  bits are allocated for quantizing its norm with  $\lambda^*$  is defined in (38). All rest bits should be allocated to the codebook  $\mathcal{C}_s$  while exploiting  $\mathbf{h}_{\text{ref}} = \frac{1}{\sqrt{M}} \mathbf{1}_{M \times 1}$  as a surrogate for the hinge vector.

**Remark 2.** The above practical bit-allocation scheme makes the proposed hierarchical quantization scheme be of low-complexity. To be specific, the relative low-dimensional block gradients makes the design complexity of codebook  $\mathcal{C}_s$  via line packing algorithm reduces significantly compared to quantizing the high-dimensional stochastic gradient as a whole. On the other hand, the design complexity is further reduced without constructing the codebook  $\mathcal{C}_h$  for the hinge vector.

## VI. LEARNING CONVERGENCE RATE ANALYSIS

Given a typical quantization scheme for the stochastic gradient, one concern related is that whether it will lead to the convergence of the learning algorithm. Thereby, in this section, the convergence rate of the learning algorithm under the proposed hierarchical quantization scheme will be theoretically investigated.

We begin our analysis in the non-convex setting, where we follow the standard assumptions of the stochastic optimization literature (see e.g., [7]). The specific assumptions are given as follows.

**Assumption 2.** (Lower Bound). *For all  $\boldsymbol{\theta}$  and some constant  $F^*$ , we have that the global objective value  $F(\boldsymbol{\theta}) \geq F^*$ .*

**Assumption 3.** (Smoothness). Let  $\bar{\mathbf{g}}(\boldsymbol{\theta})$  denote the gradient of the global objective  $F(\boldsymbol{\theta})$  evaluated at point  $\boldsymbol{\theta} = [\theta_1, \theta_2, \dots, \theta_{\text{Dim}}]^T$  with  $\text{Dim} = ML$ . Then  $\forall \boldsymbol{\theta}$  and  $\boldsymbol{\beta} = [\beta_1, \beta_2, \dots, \beta_{\text{Dim}}]^T$ , we require that for some non-negative constant vector  $\mathbf{l} = [l_1, l_2, \dots, l_{\text{Dim}}]^T$

$$|F(\boldsymbol{\beta}) - [F(\boldsymbol{\theta}) + \bar{\mathbf{g}}(\boldsymbol{\theta})^T (\boldsymbol{\beta} - \boldsymbol{\theta})]| \leq \frac{1}{2} \sum_{i=1}^{\text{Dim}} l_i (\beta_i - \theta_i). \quad (40)$$

**Assumption 4.** (Variance Bound). The stochastic gradient  $\mathbf{g}(\boldsymbol{\theta})$  is unbiased that has coordinate bounded variance:

$$\mathbb{E}[\mathbf{g}(\boldsymbol{\theta})] = \bar{\mathbf{g}}(\boldsymbol{\theta}) \quad \text{and} \quad \mathbb{E}[(\mathbf{g}(\boldsymbol{\theta})_i - \bar{\mathbf{g}}(\boldsymbol{\theta})_i)^2] \leq \sigma_i^2, \quad (41)$$

for a vector of non-negative constants  $\boldsymbol{\sigma} = [\sigma_1, \sigma_2, \dots, \sigma_{\text{Dim}}]^T$ .

Under the above three standard assumptions, we have the following result.

**Theorem 2.** (Learning Convergence Rate). Let  $N$  be the number of iterations for the federated learning algorithm,  $K$  the total number of users, and  $\eta = \frac{1}{\sqrt{l_0 N}}$  the learning rate. It follows that

$$\mathbb{E} \left[ \frac{1}{N} \sum_{n=0}^{N-1} \|\bar{\mathbf{g}}_n\|^2 \right] \leq \frac{\sqrt{l_0} \left( \frac{1}{2K} \mathbb{E} [\|\mathbf{g} - \hat{\mathbf{g}}\|^2] + \frac{\|\boldsymbol{\sigma}\|^2}{2K} + F_0 - F^* \right)}{\sqrt{N} - \frac{\sqrt{l_0}}{2K}}, \quad (42)$$

where  $F_0$  is the initial objective value and  $F^*$  is defined in Assumption 2;  $l_0 = \|\mathbf{l}\|_\infty$  with  $\mathbf{l}$  defined in Assumption 4.

*Proof:* See Appendix E. ■

Several observations can be made from (42) as follows. First, the increment of the total iteration number  $N$  leads to the convergence of the learning algorithm. Specifically, as the number of users  $K \rightarrow \infty$ , the convergence rate is asymptotically  $O\left(\frac{1}{\sqrt{N}}\right)$ . Furthermore, as the number of users  $K$  increases, the upper bound in (42) decreases. This is because that the participation of more users, called multi-user gain, makes the aggregated-and-averaged stochastic gradient closer to the true gradient, leading to a faster convergence speed.

## VII. SIMULATION RESULTS

Consider a FEEL system with one edge server and  $K = 100$  edge devices. The simulation settings are given as follows unless specified otherwise. We consider the learning task of handwritten-digit recognition using the well-known MNIST dataset that consists of 10 categories

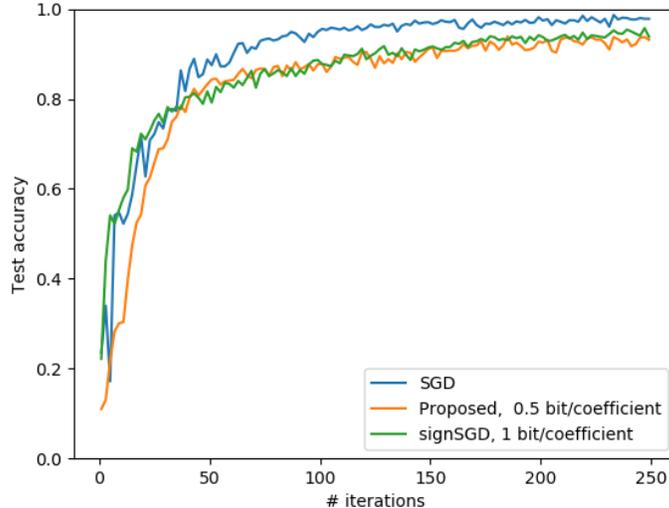


Figure 3. Performance comparison of *signSGD* and the proposed scheme.

ranging from digit “0” to “9” and a total of 60000 labeled training data samples. The classifier model is implemented using a 6-layer *convolutional neural network* (CNN) that consists of two  $5 \times 5$  convolution layers with ReLu activation (the first with 32 channels, the second with 64). Each followed with a  $2 \times 2$  max pooling, a fully connected layer with 512 units, ReLu activation, and a final softmax output layer. Furthermore, it is noted that the total number of bits used for quantizing each coefficient of the stochastic gradients is  $\frac{B_s}{L} + \frac{B_p}{\text{Dim}}$  given  $B_h = 0$  in the proposed bit-allocation scheme. Due to the fact that  $\frac{B_p}{\text{Dim}} = 0, \text{Dim} \rightarrow \infty$ , we define the number of bits per coefficient as  $\frac{B_s}{L}$  without loss of generality.

