

# BREAKING MEMORY AND COMMUNICATION BARRIERS IN MODEL-PARALLEL FINE-TUNING OF LARGE LANGUAGE MODELS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Model parallelism (MP) has emerged as a promising paradigm for distributed large language model (LLM) training across multiple computing nodes. Yet, almost all existing works about MP focus on first-order methods, which faces two persistent challenges: high communication costs from transmitting activations and gradients, and substantial memory overhead from caching them. Zeroth-order (ZO) methods, by avoiding gradient computation and storage, can naturally alleviate both memory and communication bottlenecks, but they have been largely unexplored in MP for LLM fine-tuning. In this work, we propose **SparQ**, a ZO MP framework with **S**plit layer **a**llocation informed by **Q**uantization-induced activation sparsity, designed to reduce memory and communication costs. **SparQ** builds on three key components: (1) leveraging the gradient-free nature of ZO optimization to eliminate gradient storage and transmission, significantly reducing memory and communication demands incurred by gradients; (2) applying quantization to induce activation sparsity that can be encoded with sparse representations; (3) strategically placing split layers at activation-sparse regions and using sparse representation to lower communication cost from activations almost without compromising model quality. Theoretically, **SparQ** achieves a sublinear convergence rate in non-convex settings, matching that of centralized ZO methods. Empirically, **SparQ** reduces GPU memory usage by over 3x and communication cost by 50%+ compared to state-of-the-art baselines, while maintaining comparable model performance.

## 1 INTRODUCTION

Large language models (LLMs) have demonstrated strong generalization capabilities, driving their adoption across diverse downstream tasks (Vaswani, 2017; Zhao et al., 2023; Naveed et al., 2025). However, fine-tuning such massive models remains challenging: it requires storing billions of parameters and extra information, such as activations and gradients (Kaddour et al., 2023; Han et al., 2024), often exceeding the memory capacity of a single GPU or machine. Model parallelism (MP), which partitions models across multiple nodes, is a widely used solution, with frameworks such as Megatron-LM (Shoeybi et al., 2019), ZeRO (Rajbhandari et al., 2020), DeepSpeed (Rasley et al., 2020), and MegaScale (Jiang et al., 2024) demonstrating its effectiveness. Most existing works focus on applying first-order (FO) method in the MP scenario. However, the combination of FO and MP inevitably introduce two substantial costs: (1) **High Communication Cost**: FO methods (e.g., SGD) require frequent exchange of large gradients and activations across nodes, making communication a major bottleneck; (2) **High Memory Cost**: FO methods also demand storing gradients, optimizer states, and cached activations, further straining memory and often exceeding a single node’s capacity. Scaling typically requires more GPUs and aggressive parallelization, which amplifies both communication and hardware costs.

Recently, zeroth-order (ZO) methods have gained significant attention because they only require forward passes, offering substantial memory savings (Malladi et al., 2023; Zhang et al., 2024; Zhao et al., 2025). However, existing ZO research largely focuses on optimization theory or single-node setups, but the potential of applying ZO optimization within a MP framework for fine-tuning LLMs remains unexplored. This gap naturally raises a key research question:

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

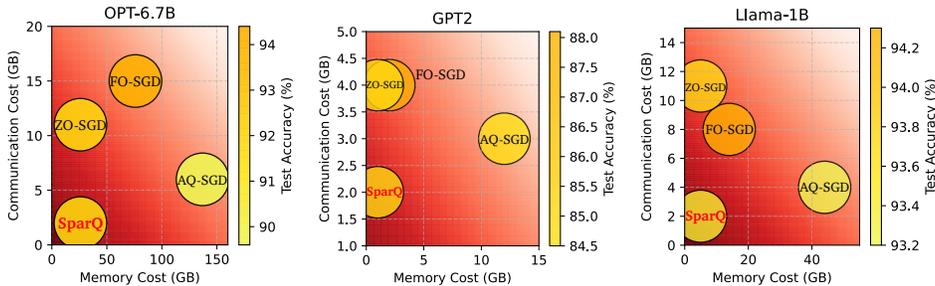


Figure 1: Illustration of Memory Cost, Communication Cost and Test Accuracy. Here we use SST-2 to fine-tune ReLU-based OPT, GELU-based GPT and SwiGLU-based Llama models.

*Q: What extra, undiscovered advantages in terms of communication and memory cost can be realized by combining ZO and MP for fine-tuning LLMs?*

The major contributions of our work are summarized as follows:

- We empirically observe employing quantization to immediate outputs after commonly used activation functions (e.g., ReLU, GELU, SwiGLU) in Transformer (Vaswani, 2017) architectures during LLM fine-tuning can result in pervasively sparse activations, with the proportion of nonzero entries across layers ranging between 1% ~ 30%.
- Inspired by such inherent sparsity, we propose **SparQ**, a ZO model-parallel framework with **Split layer allocation** informed by **Quantization-induced activation sparsity**. It includes three key components: **(1) ZO optimization for gradient-free fine-tuning**: Prior studies have shown that ZO optimization is effective for LLM fine-tuning due to the presence of low-dimensional subspaces or low-rank Hessian landscapes (Malladi et al., 2023). We exploit its gradient-free nature to eliminate the need for storing and transmitting gradients in MP-based fine-tuning, significantly reducing both memory and communication overhead introduced by gradients. **(2) Quantization-induced high activation sparsity**: In Sec. 2.4, we observe that applying quantization to intermediate outputs after activation functions in Transformer architectures induces higher sparsity than unquantized activations during LLM fine-tuning. **(3) Split layer allocation guided by activation sparsity**: SparQ strategically places split layers at these naturally sparse regions. By transmitting activations in sparse representation (i.e., only transmit nonzero entries and their indices) across partitions, communication cost is substantially reduced while maintaining model performance.
- We prove that SparQ across multiple nodes achieves a sublinear convergence rate of  $\mathcal{O}(\sqrt{d/T})$  for non-convex functions, matching the rate of centralized ZO stochastic gradient descent (ZO-SGD).
- We evaluate SparQ’s test accuracy, memory and communication costs on various LLM fine-tuning tasks and observe that: (1) SparQ achieves superior efficiency when memory and communication costs are jointly considered while maintaining comparable model performance, as shown in Fig. 1; (2) SparQ reduces communication cost by about 50% ~ 70%, compared to FO-SGD and the state-of-the-art activation compression method AQ-SGD (Wang et al., 2022).

## 2 FORMULATIONS AND PRELIMINARIES

Given the loss function  $f$ , our goal is to minimize the following objective function

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) := \mathbb{E}_{\xi \sim \mathcal{D}} [f(\xi; \mathbf{x})],$$

where  $\mathbf{x}$  is a  $d$ -dimensional model parameter, and  $\xi$  is a data sample (model input) selected from the distribution  $\mathcal{D}$ . In the following subsections, we formulate two core components of SparQ: 1) model parallelism technique, and 2) ZO optimization’s gradient estimation methods.