#### A. Performance of the Hierarchical Quantization Scheme

The effectiveness of the proposed hierarchical quantization scheme is evaluated by benchmarking against *signSGD* and SGD. The curves of the test accuracy versus the number iterations are illustrated in Fig. 3. Several observations can be made as follows. First, using fewer bits, i.e. 0.5 bit/coefficient, the performance of the proposed scheme is comparable to state-of-the-art *signSGD*, which uses 1 bit/coefficient. This attributes to the superiority of vector quantization over the scalar counterpart given the same bits at the low resolution regime. Furthermore, it can also be observed that SGD outperforms both quantization schemes because there exists quantization loss for both quantization schemes.

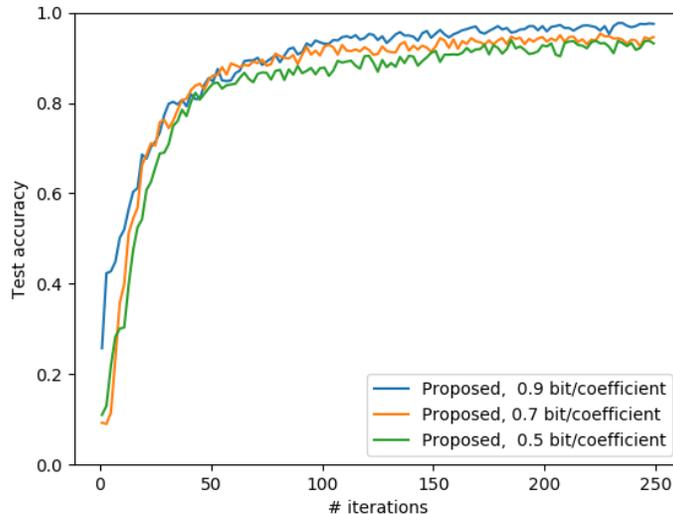


Figure 4. Effect of the codebook resolution  $B_s$  with block length  $L = 10$  and  $B_p = 26$  bits.

### B. Effect of the Codebook Resolution

Given the block length  $L$ , the effect of the codebook resolution  $B_s$  for  $\mathcal{C}_s$  is evaluated, where the resolution for  $\mathcal{C}_p$  is fixed. The curves of test accuracy versus the number of iterations by varying the codebook resolution are illustrated in Fig. 4. It can be observed that as the codebook resolution increases, the learning performance improves accordingly. This is because the increment of resolution reduces the pairwise chordal distance among codewords. Then, the resulting quantization error reduces, giving rise to a better learning performance.

### C. Effect of the Block Length

Given the fixed number of bits allocated to each block, the effects of the block length  $L$  is evaluated. In particular, the curves of test accuracy versus the number of iterations by varying the block length  $L$  are illustrated in Fig. 5. It can be observed that as  $L$  increases, the learning performance degrades accordingly in that the quantization error for the stochastic gradients enlarges. The underlying reason is that a larger  $L$  implies that the Grassmannian codebook is generated by the packing algorithm on a higher dimensional Grassmann manifold. This enlarges the pairwise chordal distance between codewords and thus the resulting quantization error is large.

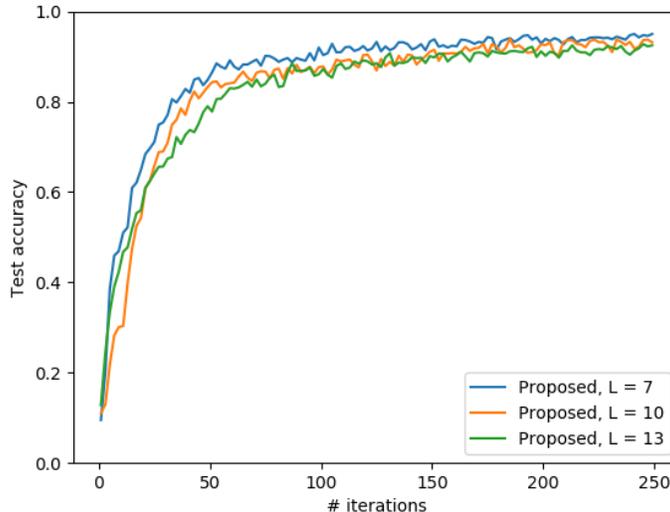


Figure 5. Effect of the block length  $L$  with  $B_s = 5$  bits and  $B_p = 26$  bits.

#### D. Effect of the Edge-Device Number

Fix the iteration number as 50, the relationship between the learning performance and the total number of edge-devices  $K$  with various block length  $L$  is illustrated in Fig. 6. It can be observed that as  $K$  increases, the learning performance improves accordingly. This is consistent with the result derived in Theorem 2. Specifically, as indicated by (42), a larger  $K$  reduces the noise variance, and also makes the aggregated-and averaged stochastic gradient closer to the true gradient, giving rise to a faster convergence speed.

### VIII. CONCLUDING REMARKS

In the context of FEEL, by investigating the statistic distribution of the normalized stochastic gradient, we propose a novel vector quantization scheme for high-dimensional stochastic gradients. This quantization scheme is of low-complexity, communication-efficient, and convergence-warranted. This work represents the first attempt to quantize high-dimensional stochastic gradients using efficient Grassmannian quantization, which is shown to be more communication-efficient than its state-of-the-art scalar counterpart. In the future, this work can be generalized into applying vector quantization to the accumulated quantization error, which is used for accelerating learning. Moreover, the vector quantization scheme can be further developed by taking the sparsity property and temporal correlation of stochastic gradients into consideration.

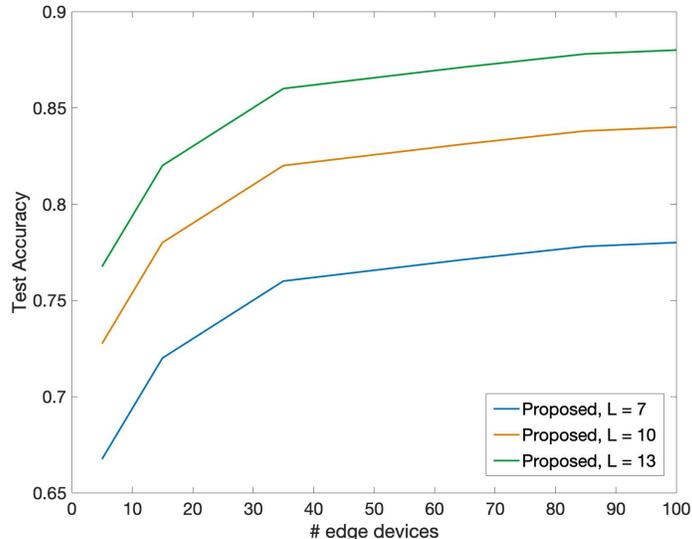


Figure 6. Effect of the edge-device number  $K$  with  $B_s = 5$  bits and  $B_\rho = 26$  bits.