### 2.1 MODEL PARALLELISM FORMULATION

We begin by formulating model parallelism. In a model-parallel setting, the entire model is partitioned into  $M$  parts, typically, with each computing node responsible for a distinct subset of the model’s parameters. For clarity and brevity, we focus on the scenario in our main paper where the model is only divided into two submodels ( $M = 2$ ), each residing on a separate node. This facilitates a

108 concise presentation of the mathematical formulations, and an extension to an arbitrary number of  
 109 partitions ( $M > 2$ ) is provided in Appendix C.3.

110  
 111 The partition of model parameters denotes as  $\mathbf{x} = \text{col}[\mathbf{x}^{(1)}, \mathbf{x}^{(2)}]$ , where  $\mathbf{x} \in \mathbb{R}^d$ ,  $\mathbf{x}^{(1)} \in \mathbb{R}^{d_1}$ ,  
 112  $\mathbf{x}^{(2)} \in \mathbb{R}^{d_2}$ , and  $d_1 + d_2 = d$ . Under this configuration, the loss function can be written as the  
 113 composition of two functions:

$$114 \quad f(\xi; \mathbf{x}) := F \left( S_2 \left( S_1(\xi; \mathbf{x}^{(1)}); \mathbf{x}^{(2)} \right) \right), \quad (1)$$

115 where  $F$  is the criterion function,  $S_i$  represents a subset of network layers located on node  $\mathcal{M}_i$   
 116 with model parameters  $\mathbf{x}^{(i)}$ . Specifically, for a given input  $\xi$ ,  $S_i(\xi; \mathbf{x}^{(i)})$  computes the forward  
 117 activation on node  $\mathcal{M}_i$ , which is then transferred to the next node  $\mathcal{M}_{i+1}$  where  $S_2(\cdot; \mathbf{x}^{(2)})$  continues  
 118 the computation. For notational simplicity, we also express the loss function using operator notation:

$$119 \quad f(\xi; \mathbf{x}) = (F \circ S_2|_{\mathbf{x}^{(2)}} \circ S_1|_{\mathbf{x}^{(1)}})(\xi), \quad (2)$$

120 where  $\circ$  denotes function composition (applying one function to the result of another).  $S_1|_{\mathbf{x}^{(1)}}$   
 121 represents the first sub-model parameterized by  $\mathbf{x}^{(1)}$  and acting on  $\xi$ .  $S_2|_{\mathbf{x}^{(2)}}$  represents the second  
 122 sub-model parameterized by  $\mathbf{x}^{(2)}$  and acting on the output of  $S_1|_{\mathbf{x}^{(1)}}$ .

123  
 124 This two-node case ( $M = 2$ ) can naturally be generalized to the multi-node case ( $M > 2$ ) by further  
 125 partitioning the model. Accordingly, the loss functions in (1) and (2) become

$$126 \quad f(\xi; \mathbf{x}) = F \left( S_M \left( \dots \left( S_2 \left( S_1(\xi; \mathbf{x}^{(1)}); \mathbf{x}^{(2)} \right); \dots \right); \mathbf{x}^{(M)} \right) \right) = (F \circ S_M|_{\mathbf{x}^{(M)}} \circ \dots \circ S_2|_{\mathbf{x}^{(2)}} \circ S_1|_{\mathbf{x}^{(1)}})(\xi).$$

127  
 128 This formulation clearly delineates the flow of data across nodes in a MP framework, laying the  
 129 groundwork for our subsequent discussions on optimization, memory and communication efficiency.

## 130 2.2 ZEROth-ORDER OPTIMIZATION FORMULATION

131  
 132 Next, let us recap the fundamentals of ZO optimization. It is a gradient-free method that approximates  
 133 the gradient by utilizing the finite difference method rather than explicit differentiation. Given a  
 134 smoothing parameter  $\mu > 0$ , the number of perturbations  $P$  and random perturbation vectors  $\mathbf{u}$  drawn  
 135 from a Gaussian distribution or a uniform ball, the ZO gradient estimate  $\hat{G}$  can be computed as:

$$136 \quad \hat{G} = \frac{1}{P} \sum_{i=1}^P g_i \cdot \mathbf{u}_i = \frac{1}{P} \sum_{i=1}^P \frac{f(\xi; \mathbf{x} + \mu \mathbf{u}_i) - f(\xi; \mathbf{x})}{\mu} \cdot \mathbf{u}_i, \quad (3)$$

137  
 138 where Eq. (3) represents a biased forward difference approach. We distribute the content about  
 139 unbiased central difference approach to Sec. C.2. These two approaches have been widely used to  
 140 estimate the gradient in a ZO context (Ghadimi & Lan, 2013; Liu et al., 2020).

## 141 2.3 MODEL PARALLELISM MEETS ZEROth-ORDER OPTIMIZATION

142  
 143 Having established the formulation for both ZO optimization and model parallelism, we now describe  
 144 how to integrate ZO optimization into a model-parallel framework. For clarity, we still focus on the  
 145 two-node case ( $M = 2$ ) here and distributed the multi-node case ( $M > 2$ ) to Appendix C.3. In this  
 146 setup, node  $\mathcal{M}_1$  is responsible solely for computing and transmitting the forward activations, while  
 147 node  $\mathcal{M}_2$  performs the gradient estimation, parameter updates and gradient scalar transmission.

148  
 149 In this work, we utilize classic Zeroth-Order Stochastic Gradient Descent (ZO-SGD) as our optimizer.  
 150 Specifically, within the second submodel on  $\mathcal{M}_2$ , the ZO gradient scalar  $g_t$  is computed using a  
 151 forward difference method as follows:

$$152 \quad g_{i,t} = \frac{F(a_{i,t}^+; \mathbf{x}_t^{(2)} + \mu \mathbf{u}_{i,t}^{(2)}) - F(a_{i,t}; \mathbf{x}_t^{(2)})}{\mu}, \quad (4)$$

153  
 154 where  $t$  is the iteration index, and  $\mathbf{u}_{i,t}^{(2)}$  is the perturbation vector generated on node  $\mathcal{M}_2$  during the  
 155  $t$ -th iteration.  $a^+$  and  $a$  are the activations obtained from the preceding sub-model on  $\mathcal{M}_1$ , defined as

$$156 \quad a_{i,t}^+ = S_1(\xi_t; \mathbf{x}_t^{(1)} + \mu \mathbf{u}_{i,t}^{(1)}), \quad a_{i,t} = S_1(\xi_t; \mathbf{x}_t^{(1)}). \quad (5)$$

157  
 158 Then, the final zeroth-order gradient estimate is given by  $\hat{G}_t = \frac{1}{P} \sum_{i=1}^P g_{i,t} \cdot \mathbf{u}_{i,t}^{(2)}$ .