## APPENDIX

### A. Proof of Lemma 1

The normalized stochastic gradient can be written as  $\mathbf{f} = \frac{\mathbf{x}}{\|\mathbf{x}\|}$  due to its uniformity on the Grassmann manifold, where elements of  $\mathbf{x}$  are i.i.d. Gaussian distributed with zero mean and unit variance. Thereby, an arbitrary block gradient can be written as  $\mathbf{v} = \left[ \frac{x_m}{\|\mathbf{x}\|}, \frac{x_{m+1}}{\|\mathbf{x}\|}, \dots, \frac{x_n}{\|\mathbf{x}\|} \right]^T$  with its norm being  $\|\mathbf{v}\| = \frac{\sqrt{\sum_{i=m}^n x_i^2}}{\|\mathbf{x}\|}$ . Then, it follows that

$$\mathbf{s} = \frac{\mathbf{v}}{\|\mathbf{v}\|} = \left[ \frac{x_m}{\sqrt{\sum_{i=m}^n x_i^2}}, \frac{x_{m+1}}{\sqrt{\sum_{i=m}^n x_i^2}}, \dots, \frac{x_n}{\sqrt{\sum_{i=m}^n x_i^2}} \right]^T. \quad (43)$$

This implies that  $\mathbf{s} = \frac{\mathbf{v}}{\|\mathbf{v}\|}$  is uniformly distributed on the Grassmann manifold.

### B. Proof of Proposition 2

To begin with, by applying the *law of large numbers*, it is easy to show that the hinge vector will converge to the constant vector  $\mathbf{h}_{\text{ref}} = \frac{1}{\sqrt{M}} \mathbf{1}_{M \times 1}$ , as  $L \rightarrow \infty$ . Next, we focus on calculating

the convergence rate. According to the definition of the chordal distance, we have

$$\begin{aligned}
P_1(r) &= \Pr(d_c(\mathbf{h}_{\text{ref}}, \mathbf{h}) > r) \\
&= \Pr\left(\sum_{i=1}^M \frac{h_i}{\sqrt{M}} < \sqrt{1-r^2}\right) \\
&= \Pr\left(\frac{1}{\sqrt{M}} \sum_{i=1}^M \sqrt{\frac{z_i}{\sum_{j=1}^M z_j}} < b\right), \tag{44}
\end{aligned}$$

where  $b^2 = 1 - r^2 = \frac{1-\epsilon_L}{1+\epsilon_L}$  with  $\epsilon_L = \frac{r^2}{2-r^2}$  and  $h_i = \sqrt{\frac{z_i}{\sum_{j=1}^M z_j}}$  with  $z_i, z_j \sim \mathcal{X}^2(L), \forall i, j$ . Since  $h_i$  and  $\sum_{j=1}^M z_j$  are independent, we have

$$\Pr(d_c(\mathbf{h}_{\text{ref}}, \mathbf{h}) > r) \cdot \Pr\left(\frac{1}{LM} \sum_{j=1}^M z_j < 1 + \epsilon_L\right) = \Pr\left(d_c(\mathbf{h}_{\text{ref}}, \mathbf{h}) > r, \frac{1}{LM} \sum_{j=1}^M z_j < 1 + \epsilon_L\right). \tag{45}$$

Since

$$\Pr\left(d_c(\mathbf{h}_{\text{ref}}, \mathbf{h}) > r, \frac{1}{LM} \sum_{j=1}^M z_j < 1 + \epsilon_L\right) \leq \Pr\left(\frac{1}{\sqrt{M}} \sum_{i=1}^M \sqrt{z_i} < b\sqrt{LM(1+\epsilon_L)}\right), \tag{46}$$

one upper bound can be derived as follows

$$\Pr(d_c(\mathbf{h}_{\text{ref}}, \mathbf{h}) > r) \leq \frac{\Pr\left(\frac{1}{\sqrt{M}} \sum_{i=1}^M \sqrt{z_i} < b\sqrt{LM(1+\epsilon_L)}\right)}{\Pr\left(\frac{1}{LM} \sum_{j=1}^M z_j < 1 + \epsilon_L\right)}. \tag{47}$$

In the following, we calculate the numerator and the denominator, respectively. First, we derive an upper bound of the numerator. Define  $y_z = \frac{1}{M} \sum_{i=1}^M \sqrt{z_i}$ , it follows that

$$\begin{aligned}
\Pr\left(\frac{1}{\sqrt{M}} \sum_{i=1}^M \sqrt{z_i} < b\sqrt{LM(1+\epsilon_L)}\right) &= \Pr\left(y_z < \sqrt{L(1-\epsilon_L)}\right) \\
&= \Pr\left(y_z - \mu_Y < -\left(\mu_Y - \sqrt{L(1-\epsilon_L)}\right)\right) \\
&\leq \Pr\left((y_z - \mu_Y)^2 > \left(\mu_Y - \sqrt{L(1-\epsilon_L)}\right)^2\right), \tag{48}
\end{aligned}$$

where the first equality holds given  $b^2 = \frac{1-\epsilon_L}{1+\epsilon_L}$ ;  $\mu_Y = \mathbb{E}[y_z] = \mathbb{E}\left[\frac{1}{M} \sum_{i=1}^M \sqrt{z_i}\right]$ . Then, by *Chebyshev's inequality*, it follows that

$$\Pr\left(\frac{1}{\sqrt{M}} \sum_{i=1}^M \sqrt{z_i} < b\sqrt{LM(1+\epsilon_L)}\right) \leq \frac{\sigma_Y^2}{\left(\mu_Y - \sqrt{L(1-\epsilon_L)}\right)^2}, \tag{49}$$

where  $\sigma_Y^2 = \mathbb{E}[y_z^2] - \mathbb{E}^2[y_z]$ . In order to derive the closed-form solution for the above bound, it suffices to calculate  $\mu_Y$  and  $\sigma_Y^2$ . Given  $z_i \sim \mathcal{X}^2(L)$ , one can have

$$\mu_Y = \mathbb{E}[\sqrt{z_i}] = \frac{\sqrt{2}\Gamma(\frac{L+1}{2})}{\Gamma(\frac{L}{2})}. \quad (50)$$

Then, by *Stirling's approximation* for gamma function, i.e.  $\Gamma(x) \approx \sqrt{2\pi}x^{x-\frac{1}{2}}e^{-x}$ , as  $x \rightarrow \infty$ , one can further have that

$$\mu_Y = \sqrt{L} \cdot e^{\frac{L}{2} \ln(1+\frac{1}{L}) - \frac{1}{2}} \stackrel{(a)}{\geq} \sqrt{L} \cdot e^{-\frac{1}{4L}} \stackrel{(b)}{\geq} \sqrt{L} \left(1 - \frac{1}{4L}\right), \quad (51)$$

where (a) follows from the fact that  $\ln(1 + \frac{1}{L}) \geq \frac{1}{L} - \frac{1}{2L^2}$  and (b) follows from the fact that  $e^{-\frac{1}{4L}} \geq 1 - \frac{1}{4L}$ . Next, we calculate  $\sigma_Y^2 = \mathbb{E}[y_z^2] - \mathbb{E}^2[y_z]$  as follows

$$\sigma_Y^2 = \mathbb{E} \left[ \frac{1}{M^2} \sum_{i=1}^M z_i + \frac{1}{M^2} \sum_{i \neq j} \sqrt{z_i z_j} \right] - \mathbb{E}^2[\sqrt{z_i}]. \quad (52)$$

Since  $z_i$  and  $z_j$  are independent  $\mathcal{X}^2(L)$  distributed random variables, the above equation can be further simplified as

$$\sigma_Y^2 = \frac{L}{M} - \frac{1}{M} \mu_Y^2 \leq \frac{1}{2M} + O\left(\frac{1}{M} L^{-1}\right), \quad (53)$$

where the inequality follows from (51).

In the following, we aim to derive a lower bound on  $(\mu_Y - \sqrt{L(1 - \epsilon_L)})^2$ . By (51), one can have

$$\left(\mu_Y - \sqrt{L(1 - \epsilon_L)}\right)^2 \geq \left(\sqrt{L} - \frac{1}{4\sqrt{L}} - \sqrt{L(1 - \epsilon_L)}\right)^2 \quad (54)$$

Due to the fact that  $\sqrt{1 - \epsilon_L} \leq 1 - \frac{\epsilon_L}{2}$  and further write  $\epsilon_L = \frac{L^{-\frac{1}{4}}}{2}$ , the following result holds.

$$\left(\mu_Y - \sqrt{L(1 - \epsilon_L)}\right)^2 \geq \left(\frac{L^{\frac{1}{4}}}{4} - \frac{1}{4\sqrt{L}}\right)^2 \geq \frac{1}{2M} L^{\frac{1}{2}}, \quad (55)$$

where the second inequality follows from the fact that  $\frac{L^{\frac{1}{4}}}{4} - \frac{1}{4\sqrt{L}} \geq \frac{L^{\frac{1}{4}}}{10} \geq \frac{1}{\sqrt{2M}} L^{\frac{1}{4}}$  with  $M \geq 50$ . It is further noted that the condition, i.e.  $M \geq 50$ , can always hold in our scenario. Substitute (53) and (55) into (49), one can have that

$$\Pr \left( \frac{1}{\sqrt{M}} \sum_{i=1}^M \sqrt{z_i} < b\sqrt{LM(1 + \epsilon_L)} \right) \leq L^{-\frac{1}{2}} + O\left(L^{-\frac{3}{2}}\right). \quad (56)$$

Next, we calculate the denominator. It follows from the fact  $\sum_{i=1}^M z_i = \sum_{i=1}^{LM} y_i$  with  $y_i \sim \mathcal{X}^2(1)$  that

$$\Pr\left(\frac{1}{LM} \sum_{j=1}^M z_j < 1 + \epsilon_L\right) = 1 - \Pr\left(\sum_{i=1}^{LM} y_i \geq LM(1 + \epsilon_L)\right). \quad (57)$$

Then, by applying the *Chernoff bound*, the following result holds

$$\Pr\left(\sum_{i=1}^{LM} y_i \geq LM(1 + \epsilon_L)\right) \leq e^{ML \cdot \min_t \{[-\frac{1}{2} \ln(1-2t) - t(1+\epsilon_L)]\}}, \quad t \in [0, \frac{1}{2}), \quad (58)$$

where the minimum is obtained by setting  $t = \frac{1}{2}(1 - \frac{1}{1+\epsilon_L})$ . Substituting into the above inequality, one can have that

$$\Pr\left(\sum_{i=1}^{LM} y_i \geq LM(1 + \epsilon_L)\right) \leq e^{ML[\frac{1}{2} \ln(1+\epsilon_L) - \frac{1+\epsilon_L}{2} + \frac{1}{2}]}. \quad (59)$$

Due to the fact that  $\frac{1}{2} \ln(1 + \epsilon_L) - \frac{1+\epsilon_L}{2} + \frac{1}{2} \leq \frac{\ln 2 - 1}{2} \epsilon_L^2, \forall \epsilon_L \in [0, 1]$ , one can further have

$$\Pr\left(\sum_{i=1}^{LM} y_i \geq LM(1 + \epsilon_L)\right) \leq e^{-\frac{(1-\ln 2)M}{2} \epsilon_L^2 L}. \quad (60)$$

Thereby, the denominator can be bounded as

$$\Pr\left(\frac{1}{LM} \sum_{j=1}^M z_j < 1 + \epsilon_L\right) \geq 1 - e^{-\frac{(1-\ln 2)M}{2} \epsilon_L^2 L}. \quad (61)$$

Since  $\epsilon_L = \frac{L^{-\frac{1}{4}}}{2}$  and substitute (61) and (56) into (47),  $P_1(r)$  can be bounded as

$$P_1(r) \leq \frac{L^{-\frac{1}{2}} + O(L^{-\frac{3}{2}})}{1 + o(1)} \approx L^{-\frac{1}{2}} + O(L^{-\frac{3}{2}}). \quad (62)$$

This completes the whole proof.

### C. Proof of Lemma 7

Directly calculate (27) is difficult. Instead, we aim to bound the four terms in (27), respectively, in the following.

1) *Calculation of  $P_1(r)$* : By Proposition 2, one can have that  $P_1(r) \leq L^{-\frac{1}{2}} + O(L^{-\frac{3}{2}})$ .

2) *Calculation of  $P_2(\alpha, r)$* : define  $\zeta = \Pr(d_c(\mathbf{c}, \mathbf{h}_{\text{ref}}) \leq (1 + \alpha)r)$ , it follows that

$$P_2(\alpha, r) = (1 - \zeta)^{N'} = e^{N' \ln(1 - \zeta)} \leq e^{-N' \zeta}. \quad (63)$$

Moreover, since the codewords are uniformly distributed on the Grassmann manifold, one can calculate  $\zeta$  as follows

$$\zeta = \frac{A_{\mathcal{B}_{\mathbf{h}_{\text{ref}}}((1+\alpha)r)}}{A_{\text{Manifold}}} = \frac{[(1 + \alpha)r]^{\frac{M-1}{2}}}{2\sqrt{\pi} \left(\frac{M-1}{2}\right)^{\frac{1}{2}}}, \quad (64)$$

where  $A_{\mathcal{B}_{\mathbf{h}_{\text{ref}}}((1+\alpha)r)}$  and  $A_{\text{Manifold}}$  denote the surface areas of the ball  $\mathcal{B}_{\mathbf{h}_{\text{ref}}}((1 + \alpha)r)$  and the Grassmann manifold, respectively. Furthermore, let  $1 + \alpha = r^{-\frac{1}{2}}$  with  $r \approx L^{-\frac{1}{8}}$ , one can have that

$$\zeta = \frac{[(1 + \alpha)r]^{\frac{M-1}{2}}}{2\sqrt{\pi} \left(\frac{M-1}{2}\right)^{\frac{1}{2}}} = \frac{1}{2\sqrt{\pi}} \left(\frac{1}{L}\right)^{\frac{M-1}{32}} \left(\frac{M-1}{2}\right)^{-\frac{1}{2}}. \quad (65)$$

Next, set  $N' = 2^{B_h+1} L^{\frac{M-1}{4}+1} \left(\frac{M-1}{2}\right)^{\frac{1}{2}} \sqrt{\pi}$ , it follows from (63) that

$$P_2(\alpha, r) \leq e^{-N' \zeta} = e^{-2^{B_h} L^{1+\frac{7(M-1)}{32}} \sqrt{\pi}} \leq e^{-L}. \quad (66)$$