159  
 160 Integrating ZO optimization into a model-parallel framework offers a promising avenue to alleviate  
 161 memory burdens associated with large-scale training. Unlike traditional FO methods that require  
 storing gradients, ZO optimization relies exclusively on forward passes, thereby eliminating memory

usage from storing gradients. However, this integration introduces a unique communication challenge. Since ZO optimization often depends on multiple perturbations to approximate gradients accurately, each computing node must transmit several forward activations across split layer per iteration. In contrast, FO methods generally exchange only a single forward activation and one backward gradient per iteration. Thus, while ZO optimization enables significant memory savings, its deployment in MP necessitates careful optimization of communication overhead to maintain training stability and performance. In Sec. 3, we will introduce how we address communication bottleneck.

### 2.4 KEY OBSERVATIONS

We first summarize several key observations from Fig. 2 that motivate SparQ.

**(1) High Sparsity of Original ReLU-Based Activations in LLM Fine-Tuning.** Li et al. (2023) reported that ReLU-based activations exhibit high sparsity during LLM pre-training. Consistent with this finding, we observe a similar pattern in LLM fine-tuning in Fig. 2a: the proportion of nonzero values in nearly all ReLU-based activations remains below 10%. **(2) Low Sparsity of Original SwiGLU- and GELU-based Activations in LLM Fine-Tuning.** Smoother activation functions, such as SwiGLU and GELU, produce predominantly small but nonzero values, resulting in consistently dense activations with 99% ~ 100% nonzero entries. **(3) Quantization Induces High Sparsity across Various Activations.** Applying quantization can dramatically enhance activation sparsity during fine-tuning. Specifically, by using 4-bit quantization, ReLU-based activations become even sparser (below 5% nonzero entries), SwiGLU-based activations stabilize below 20%, and nearly all GELU-based activations fall under 5% nonzero entries. This high sparsity motivates us to split the model immediately after the activation functions since activations can be efficiently encoded in sparse representation (described in Sec. C.1), significantly reducing the communication burden associated with their transmission because we only need to transmit the nonzero values and their indices. The specific split layer selections in our experiments are described in Sec. B.1.

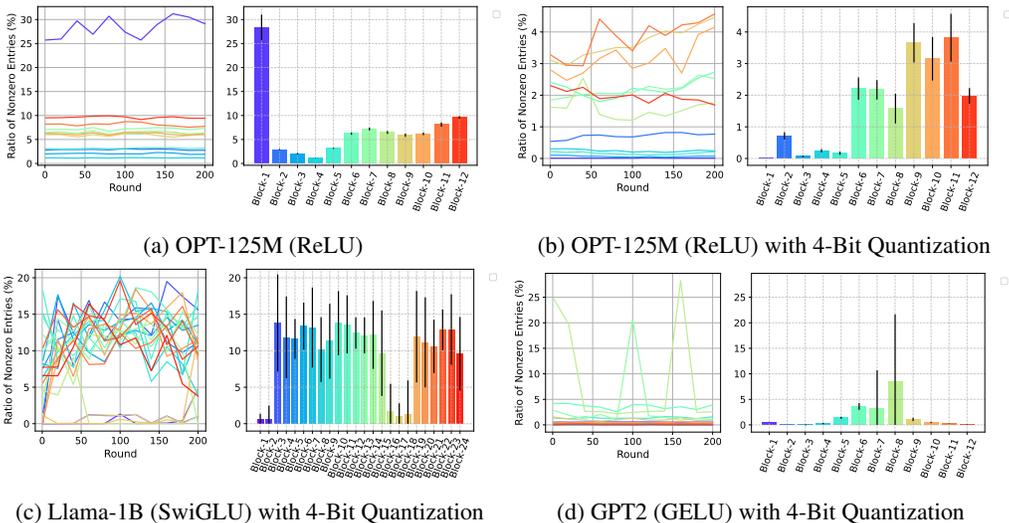


Figure 2: Sparsity Levels of Various Activations Here we do not show the sparsity level of SwiGLU and GELU because they are very dense. Specifically, their ratio of nonzero entries are 99%-100%.

### 3 SPARQ FRAMEWORK DESIGN

We introduce SparQ, a zeroth-order model-parallel framework with with Split layer allocation informed by Quantization-induced activation sparsity, reducing communication cost and minimizing memory costs for LLM fine-tuning in model parallelism. It employs ZO optimization to bypass the need for storing and transmitting backward gradients. It is crucial to highlight our split strategy that is closely related to communication reduction. By capitalizing on the quantization-induced high sparsity of activations immediately after activation functions, we partition the model at these sparse areas to naturally cut down on communication costs.

In Alg. 1, we illustrate how SparQ operates on a two-GPU setup ( $M = 2$ ) for brevity and distribute the multi-node case ( $M > 2$ ) to Appendix C.3 due to the limited space. The entire model is partitioned

**Algorithm 1** SparQ with Forward Difference Method ( $M = 2$ )

---

```

216
217
218 1: Initialize: split model immediately after activation functions to get submodels  $S_1$  and  $S_2$ , model
219   parameter  $\mathbf{x}_0 = \text{col}[\mathbf{x}_0^{(1)}, \mathbf{x}_0^{(2)}]$ , learning rate  $\eta$ , smoothing parameter  $\mu$ , iterations  $T$ .
220 2: for  $t = 0, 1, \dots, T - 1$  do
221 3:   On node  $\mathcal{M}_1$ :
222 4:     Sample a random seed  $s$  and a data sample  $\xi_t$ .
223 5:     Compute  $a_t = S_1(\xi_t; \mathbf{x}_t^{(1)})$ 
224 6:      $\mathbf{x}_t^{(1)} \leftarrow \text{Perturb}(\mathbf{x}_t^{(1)}, \mu, s)$ 
225 7:     Compute  $a_t^+ = S_1(\xi_t; \mathbf{x}_t^{(1)})$ 
226 8:      $\mathbf{x}_t^{(1)} \leftarrow \text{Perturb}(\mathbf{x}_t^{(1)}, -\mu, s)$ 
227 9:     Send  $\mathbb{C}(a_t^+)$  and  $\mathbb{C}(a_t)$ , which are sparse representations of quantized  $a_t^+$  and  $a_t$ , to  $\mathcal{M}_2$ .
228 10:  On node  $\mathcal{M}_2$ :
229 11:   Compute  $f_t = F(S_2(\mathbb{C}(a_t); \mathbf{x}_t^{(2)}))$ 
230 12:    $\mathbf{x}_t^{(2)} \leftarrow \text{Perturb}(\mathbf{x}_t^{(2)}, \mu, s)$ 
231 13:   Compute  $f_t^+ = F(S_2(\mathbb{C}(a_t^+); \mathbf{x}_t^{(2)}))$ 
232 14:    $\mathbf{x}_t^{(2)} \leftarrow \text{Perturb}(\mathbf{x}_t^{(2)}, -\mu, s)$ 
233 15:    $\hat{g}_t = \frac{1}{\mu}(f_t^+ - f_t)$ 
234 16:    $\mathbf{x}_{t+1}^{(2)} \leftarrow \text{Update}(\mathbf{x}_t^{(2)}, \eta, \hat{g}_t, s)$ 
235 17:   Send  $\hat{g}_t$  back to  $\mathcal{M}_1$ . ► Only a scalar is transmitted in backward propagation
236 18:  On node  $\mathcal{M}_1$ :
237 19:    $\mathbf{x}_{t+1}^{(1)} \leftarrow \text{Update}(\mathbf{x}_t^{(1)}, \eta, \hat{g}_t, s)$ 
238 20: end for
239
240
241 21: Function  $\text{Perturb}(\mathbf{x}, \mu, s)$ :
242 22:   Use random seed  $s$  to reset random number generator
243 23:   for  $x_i \in \mathbf{x}$  do
244 24:      $x_i \leftarrow x_i + \mu \cdot u$ , where  $u \sim \mathcal{N}(0, 1)$ .
245 25:   end for
246 26:   return  $\mathbf{x}$ 
247
248 27: Function  $\text{Update}(\mathbf{x}, \eta, \hat{g}, s)$ :
249 28:   Use random seed  $s$  to reset random number generator
250 29:   for  $x_i \in \mathbf{x}$  do
251 30:      $x_i \leftarrow x_i - \eta \cdot \hat{g} \cdot u$ , where  $u \sim \mathcal{N}(0, 1)$  ► Standard Zeroth-Order SGD
252 31:   end for
253 32:   return  $\mathbf{x}$ 

```

---

into two segments,  $S_1$  and  $S_2$ , which are deployed on two GPUs  $\mathcal{M}_1$  and  $\mathcal{M}_2$ , respectively. The detailed procedure in the  $t$ -th iteration is as follows:

- **Step 1:** On node  $\mathcal{M}_1$ , the submodel  $S_1$  computes activations  $a_t^+$  and  $a_t$  (Lines 5-7). Here, we perturb  $\mathbf{x}_t^{(1)}$  element-wise following (Malladi et al., 2023) to reduce memory usage. By default, we use a biased forward difference approach in Alg. 1 and distribute the unbiased central difference approach to Appendix C.2. To lower communication cost, we design a compressor  $\mathbb{C}$  that first applies quantization (e.g., 4 bit) to induce high sparsity in  $a_t$  and  $a_t^+$  and encodes them with sparse representations (i.e., only nonzero values and their indices). Subsequently, the compressed activations  $\mathbb{C}(a_t)$  and  $\mathbb{C}(a_t^+)$  are transmitted to the next node  $\mathcal{M}_2$ .
- **Step 2:** Upon receiving  $\mathbb{C}(a_t)$  and  $\mathbb{C}(a_t^+)$  from  $\mathcal{M}_1$ , node  $\mathcal{M}_2$  feed them into submodel  $S_2$  to compute the ZO gradient scalar  $\hat{g}_t$  (Lines 11-15). The perturbation of  $\mathbf{x}_t^{(2)}$  is performed in the same manner as in Step 1. Then, submodel parameters at  $\mathcal{M}_2$  are updated via standard ZO-SGD (Line 16). Finally, only the scalar  $\hat{g}_t$  is communicated back to  $\mathcal{M}_1$ , instead of a high-dimensional first-order gradient, yielding substantial communication savings.
- **Step 3:** Once  $\mathcal{M}_1$  receives  $\hat{g}_t$ , it updates sub-model parameters using the same ZO-SGD procedure applied on  $\mathcal{M}_2$  (Line 19).

Here, we highlight the comparison of communication costs between central and forward difference methods. Given  $P$  perturbations, in each round, utilizing the central difference method transmits  $2P$  activations and  $P$  gradient scalars, whereas utilizing the forward difference method transmits only  $P + 1$  activations and  $P$  gradient scalars. Table 2 empirically shows that the forward difference method achieves lower communication costs while maintaining comparable test accuracy compared to the central difference method.

In Fig. 3, we illustrate the execution timelines of our ZO method and conventional FO methods. The critical difference lies in the backward communication pattern: FO methods require layer-wise transmission of high-dimensional gradients during backpropagation, incurring substantial communication and time costs. In contrast, our ZO method eliminates gradient tensors entirely and instead communicates only a single scalar per iteration, which can be broadcast to all nodes simultaneously. This design not only reduces the per-iteration communication complexity from dimension-dependent to dimension-free, but also eliminates sequential gradient exchanges, thereby reducing communication cost and shortening the overall execution timeline.

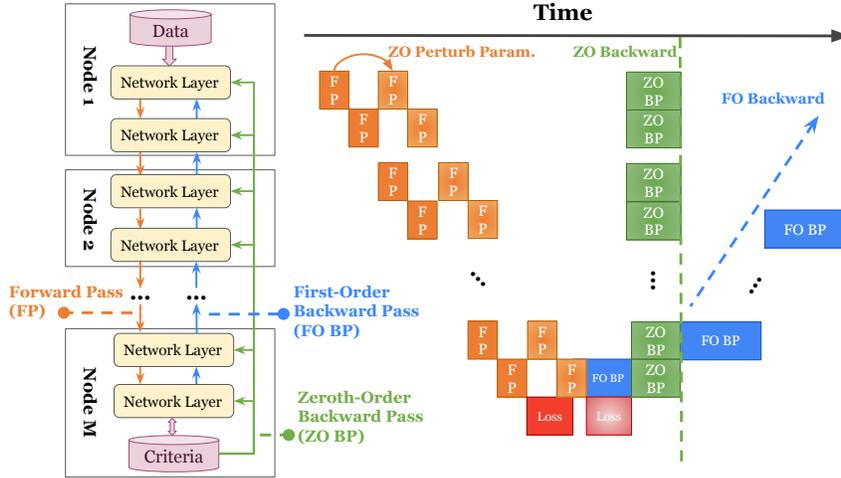


Figure 3: Execution Timeline Comparison of Zeroth-Order (e.g., SparQ) and First-Order Methods.

#### 4 THEORETICAL ANALYSIS

We start with all assumptions used in this work. Note that  $f(\cdot)$  is the loss function,  $\nabla f(\cdot)$  is the first-order gradient,  $\hat{\nabla} f(\cdot)$  is the zeroth-order gradient. For simplicity, we only show two-node case in the following main paper and distribute the  $M > 2$  case to Appendix D.3.