3) *Calculation of  $P_3(\alpha, r)$* : define  $X$  the random variable, which indicates the total number of codewords lying inside  $\mathcal{B}_{\mathbf{h}_{\text{ref}}}((1 + \alpha)r)$ . Thereby, it follows that  $X \sim \text{Binomial}(N', \zeta)$ . Given large  $N'$ ,  $X$  can be further approximated as a gaussian distributed random variable, i.e.  $X \sim \mathcal{N}(N'\zeta, N'\zeta(1 - \zeta))$ . Then, one can have that

$$\begin{aligned} P_3(\alpha, r) &= \Pr(X \geq N) \\ &= \Pr\left(\frac{X - N'\zeta}{\sqrt{N'\zeta(1 - \zeta)}} \geq \frac{N - N'\zeta}{\sqrt{N'\zeta(1 - \zeta)}}\right) \\ &\leq e^{-\frac{(N - N'\zeta)^2}{2N'\zeta(1 - \zeta)}} \\ &\leq e^{-\left(\frac{N'\zeta}{2} - N\right)}. \end{aligned} \quad (67)$$

Recall that  $N' = 2^{B_h+1} L^{\frac{M-1}{4}+1} \left(\frac{M-1}{2}\right)^{\frac{1}{2}} \sqrt{\pi}$  and  $\zeta = \frac{1}{2\sqrt{\pi}} \left(\frac{1}{L}\right)^{\frac{M-1}{32}} \left(\frac{M-1}{2}\right)^{-\frac{1}{2}}$ , one can have that

$$\frac{N'\zeta}{2} - N > 2^{B_h} \left(\frac{L}{2} - 1\right) \geq \frac{L}{2} - 1. \quad (68)$$

Substitute the above inequality into (67), the following result holds

$$P_3(\alpha, r) \leq e^{-\left(\frac{L}{2} - 1\right)}. \quad (69)$$

4) *Calculation of  $P_4(\alpha, r) = (1 - P_1(r)) (1 + \frac{2}{\alpha}) D(N')$* : since  $N'$  codewords are uniformly distributed on the Grassmann manifold, one can have

$$\begin{aligned}
D(N') &\leq e^{-\frac{2}{M-1} \ln N'} \\
&= e^{-\frac{2}{M-1} \ln \left( 2^{B_h+1} L^{\frac{M-1}{4}+1} \left(\frac{M-1}{2}\right)^{\frac{1}{2}} \sqrt{\pi} \right)} \\
&= e^{-\frac{2}{M-1} \left( \ln 2^{B_h+1} + \ln L^{\frac{M-1}{4}+1} + \ln \left(\frac{M-1}{2}\right)^{\frac{1}{2}} + \ln 2\sqrt{\pi} \right)} \\
&\leq 2^{-\frac{2B_h}{M-1} L^{-\frac{1}{2} - \frac{2}{M-1}}} \\
&\leq 2^{-\frac{2B_h}{M-1} L^{-\frac{1}{2}}}, \tag{70}
\end{aligned}$$

where the second inequality follows from the facts that  $e^{-\frac{1}{M-1} \ln \left(\frac{M-1}{2}\right)} \leq 1$  and  $e^{-\frac{2}{M-1} \ln 2\sqrt{\pi}} \leq 1$ . Then, substitute  $\alpha = r^{-\frac{1}{2}} - 1$  into  $P_4(\alpha, r)$ , one can have

$$P_4(\alpha, r) \leq \left(1 + \frac{2}{\alpha}\right) D(N') \leq \beta_L 2^{-\frac{2B_h}{M-1}} L^{-\frac{1}{2}}, \tag{71}$$

where  $\beta_L = \frac{1+r^{\frac{1}{2}}}{1-r^{\frac{1}{2}}}$  with  $r = L^{-\frac{1}{8}}$ . Taking the summation on the derived four upper bounds, (28) is straightforward.

#### D. Proof of Theorem 1

Following from the assumption that the elements of stochastic gradient  $\mathbf{g}$  are i.i.d Gaussian distributed with 0 mean and unit variance as aforementioned, one can have

$$B \geq \frac{ML}{2} \log_2 \frac{1}{D}, \tag{72}$$

where the inequality follows from the *converse theorem*;  $B$  denotes the total number of bits;  $D = \frac{\mathbb{E}[\|\mathbf{g} - \hat{\mathbf{g}}\|^2]}{ML} = \mathbb{E}[(g_i - \hat{g}_i)^2], \forall i$ . Rearranging the terms, it follows that

$$\ln \frac{\mathbb{E}[\|\mathbf{g} - \hat{\mathbf{g}}\|^2]}{ML} \geq -\frac{2 \ln 2}{ML} B. \tag{73}$$

Next, we focus on calculation of the upper bound. To begin with, one can show that

$$\frac{\frac{\rho_{\max}^2}{4} 2^{-2B_\rho^*}}{ML \cdot 2^{-\frac{2(B_s^*-1)}{L-1}}} \leq 1, \tag{74}$$

where  $B_\rho^*$  and  $B_s^*$  is defined in (35) and (36), respectively. Thereby, one upper bound of  $\ln \frac{\mathbb{E}[\|\mathbf{g} - \hat{\mathbf{g}}\|^2]}{ML}$  is derived in (75) at the top of next page, where  $\beta_L = \frac{1+L^{-\frac{1}{16}}}{1-L^{-\frac{1}{16}}}$ . Furthermore, since we consider the quantization problem in the low resolution regime, i.e.  $B \leq ML$ , implying

$$\begin{aligned}
\ln \frac{\mathbb{E} [\|\mathbf{g} - \widehat{\mathbf{g}}\|^2]}{ML} &\leq \ln \left( 2 \cdot 2^{-\frac{2(B_s^*-1)}{L-1}} + \beta_L 2^{-\frac{2B_h^*}{M-1}} L^{-\frac{1}{2}} + L^{-\frac{1}{2}} + O(L^{-\frac{3}{2}}) \right) \\
&= \ln \left( 2 \cdot 2^{-\frac{2(B_s^*-1)}{L-1}} \left[ 1 + \frac{\beta_L}{2} L^{-\frac{1}{2}} 2^{\frac{2(B_s^*-1)}{L-1} - \frac{2B_h^*}{M-1}} + \frac{1}{2} L^{-\frac{1}{2}} 2^{\frac{2(B_s^*-1)}{L-1}} + O(L^{-\frac{3}{2}}) \right] \right) \\
&\leq -\frac{2 \ln 2}{L-1} B_s^* + \frac{L+1}{L-1} \ln 2 \\
&\quad + \ln \left( 1 + \frac{\beta_L}{2} L^{-\frac{1}{2}} 2^{\frac{2(B_s^*-1)}{L-1} - \frac{2B_h^*}{M-1}} + \frac{1}{2} L^{-\frac{1}{2}} 2^{\frac{2(B_s^*-1)}{L-1}} + O(L^{-\frac{3}{2}}) \right) \tag{75}
\end{aligned}$$


---

$$\begin{aligned}
&\ln \left( 1 + \frac{\beta_L}{2} L^{-\frac{1}{2}} 2^{\frac{2(B_s^*-1)}{L-1} - \frac{2B_h^*}{M-1}} + \frac{1}{2} L^{-\frac{1}{2}} 2^{\frac{2(B_s^*-1)}{L-1}} + O(L^{-\frac{3}{2}}) \right) \\
&\leq \ln \left( 1 + 2\beta_L L^{-\frac{1}{2}} + 2L^{-\frac{1}{2}} + O(L^{-\frac{3}{2}}) \right) \\
&\leq 2(\beta_L + 1)L^{-\frac{1}{2}} - 2(\beta_L + 1)^2 L^{-1} + O(L^{-\frac{3}{2}}). \tag{76}
\end{aligned}$$