**Assumption 1 (Unbiased Function Estimation)**  $\mathbb{E}[f(\xi_t; \mathbf{x})] = f(\mathbf{x})$ .

**Assumption 2 (Unbiased Zeroth-Order Stochastic Gradients with Bounded Variance)** The stochastic gradient  $\hat{\nabla} f(\mathbf{x}_k; \xi_k)$  is unbiased, and its variance is bounded, so we have

$$\mathbb{E}[\hat{\nabla} f(\xi_k; \mathbf{x}_k)] = \nabla f(\mathbf{x}_k) \quad \text{and} \quad \mathbb{E}[\|\hat{\nabla} f(\xi_k; \mathbf{x}_k) - \nabla f(\mathbf{x}_k)\|^2] \leq \sigma^2. \quad (6)$$

**Assumption 3 (Gradient Lipschitz Condition)** For the loss function and each composition function, their gradients are Lipschitz-continuous such that for any  $\xi_t$ ,

$$\begin{aligned} \|\nabla f(\xi_t; \mathbf{x}) - \nabla f(\xi_t; \mathbf{y})\| &\leq L\|\mathbf{x} - \mathbf{y}\| \\ \|\nabla(S_1|_{\mathbf{x}^{(1)}})(\xi_t) - \nabla(S_1|_{\mathbf{y}^{(1)}})(\xi_t)\| &\leq L_1\|\mathbf{x}^{(1)} - \mathbf{y}^{(1)}\| \\ \|\nabla(F \circ S_2|_{\mathbf{x}^{(2)}})(\zeta_t) - \nabla(F \circ S_2|_{\mathbf{y}^{(2)}})(\zeta_t)\| &\leq L_2\|\mathbf{x}^{(2)} - \mathbf{y}^{(2)}\| \end{aligned}$$

Further, the gradients are bounded:  $\|\nabla(S_1|_{\mathbf{x}^{(1)}})(\xi_t)\| \leq C_{S_1}$  and  $\|\nabla(F \circ S_2|_{\mathbf{x}^{(2)}})(\zeta_t)\| \leq C_{S_2}$ .

**Assumption 4 (Bounded Output of Split Layer)**  $\|S_1(\xi; \mathbf{x})\| \leq L_{S_1}, \forall \mathbf{x}, \xi$ .

**Assumption 5 (Unbiased Compressor  $\mathbb{C}(\cdot)$ )**  $\|\mathbf{x} - \mathbb{C}(\mathbf{x})\| \leq \kappa\|\mathbf{x}\|$ , where  $\kappa \in [0, 1)$ .

Under the assumptions above, we derive the lemma and obtain the `SparQ`'s nonconvex convergence bound the two-node case as follows.

**Lemma 1 (Distance between Gradients of Uncompressed and Compressed Activations)** For Algorithm 1 with  $M = 2$ , the difference between gradients of uncompressed and compressed activations can be bounded as:  $\|\nabla f(\mathbf{x}_t) - \nabla \hat{f}(\mathbf{x}_t)\|^2 \leq (1 + C_{S_1}^2)L_2^2\kappa^2L_{S_1}^2$ .

**Remark 1** From Lemma 1, we observe that the upper bound is a constant, independent of learning rate  $\eta$  and smoothing parameter  $\mu$ . This constant bound could, in principle, be further reduced. For example, AQ-SGD (Wang et al., 2022) applies a error-feedback technique on the compressed activation, but its memory cost is huge since it requires storing the previous information for each data sample on both nodes, which is mostly infeasible in practice. Our paper considers the memory limitation scenario, so we opt to compress the activation directly without any extra memory cost.

**Theorem 1 (Convergence of `SparQ` under Non-Convexity,  $M = 2$ )** Under the assumptions 1, 2, 3, 4 and 5, supposing that  $\eta = \mathcal{O}(1/\sqrt{Td})$ ,  $\mu \leq 1/(d + 6)\sqrt{T}$  and  $D = f(\mathbf{x}_0) - f(\mathbf{x}^*)$ , then the sequence  $\{\mathbf{x}_t\}$  generated by `SparQ` satisfies

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(\mathbf{x}_t)\|^2 = \mathcal{O}\left(D\sqrt{d/T}\right) + \mathcal{O}\left(\sqrt{d/T}\sigma^2\right) + \mathcal{O}\left((1 + C_{S_1}^2)L_2^2\kappa^2L_{S_1}^2\sqrt{d/T}\right).$$

**Remark 2** The first term in the above bound is associated with the distance between the initial point and the optimal point. The second term is related to the variance of stochastic gradients. Both terms have the  $\mathcal{O}(\sqrt{d/T})$  convergence rate matching with the standard ZO method in the non-convex scenario. The third one is an extra term caused by the extra compression in the activation. It does not impact the overall convergence rate asymptotically. When using a lossless compressor ( $\kappa = 0$ ) like sparse representation, then the third term disappears, and `SparQ` can achieve a convergence rate of  $\mathcal{O}(\sqrt{d/T})$  under the non-convex condition.

## 5 EXPERIMENTS

**Models & Datasets.** In our evaluation, we utilize three NLP datasets: SST2 (Socher et al., 2013) for sentiment classification, WIC (Pilehvar & Camacho-Collados, 2019) for context-sensitive word embeddings evaluation and RTE (Bowman et al., 2015) for textual entailment recognition. For each of them, we fine-tune OPT-125M, OPT-1.3B, OPT-6.7B, GPT2 and Llama-1B models and monitor their test accuracy, peak GPU memory usage and communication cost.

**Baselines.** To comprehensively evaluate the performance, we compare `SparQ` with several baselines: first-order SGD (FO-SGD), AQ-SGD (Wang et al., 2022) and ZO-SGD (i.e., MeZO (Malladi et al., 2023)). FO-SGD represents a centralized first-order method. AQ-SGD represents a first-order model parallel method with activation (using error feedback) and gradient compression. MeZO represents a recent memory-efficient ZO method without considering communication cost. For compression level, we use 4-bit quantization to compress activations for AQ-SGD and `SparQ`, and we employ 8-bit quantization to compress backward gradients for AQ-SGD.

**Ablation Experiments.** 1) splitting between blocks and splitting immediately after activation functions to explore the impact of split positions in Fig. 4; 2) using different quantization degrees (e.g., 1-bit, 2-bit, 4-bit and 8-bit) to investigate the influence of quantization levels in Table. 1.

### EXPERIMENT RESULTS:

**Split Between Blocks v.s. Split Immediately After Activation Functions.** Fig. 4 shows our splitting strategy. As shown in (a)-(c), placing the split after activation functions leads to trainable models. Specifically, ReLU induces natural sparsity that is well-preserved after 4-bit quantization, while SwiGLU and GELU, though producing dense activations, still maintain sufficient information for effective training. In contrast, (d) shows that splitting between blocks severely damages the activation representation after quantization, making the model untrainable. These results demonstrate that our splitting strategy (a-c) effectively balances quantization and trainability, validating the design choice.

**`SparQ`'s Performance with Different Quantization Levels.** Table 1 reports `SparQ`'s performance under four quantization schemes: 1-bit, 2-bit, 4-bit and 8-bit, under a single split setting ( $M = 2$ ). We make two key observations: (1) stronger quantization (fewer bits) consistently reduces communication overhead; (2) accuracy degradation is relatively minor, even under aggressive quantization. It is worth

378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393

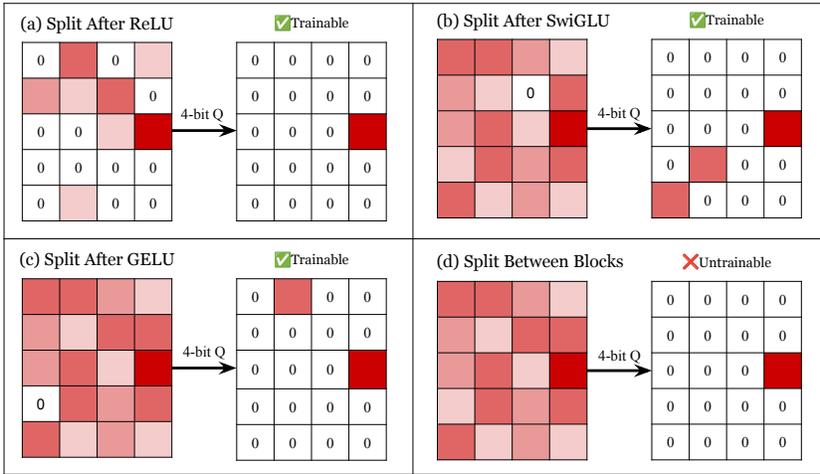


Figure 4: An Illustration of 4-Bit Quantized Activations Comparison. For each strategy, the matrix on the left represents the forward pass (e.g., activation).

394  
395  
396  
397  
398  
399  
400  
401

noting that these results reflect the impact of quantization in the case of a single model split. When the model is split into multiple segments, the effect of quantization on both accuracy and communication overhead becomes more pronounced. Balancing accuracy and efficiency, we therefore select 4-bit quantization as the default configuration in subsequent experiments.

Table 1: Activation Quantization Levels’ Impact on  $S_{parQ}$ ’s Test Accuracy and **Communication Overhead**. Hyper-parameter setup: batch size=32,  $M = 2$ .

402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412

Model	Dataset	1-bit Quantization	2-bit Quantization	4-bit Quantization	8-bit Quantization
OPT-1.3B (ReLU)	SST-2	89.84% (0.31 GB)	89.97% (0.58 GB)	92.34% (1.20 GB)	92.04% (2.32 GB)
	WIC	54.75% (0.78 GB)	55.50% (1.57 GB)	55.62% (3.23 GB)	56.07% (6.26 GB)
	RTE	57.03% (2.49 GB)	57.54% (4.34 GB)	57.13% (7.18 GB)	58.25% (14.06 GB)
GPT2 (GELU)	SST-2	84.04% (0.52 GB)	84.15% (0.91 GB)	84.82% (1.55 GB)	85.12% (2.64 GB)
	WIC	51.59% (1.18 GB)	52.20% (2.33 GB)	52.70% (3.79 GB)	53.05% (6.88 GB)
	RTE	51.26% (3.11 GB)	52.01% (5.48 GB)	52.37% (8.93 GB)	52.78% (14.72 GB)
Llama-1B (SwiGLU)	SST-2	89.32% (0.65 GB)	92.13% (1.47 GB)	93.44% (2.42 GB)	93.69% (5.64 GB)
	WIC	52.77% (1.89 GB)	53.35% (3.53 GB)	53.63% (7.31 GB)	53.82% (10.22 GB)
	RTE	53.62% (3.47 GB)	54.58% (6.83 GB)	55.21% (12.46 GB)	56.17% (24.74 GB)

413  
414  
415

**Test Accuracy of  $S_{parQ}$  with 4-Bit Quantization Matches ZO-SGD.** Table 2 indicates that FO-SGD achieves the highest test accuracy among all methods. Notably,  $S_{parQ}$  with 4-bit quantization attains performance comparable to ZO-SGD and even outperforms AQ-SGD on the SST2 and WIC.

416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426

**$S_{parQ}$  Achieves the Best Efficiency in Both Memory and Communication.** As reported in Table 2,  $S_{parQ}$  achieves the lowest overhead when considering both memory and communication costs. On the memory side, compared with FO-SGD,  $S_{parQ}$  reduces total peak usage by about 30 ~ 70%, while also avoiding the substantial extra per-sample storage required by AQ-SGD. This shows that  $S_{parQ}$  can effectively scale to larger models where FO methods often encounter memory bottlenecks. On the communication side,  $S_{parQ}$  stably outperforms both FO-SGD and AQ-SGD, with the forward-difference variant consistently achieving the lowest communication cost across all model scales. On average,  $S_{parQ}$  reduces communication overhead by more than 50% compared to AQ-SGD and over 70% compared to FO-SGD. Taken together, these results establish  $S_{parQ}$  as the most efficient method overall, striking a favorable balance between maintaining accuracy and minimizing both memory and communication demands.

427  
428  
429  
430  
431

## 6 RELATED WORK

**Model Parallelism (MP).** MP (Dean et al., 2012) is a fundamental technique in distributed deep learning that partitions a deep neural network into disjoint segments, each assigned to a separate computing node (e.g., a GPU or machine). Building on MP, various algorithms have been proposed. Among these, AQ-SGD (Wang et al., 2022) is particularly relevant to this work. It employs pipeline parallelism and can function effectively over slow networks but introduces a huge extra memory cost.

Table 2: Test Accuracy, **Memory Cost**, and **Communication Cost**. 1) Memory cost here means total peak memory usage. For AQ-SGD, the value before + is the total peak message usage on two GPUs, and the value after + is the extra memory usage to store per-sample messages, which are on SSD or CPU, as Wang et al. (2022); 2) Key hyper-parameter setup:  $P = 5$  for all ZO methods, batch size= 32; 3) "4Q" means 4-bit quantization. "forward" means forward difference method.

Model	Dataset	FO-SGD	AQ-SGD	ZO-SGD	Ours (4Q, forward)
OPT-125M	SST-2	87.5% (2 GB, 2.8 GB)	82.3% (2+11 GB, 1.1 GB)	85.2% (1 GB, 3.1 GB)	84.5% (1 GB, 0.5 GB)
	WIC	54.1% (3 GB, 7.2 GB)	53.9% (3+5 GB, 2.7 GB)	53.6% (2 GB, 8.3 GB)	53.4% (2 GB, 1.2 GB)
	RTE	56.5% (11 GB, 15.7 GB)	54.3% (11+10 GB, 5.9 GB)	53.5% (6 GB, 17.1 GB)	53.2% (6 GB, 2.6 GB)
OPT-1.3B	SST-2	91.7% (15 GB, 8 GB)	84.2% (15+31 GB, 3 GB)	90.6% (6 GB, 8 GB)	92.3% (6 GB, 1 GB)
	WIC	63.5% (16 GB, 19 GB)	56.7% (16+9 GB, 7 GB)	55.8% (7 GB, 21 GB)	55.6% (7 GB, 3 GB)
	RTE	70.8% (41 GB, 42 GB)	60.2% (41+34 GB, 16 GB)	57.3% (11 GB, 47 GB)	57.1% (11 GB, 7 GB)
OPT-6.7B	SST-2	94.4% (76 GB, 15 GB)	89.6% (76+61 GB, 6 GB)	92.1% (26 GB, 11 GB)	92.0% (26 GB, 2 GB)
	WIC	65.8% (78 GB, 39 GB)	56.8% (78+19 GB, 15 GB)	58.7% (27 GB, 31 GB)	57.8% (27 GB, 5 GB)
	RTE	71.1% (202 GB, 83 GB)	64.3% (202+69 GB, 31 GB)	63.5% (29 GB, 63 GB)	63.1% (29 GB, 9 GB)
GPT2	SST-2	88.1% (2 GB, 4 GB)	84.5% (2+10 GB, 3 GB)	84.9% (1 GB, 4 GB)	84.8% (1 GB, 2 GB)
	WIC	61.3% (2 GB, 22 GB)	55.6% (2+5 GB, 9 GB)	52.5% (1 GB, 25 GB)	52.7% (1 GB, 4 GB)
	RTE	63.1% (5 GB, 46 GB)	55.9% (5+10 GB, 20 GB)	52.3% (3 GB, 52 GB)	52.3% (3 GB, 9 GB)
Llama-1B	SST-2	94.3% (14 GB, 8 GB)	93.2% (14+30 GB, 4 GB)	93.7% (5 GB, 11 GB)	93.4% (5 GB, 2 GB)
	WIC	60.4% (16 GB, 18 GB)	55.7% (16+10 GB, 8 GB)	53.6% (6 GB, 20 GB)	53.6% (6 GB, 7 GB)
	RTE	64.7% (25 GB, 41 GB)	59.1% (25+30 GB, 17 GB)	55.8% (8 GB, 42 GB)	55.2% (8 GB, 12 GB)