---

$B_s < L$ , one can have that  $2^{\frac{2(B_s^*-1)}{L-1}} < 4$ . Then, it follows that one upper bound of the third term in (75) can be derived in (76). Substituting (76) into (75) with (36) involved, we have

$$\ln \frac{\mathbb{E} [\|\mathbf{g} - \widehat{\mathbf{g}}\|^2]}{ML} \leq c_{\text{gap}} - \frac{2 \ln 2}{ML} B + O\left(L^{-\frac{3}{2}}\right), \quad L \rightarrow \infty, \tag{77}$$

where  $c_{\text{gap}} = \ln 2 - \ln \frac{2L}{L-1} + \frac{2}{ML} \ln \frac{ML + \sqrt{ML}}{2} + \frac{L-1}{L} \ln \frac{2L}{L-1} + \frac{2}{L} \ln 2 + \frac{2(M-1)}{ML} \ln \frac{2}{M-1} + \frac{\ln 2}{ML} + \frac{M-1}{ML} \ln \beta_L M \sqrt{L} + 2(\beta_L + 1)L^{-\frac{1}{2}} - 2(\beta_L + 1)^2 L^{-1}$  is a constant given  $M$  and  $L$ . This completes the whole proof.

### E. Proof of Theorem 2

Take Assumption 3, one can have that

$$F_{n+1} - F_n \leq \bar{\mathbf{g}}_n^T (\boldsymbol{\theta}_{n+1} - \boldsymbol{\theta}_n) + \sum_{i=1}^{\text{Dim}} \frac{l_i}{2} (\boldsymbol{\theta}_{n+1} - \boldsymbol{\theta}_n)_i^2, \tag{78}$$

where  $F_n$  denotes the global objective at the  $n$ -th iteration;  $\bar{\mathbf{g}}_n = \nabla F_n$  is the gradient of the global objective and  $\boldsymbol{\theta}_n = \boldsymbol{\theta}[n]$  is the model parameter;  $(\boldsymbol{\theta}_{n+1} - \boldsymbol{\theta}_n)_i^2$  denotes the square of the  $i$ -th coefficient of  $\boldsymbol{\theta}_{n+1} - \boldsymbol{\theta}_n$ .

By taking the expectation under the current model  $\boldsymbol{\theta}_n$  on the both side of the above inequality, it follows that

$$\begin{aligned} \mathbb{E} [F_{n+1} - F_n | \boldsymbol{\theta}_n] &\leq \mathbb{E} [\bar{\mathbf{g}}_n^T (\boldsymbol{\theta}_{n+1} - \boldsymbol{\theta}_n) | \boldsymbol{\theta}_n] \\ &\quad + \mathbb{E} \left[ \sum_{i=1}^{\text{Dim}} \frac{l_i}{2} (\boldsymbol{\theta}_{n+1} - \boldsymbol{\theta}_n)_i^2 | \boldsymbol{\theta}_n \right]. \end{aligned} \quad (79)$$

In the following, we treat  $\mathbb{E} [\bar{\mathbf{g}}_n^T (\boldsymbol{\theta}_{n+1} - \boldsymbol{\theta}_n) | \boldsymbol{\theta}_n]$  and  $\mathbb{E} \left[ \sum_{i=1}^{\text{Dim}} \frac{l_i}{2} (\boldsymbol{\theta}_{n+1} - \boldsymbol{\theta}_n)_i^2 | \boldsymbol{\theta}_n \right]$ , separately.

**(1) Calculation of  $\mathbb{E} [\bar{\mathbf{g}}_n^T (\boldsymbol{\theta}_{n+1} - \boldsymbol{\theta}_n) | \boldsymbol{\theta}_n]$ :**

Plug in the model-update step (7), one can have that

$$\mathbb{E} [\bar{\mathbf{g}}_n^T (\boldsymbol{\theta}_{n+1} - \boldsymbol{\theta}_n) | \boldsymbol{\theta}_n] = \mathbb{E} \left[ \frac{-\eta}{K} \bar{\mathbf{g}}_n^T \sum_{k=1}^K \widehat{\mathbf{g}}_n^{(k)} | \boldsymbol{\theta}_n \right], \quad (80)$$

where  $\widehat{\mathbf{g}}_n^{(k)} = \left\| \widehat{\mathbf{g}}_n^{(k)} \right\| \widehat{\mathbf{f}}_n^{(k)}$  denotes the quantized version of the stochastic gradient from the  $k$ -th edge device at the  $n$ -th iteration. Let  $\widehat{\mathbf{g}}_n^{(k)} = \mathbf{g}_n^{(k)} - \Delta_n^{(k)}$  with  $\|\Delta_n^{(k)}\|^2$  denoting the quantization error, the above equation can be rewritten as

$$\mathbb{E} [\bar{\mathbf{g}}_n^T (\boldsymbol{\theta}_{n+1} - \boldsymbol{\theta}_n) | \boldsymbol{\theta}_n] = \mathbb{E} \left[ \frac{-\eta}{K} \bar{\mathbf{g}}_n^T \sum_{k=1}^K (\mathbf{g}_n^{(k)} - \Delta_n^{(k)}) | \boldsymbol{\theta}_n \right]. \quad (81)$$

Furthermore, given  $\|\widehat{\mathbf{g}}_n^{(k)}\| = \|\mathbf{g}_n^{(k)}\| - \Delta\rho_n^{(k)}$  and  $\widehat{\mathbf{f}}_n^{(k)} = \mathbf{f}_n^{(k)} - \Delta\mathbf{f}_n^{(k)}$  with  $\mathbb{E}[\Delta\rho_n^{(k)}] = 0$  and  $\mathbb{E}[\Delta\mathbf{f}_n^{(k)}] = \mathbf{0} \in \mathbb{R}^{\text{Dim} \times 1}$ ,  $\forall n, k$ , it follows that

$$\mathbb{E}[\widehat{\mathbf{g}}_n^{(k)}] = \mathbb{E} [(\|\mathbf{g}_n^{(k)}\| - \Delta\rho_n^{(k)}) (\mathbf{f}_n^{(k)} - \Delta\mathbf{f}_n^{(k)})] = \mathbb{E}[\mathbf{g}_n^{(k)}]. \quad (82)$$

Due to the fact that  $\mathbb{E}[\widehat{\mathbf{g}}_n^{(k)}] = \mathbb{E}[\mathbf{g}_n^{(k)}] - \mathbb{E}[\Delta_n^{(k)}]$ , one can conclude that  $\mathbb{E}[\Delta_n^{(k)}] = \mathbf{0} \in \mathbb{R}^{\text{Dim} \times 1}$ .

Thereby, it follows from (81) that

$$\mathbb{E} [\bar{\mathbf{g}}_n^T (\boldsymbol{\theta}_{n+1} - \boldsymbol{\theta}_n) | \boldsymbol{\theta}_n] = \mathbb{E} \left[ \frac{-\eta}{K} \bar{\mathbf{g}}_n^T \sum_{k=1}^K \mathbf{g}_n^{(k)} | \boldsymbol{\theta}_n \right] = -\eta \|\bar{\mathbf{g}}_n\|^2, \quad (83)$$

where the second equality follows from Assumption 4 that  $\mathbb{E}[\mathbf{g}_n^{(k)}] = \bar{\mathbf{g}}_n, \forall k$ .