**Zeroth-Order (ZO) Optimization.** ZO optimization is a gradient-free approach that estimates gradients using only differences in function values and random perturbation vectors (Liu et al., 2020), in contrast to first-order methods, which rely on explicitly computed gradients. Prior research has demonstrated the efficacy of ZO optimization in black-box attack (Kariyappa et al., 2021; Yu et al., 2024), reinforcement learning (Pan et al., 2022; Jing et al., 2024), communication savings (Fang et al., 2022; Qin et al., 2024; Li et al., 2025a;b), etc. In addition, ZO optimization has been shown to lower memory consumption during LLM fine-tuning. For example, MeZO (Malladi et al., 2023) and LOZO (Chen et al., 2025) utilize ZO optimization to perform solely forward passes, thereby eliminating the need to store gradients from backward propagation. Yet, the integration of ZO optimization within MP frameworks remains largely unexplored.

**Activation Compression & Sparsity.** Compression is widely adopted to reduce communication and memory overhead in distributed systems. Popular compression techniques, such as quantization (Gray & Neuhoff, 1998; Alistarh et al., 2017; Horváth et al., 2023; Lin et al., 2024) and sparsification (Wangni et al., 2018; Yang et al., 2021; Yoon & Oh, 2023), have primarily targeted gradients and weights. Recently, attention has shifted toward compressing activations (Evans & Aamodt, 2021; Liu et al., 2021; 2022; Chen et al., 2021; Eliassen & Selvan, 2024; Lin et al., 2024). For example, ActNN (Chen et al., 2021) and ALAM (Woo et al., 2024) quantize activations to improve memory efficiency. Moreover, AQ-SGD (Wang et al., 2022) quantizes backward gradients and the changes of forward activation to enable communication-efficient training in pipeline parallel architectures. However, its reliance on error-feedback mechanism to guarantee convergence necessitates storing per-sample information on CPUs or SSDs, thereby imposing a significant extra memory burden.

## 7 CONCLUSION

In summary, we introduce `SparQ`, a ZO model-parallel framework with split layer allocation informed by quantization-induced activation sparsity, reducing communication overhead and minimizes memory costs for LLM fine-tuning under MP. Our theoretical analysis establishes a sublinear convergence rate in non-convex settings, and empirical results show that `SparQ` reduces GPU memory consumption by more than  $3\times$  and communication cost by over 50% compared to state-of-the-art baselines. These findings highlight the potential of ZO, sparsity-guided frameworks for scaling LLM fine-tuning, and open up promising directions for future research on distributed optimization.

## REFERENCES

- 486  
487  
488 Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. Qsgd: Communication-  
489 efficient sgd via gradient quantization and encoding. *Advances in Neural Information Processing*  
490 *Systems*, 30, 2017.
- 491 Samuel Bowman, Gabor Angeli, Christopher Potts, and Christopher D Manning. A large annotated  
492 corpus for learning natural language inference. In *Proceedings of the 2015 Conference on Empirical*  
493 *Methods in Natural Language Processing*, pp. 632–642, 2015.
- 494 Jianfei Chen, Lianmin Zheng, Zhewei Yao, Dequan Wang, Ion Stoica, Michael Mahoney, and Joseph  
495 Gonzalez. Actnn: Reducing training memory footprint via 2-bit activation compressed training. In  
496 *International Conference on Machine Learning*, pp. 1803–1813. PMLR, 2021.
- 497 Yiming Chen, Yuan Zhang, Liyuan Cao, Kun Yuan, and Zaiwen Wen. Enhancing zeroth-order  
498 fine-tuning for language models with low-rank structures. In *The Thirteenth International Confer-*  
499 *ence on Learning Representations*, 2025. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=9BiVepgmWW)  
500 [9BiVepgmWW](https://openreview.net/forum?id=9BiVepgmWW).
- 501 Jeffrey Dean, Greg Corrado, Rajat Monga, Kai Chen, Matthieu Devin, Mark Mao, Marc’ aurelio  
502 Ranzato, Andrew Senior, Paul Tucker, Ke Yang, et al. Large scale distributed deep networks.  
503 *Advances in Neural Information Processing Systems*, 25, 2012.
- 504 Sebastian Eliassen and Raghavendra Selvan. Activation compression of graph neural networks  
505 using block-wise quantization with improved variance minimization. In *ICASSP 2024-2024 IEEE*  
506 *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 7430–7434.  
507 IEEE, 2024.
- 508 R David Evans and Tor Aamodt. Ac-gc: Lossy activation compression with guaranteed convergence.  
509 *Advances in Neural Information Processing Systems*, 34:27434–27448, 2021.
- 510 Wenzhi Fang, Ziyi Yu, Yuning Jiang, Yuanming Shi, Colin N Jones, and Yong Zhou. Communication-  
511 efficient stochastic zeroth-order optimization for federated learning. *IEEE Transactions on Signal*  
512 *Processing*, 70:5058–5073, 2022.
- 513 Saeed Ghadimi and Guanghui Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic  
514 programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013.
- 515 Robert M. Gray and David L. Neuhoff. Quantization. *IEEE Transactions on Information Theory*, 44  
516 (6):2325–2383, 1998.
- 517 Zeyu Han, Chao Gao, Jinyang Liu, Jeff Zhang, and Sai Qian Zhang. Parameter-efficient fine-tuning  
518 for large models: A comprehensive survey. *Transactions on Machine Learning Research*, 2024.  
519 ISSN 2835-8856. URL <https://openreview.net/forum?id=1IsCS8b6zj>.
- 520 Samuel Horváth, Dmitry Kovalev, Konstantin Mishchenko, Peter Richtárik, and Sebastian Stich.  
521 Stochastic distributed learning with gradient quantization and double-variance reduction. *Opti-*  
522 *mization Methods and Software*, 38(1):91–106, 2023.
- 523 Ziheng Jiang, Haibin Lin, Yinmin Zhong, Qi Huang, Yangrui Chen, Zhi Zhang, Yanghua Peng,  
524 Xiang Li, Cong Xie, Shibiao Nong, et al. {MegaScale}: Scaling large language model training  
525 to more than 10,000 {GPUs}. In *21st USENIX Symposium on Networked Systems Design and*  
526 *Implementation (NSDI 24)*, pp. 745–760, 2024.
- 527 Gangshan Jing, He Bai, Jemin George, Aranya Chakraborty, and Piyush K Sharma. Asynchronous  
528 distributed reinforcement learning for lqr control via zeroth-order block coordinate descent. *IEEE*  
529 *Transactions on Automatic Control*, 2024.
- 530 Jean Kaddour, Joshua Harris, Maximilian Mozes, Herbie Bradley, Roberta Raileanu, and Robert  
531 McHardy. Challenges and applications of large language models. *arXiv preprint arXiv:2307.10169*,  
532 2023.
- 533 Sanjay Kariyappa, Atul Prakash, and Moinuddin K Qureshi. Maze: Data-free model stealing attack  
534 using zeroth-order gradient estimation. In *Proceedings of the IEEE/CVF Conference on Computer*  
535 *Vision and Pattern Recognition*, pp. 13814–13823, 2021.