**(2) Calculation of  $\mathbb{E} \left[ \sum_{i=1}^{\text{Dim}} \frac{l_i}{2} (\boldsymbol{\theta}_{n+1} - \boldsymbol{\theta}_n)_i^2 | \boldsymbol{\theta}_n \right]$ :**

Let  $l_0 = \|\mathbf{1}\|_\infty$  with  $\mathbf{1}$  defined in Assumption 4, one can have that

$$\mathbb{E} \left[ \sum_{i=1}^{\text{Dim}} \frac{l_i}{2} (\boldsymbol{\theta}_{n+1} - \boldsymbol{\theta}_n)_i^2 | \boldsymbol{\theta}_n \right] \leq \mathbb{E} \left[ \frac{l_0}{2} \|\boldsymbol{\theta}_{n+1} - \boldsymbol{\theta}_n\|^2 | \boldsymbol{\theta}_n \right] \quad (84)$$

Since  $\boldsymbol{\theta}_{n+1} - \boldsymbol{\theta}_n = \frac{\eta}{K} \sum_{k=1}^K \widehat{\mathbf{g}}_n^{(k)}$  with  $\widehat{\mathbf{g}}_n^{(k)}$  denoting the quantized stochastic gradient from the  $k$ -th local device at the  $n$ -th iteration, it follows that

$$\begin{aligned} \mathbb{E} \left[ \sum_{i=1}^{\text{Dim}} \frac{l_i}{2} (\boldsymbol{\theta}_{n+1} - \boldsymbol{\theta}_n)_i^2 | \boldsymbol{\theta}_n \right] &\leq \mathbb{E} \left[ \frac{l_0}{2} \|\boldsymbol{\theta}_{n+1} - \boldsymbol{\theta}_n\|^2 | \boldsymbol{\theta}_n \right] \\ &= \mathbb{E} \left[ \frac{\eta^2 l_0}{2K^2} \left\| \sum_{k=1}^K \widehat{\mathbf{g}}_n^{(k)} \right\|^2 | \boldsymbol{\theta}_n \right]. \end{aligned} \quad (85)$$

Due to the fact that arbitrary two different high-dimensional vectors are quasi-orthogonal, the following result holds

$$\left\| \sum_{k=1}^K \widehat{\mathbf{g}}_n^{(k)} \right\|^2 = \left( \sum_{k=1}^K \widehat{\mathbf{g}}_n^{(k)} \right)^T \left( \sum_{k=1}^K \widehat{\mathbf{g}}_n^{(k)} \right) \approx \sum_{k=1}^K \|\widehat{\mathbf{g}}_n^{(k)}\|^2. \quad (86)$$

Substituting (86) into (85), the following inequality follows

$$\mathbb{E} \left[ \sum_{i=1}^{\text{Dim}} \frac{l_i}{2} (\boldsymbol{\theta}_{n+1} - \boldsymbol{\theta}_n)_i^2 | \boldsymbol{\theta}_n \right] \leq \frac{\eta^2 l_0}{2K^2} \sum_{k=1}^K \mathbb{E} \left[ \|\widehat{\mathbf{g}}_n^{(k)}\|^2 | \boldsymbol{\theta}_n \right]. \quad (87)$$

Moreover, since  $\widehat{\mathbf{g}}_n^{(k)} = \mathbf{g}_n^{(k)} - \boldsymbol{\Delta}_n^{(k)}$ , one can have that

$$\|\widehat{\mathbf{g}}_n^{(k)}\|^2 = (\mathbf{g}_n^{(k)} - \boldsymbol{\Delta}_n^{(k)})^T (\mathbf{g}_n^{(k)} - \boldsymbol{\Delta}_n^{(k)}) \stackrel{(a)}{\approx} \|\mathbf{g}_n^{(k)}\|^2 + \|\boldsymbol{\Delta}_n^{(k)}\|^2,$$

where (a) follows from the same argument that high-dimensional vectors are quasi-orthogonal.

Then, the inequality (87) can be simplified as

$$\begin{aligned} \mathbb{E} \left[ \sum_{i=1}^{\text{Dim}} \frac{l_i}{2} (\boldsymbol{\theta}_{n+1} - \boldsymbol{\theta}_n)_i^2 | \boldsymbol{\theta}_n \right] &\leq \frac{\eta^2 l_0}{2K^2} \sum_{k=1}^K \mathbb{E} \left[ \|\mathbf{g}_n^{(k)}\|^2 | \boldsymbol{\theta}_n \right] \\ &\quad + \frac{\eta^2 l_0}{2K^2} \sum_{k=1}^K \mathbb{E} \left[ \|\mathbf{g} - \widehat{\mathbf{g}}\|^2 \right], \end{aligned} \quad (88)$$

where  $\mathbb{E}[\|\mathbf{g} - \widehat{\mathbf{g}}\|^2] = \mathbb{E}[\|\boldsymbol{\Delta}_n^{(k)}\|^2 | \boldsymbol{\theta}_n], \forall n, k$ . Next, by Assumption 4, it can be obtained that  $\mathbb{E}[\|\mathbf{g}_n^{(k)}\|^2 | \boldsymbol{\theta}_n] = \|\boldsymbol{\sigma}\|^2 + \|\bar{\mathbf{g}}_n\|^2$  with  $\bar{\mathbf{g}}_n$  denoting the gradient of at the  $n$ -th iteration. Then, one can have that

$$\begin{aligned} \mathbb{E} \left[ \sum_{i=1}^{\text{Dim}} \frac{l_i}{2} (\boldsymbol{\theta}_{n+1} - \boldsymbol{\theta}_n)_i^2 | \boldsymbol{\theta}_n \right] &\leq \frac{\eta^2 l_0}{2K} (\|\boldsymbol{\sigma}\|^2 + \|\bar{\mathbf{g}}_n\|^2) \\ &\quad + \frac{\eta^2 l_0}{2K} \mathbb{E} \left[ \|\mathbf{g} - \widehat{\mathbf{g}}\|^2 \right]. \end{aligned} \quad (89)$$

$$\begin{aligned}
F_0 - F^* &\geq F_0 - \mathbb{E}[F_N] \\
&= \mathbb{E} \left[ \sum_{n=0}^{N-1} (F_n - F_{n+1}) \right] \\
&\geq \mathbb{E} \left[ \sum_{n=0}^{N-1} \left( \eta \|\bar{\mathbf{g}}_n\|^2 - \frac{\eta^2 l_0}{2K} \mathbb{E} [\|\mathbf{g} - \hat{\mathbf{g}}\|^2] - \frac{\eta^2 l_0}{2K} (\|\boldsymbol{\sigma}\|^2 + \|\bar{\mathbf{g}}_n\|^2) \right) \right] \tag{91}
\end{aligned}$$

By substituting (83) and (89) into (79), the following result holds

$$\begin{aligned}
\mathbb{E} [F_{n+1} - F_n | \boldsymbol{\theta}_n] &\leq -\eta \|\bar{\mathbf{g}}_n\|^2 + \frac{\eta^2 l_0}{2K} \mathbb{E} [\|\mathbf{g} - \hat{\mathbf{g}}\|^2] \\
&\quad + \frac{\eta^2 l_0}{2K} (\|\boldsymbol{\sigma}\|^2 + \|\bar{\mathbf{g}}_n\|^2). \tag{90}
\end{aligned}$$

Next, extend the expectation over randomness in the trajectory, and perform a telescoping sum over the all the iterations, one lower bound on  $F_0 - F^*$  is derived in (91).

Substituting  $\eta = \frac{1}{\sqrt{l_0 N}}$  into the above inequality and rearrange the terms, it follows that

$$\mathbb{E} \left[ \frac{1}{N} \sum_{n=0}^{N-1} \|\bar{\mathbf{g}}_n\|^2 \right] \leq \frac{\sqrt{l_0} \left( \frac{1}{2K} \mathbb{E} [\|\mathbf{g} - \hat{\mathbf{g}}\|^2] + \frac{\|\boldsymbol{\sigma}\|^2}{2K} + F_0 - F^* \right)}{\sqrt{N} - \frac{\sqrt{l_0}}{2K}}. \tag{92}$$

This completes the whole proof.

## REFERENCES

- [1] G. Zhu, D. Liu, Y. Du, C. You, J. Zhang, and K. Huang, “Towards an intelligent edge: Wireless communication meets machine learning,” 2018. [Online]. Available: <https://arxiv.org/pdf/1809.00343.pdf>.
- [2] J. Konečný, H. B. McMahan, F. X. Yu, P. Richtárik, A. T. Suresh, and D. Bacon, “Federated learning: Strategies for improving communication efficiency,” 2016. [Online]. <https://arxiv.org/pdf/1610.05492.pdf>.
- [3] R. Tandon, Q. Lei, A. G. Dimakis, and N. Karampatziakis, “Gradient coding: Avoiding stragglers in distributed learning,” in *Proc. of Intl. Conf. Mach. Learning (ICML)*, pp. 3368–3376, Jul. 2017.
- [4] T. Chen, G. Giannakis, T. Sun, and W. Yin, “LAG: Lazily aggregated gradient for communication-efficient distributed learning,” in *Proc. of Adv. Neural Inf. Proc. Systems (NIPS)*, pp. 5050–5060, Dec. 2018.
- [5] H. Zhang, J. Li, K. Kara, D. Alistarh, J. Liu, and C. Zhang, “ZIPML: Training linear models with end-to-end low precision, and a little bit of deep learning,” in *Proc. of the 34th Intl. Conf. Mach. Learning (ICML)*, pp. 4035–4043, Aug. 2017.
- [6] D. Alistarh, D. Grubic, J. Li, R. Tomioka, and M. Vojnovic, “QSGD: Communication-efficient SGD via gradient quantization and encoding,” in *Proc. of Adv. Neural Inf. Proc. Systems (NIPS)*, pp. 1709–1720, Dec. 2017.
- [7] J. Bernstein, Y.-X. Wang, K. Azizzadenesheli, and A. Anandkumar, “signSGD: Compressed optimisation for non-convex problems,” in *Proc. of Intl. Conf. Mach. Learning (ICML)*, vol. 80, pp. 559–568, Jul. 2018.

- [8] J. Wu, W. Huang, J. Huang, and T. Zhang, "Error compensated quantized SGD and its applications to large-scale distributed optimization," in *Proc. of Intl. Conf. Mach. Learning (ICML)*, pp. 5321–5329, Jul. 2018.
- [9] S. Zheng, Z. Huang, and J. T. Kwok, "Communication-efficient distributed blockwise momentum SGD with error-feedback," 2019. [Online]. <https://arxiv.org/pdf/1905.10936.pdf>.
- [10] D. P. Kingma and J. Ba, "ADAM: A method for stochastic optimization," 2014. [Online]. <https://arxiv.org/pdf/1412.6980.pdf>
- [11] A. Gersho and R. M. Gray, *Vector quantization and signal compression*. Springer Science & Business Media, 2012, vol. 159.
- [12] D. J. Love, R. W. Heath, Jr, V. K. Lau, D. Gesbert, B. D. Rao, and M. Andrews, "An overview of limited feedback in wireless communication systems," *IEEE J. Sel. Areas Commun.*, vol. 26, no. 8, pp. 1341–1365, 2008.
- [13] K. K. Mukkavilli, A. Sabharwal, E. Erkip, and B. Aazhang, "On beamforming with finite rate feedback in multiple-antenna systems," *IEEE Trans. Inf. Theory*, vol. 49, no. 10, pp. 2562–2579, 2003.
- [14] D. J. Love, R. W. Heath, Jr, and T. Strohmer, "Grassmannian beamforming for multiple-input multiple-output wireless systems," *IEEE Trans. Inf. Theory*, vol. 49, no. 10, p. 2735, 2003.
- [15] V. Raghavan, R. W. Heath, Jr, and A. M. Sayeed, "Systematic codebook designs for quantized beamforming in correlated MIMO channels," *IEEE J. Sel. Areas Commun.*, vol. 25, no. 7, pp. 1298–1310, 2007.
- [16] T. Kim, D. J. Love, and B. Clerckx, "MIMO systems with limited rate differential feedback in slowly varying channels," *IEEE Trans. Commun.*, vol. 59, no. 4, pp. 1175–1189, 2011.
- [17] K. Huang, R. W. Heath, Jr, and J. G. Andrews, "Limited feedback beamforming over temporally-correlated channels," *IEEE Trans. Sig. Process.*, vol. 57, no. 5, pp. 1959–1975, 2009.
- [18] P. Xia and G. B. Giannakis, "Design and analysis of transmit-beamforming based on limited-rate feedback," *IEEE Trans. Sig. Process.*, vol. 54, no. 5, pp. 1853–1863, 2006.
- [19] J. Nam, J.-Y. Ahn, A. Adhikary, and G. Caire, "Joint spatial division and multiplexing: Realizing massive MIMO gains with limited channel state information," in *Proc. of Annual Conf. Inf. Sciences and Systems (CISS)*, pp. 1–6, Mar. 2012.
- [20] J. Choi, Z. Chance, D. J. Love, and U. Madhow, "Noncoherent trellis coded quantization: A practical limited feedback technique for massive MIMO systems," *IEEE Trans. Commun.*, vol. 61, no. 12, pp. 5016–5029, 2013.
- [21] M. S. Sim, J. Park, C.-B. Chae, and R. W. Heath, Jr, "Compressed channel feedback for correlated massive MIMO systems," *J. Commun. and Networks*, vol. 18, no. 1, pp. 95–104, 2016.
- [22] I. S. Dhillon, J. R. Heath, T. Strohmer, and J. A. Tropp, "Constructing packings in Grassmannian manifolds via alternating projection," *Experimental mathematics*, vol. 17, no. 1, pp. 9–35, 2008.
- [23] V. Lau, Y. Liu, and T.-A. Chen, "On the design of MIMO block-fading channels with feedback-link capacity constraint," *IEEE Trans. Commun.*, vol. 52, no. 1, pp. 62–70, 2004.
- [24] J. H. Conway and N. J. A. Sloane, *Sphere packings, lattices and groups*. Springer Science & Business Media, vol. 290, 2013.
- [25] W. Dai, Y. Liu, and B. Rider, "Quantization bounds on Grassmann manifolds and applications to MIMO communications," *IEEE Trans. Inf. Theory*, vol. 54, no. 3, pp. 1108–1123, 2008.