- 540 Zhe Li, Bicheng Ying, Zidong Liu, Chaosheng Dong, and Haibo Yang. Achieving dimension-free  
541 communication in federated learning via zeroth-order optimization. In *The Thirteenth International*  
542 *Conference on Learning Representations*, 2025a.
- 543
- 544 Zhe Li, Bicheng Ying, Zidong Liu, Chaosheng Dong, and Haibo Yang. Reconciling hessian-informed  
545 acceleration and scalar-only communication for efficient federated zeroth-order fine-tuning. *arXiv*  
546 *preprint arXiv:2506.02370*, 2025b.
- 547
- 548 Zonglin Li, Chong You, Srinadh Bhojanapalli, Daliang Li, Ankit Singh Rawat, Sashank J. Reddi,  
549 Ke Ye, Felix Chern, Felix Yu, Ruiqi Guo, and Sanjiv Kumar. The lazy neuron phenomenon:  
550 On emergence of activation sparsity in transformers. In *The Eleventh International Confer-*  
551 *ence on Learning Representations*, 2023. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=TJ2nxc1Yck-)  
552 [TJ2nxc1Yck-](https://openreview.net/forum?id=TJ2nxc1Yck-).
- 553 Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Wei-Ming Chen, Wei-Chen Wang, Guangxuan  
554 Xiao, Xingyu Dang, Chuang Gan, and Song Han. Awq: Activation-aware weight quantization for  
555 on-device llm compression and acceleration. *Proceedings of Machine Learning and Systems*, 6:  
556 87–100, 2024.
- 557
- 558 Sijia Liu, Pin-Yu Chen, Bhavya Kailkhura, Gaoyuan Zhang, Alfred O Hero III, and Pramod K  
559 Varshney. A primer on zeroth-order optimization in signal processing and machine learning:  
560 Principals, recent advances, and applications. *IEEE Signal Processing Magazine*, 37(5):43–54,  
561 2020.
- 562 Xiaoxuan Liu, Lianmin Zheng, Dequan Wang, Yukuo Cen, Weize Chen, Xu Han, Jianfei Chen,  
563 Zhiyuan Liu, Jie Tang, Joey Gonzalez, et al. Gact: Activation compressed training for generic  
564 network architectures. In *International Conference on Machine Learning*, pp. 14139–14152.  
565 PMLR, 2022.
- 566
- 567 Zirui Liu, Kaixiong Zhou, Fan Yang, Li Li, Rui Chen, and Xia Hu. Exact: Scalable graph neural  
568 networks training via extreme activation compression. In *International Conference on Learning*  
569 *Representations*, 2021.
- 570
- 571 Sadhika Malladi, Tianyu Gao, Eshaan Nichani, Alex Damian, Jason D Lee, Danqi Chen, and Sanjeev  
572 Arora. Fine-tuning language models with just forward passes. *Advances in Neural Information*  
573 *Processing Systems*, 36:53038–53075, 2023.
- 574
- 575 Humza Naveed, Asad Ullah Khan, Shi Qiu, Muhammad Saqib, Saeed Anwar, Muhammad Usman,  
576 Naveed Akhtar, Nick Barnes, and Ajmal Mian. A comprehensive overview of large language  
577 models. *ACM Transactions on Intelligent Systems and Technology*, 16(5):1–72, 2025.
- 578
- 579 Yurii Nesterov and Vladimir Spokoiny. Random gradient-free minimization of convex functions.  
580 *Foundations of Computational Mathematics*, 17:527–566, 2017.
- 581
- 582 Ling Pan, Longbo Huang, Tengyu Ma, and Huazhe Xu. Plan better amid conservatism: Offline multi-  
583 agent reinforcement learning with actor rectification. In *International Conference on Machine*  
584 *Learning*, pp. 17221–17237. PMLR, 2022.
- 585
- 586 Mohammad Taher Pilehvar and Jose Camacho-Collados. Wic: the word-in-context dataset for  
587 evaluating context-sensitive meaning representations. In *Proceedings of the 2019 Conference of*  
588 *the North American Chapter of the Association for Computational Linguistics: Human Language*  
589 *Technologies, Volume 1 (Long and Short Papers)*, pp. 1267–1273, 2019.
- 590
- 591 Zhen Qin, Daoyuan Chen, Bingchen Qian, Bolin Ding, Yaliang Li, and Shuiguang Deng. Federated  
592 full-parameter tuning of billion-sized language models with communication cost under 18 kilobytes.  
593 In *International Conference on Machine Learning*, pp. 41473–41497. PMLR, 2024.
- 594
- 595 Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. Zero: Memory optimizations  
596 toward training trillion parameter models. In *SC20: International Conference for High Performance*  
597 *Computing, Networking, Storage and Analysis*, pp. 1–16. IEEE, 2020.

- 594 Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. DeepSpeed: System optimiza-  
595 tions enable training deep learning models with over 100 billion parameters. In *Proceedings of*  
596 *the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp.  
597 3505–3506, 2020.
- 598 Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catan-  
599 zaro. Megatron-lm: Training multi-billion parameter language models using model parallelism.  
600 *arXiv preprint arXiv:1909.08053*, 2019.
- 601 Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng,  
602 and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment  
603 treebank. In *Proceedings of The 2013 Conference on Empirical Methods in Natural Language*  
604 *Processing*, pp. 1631–1642, 2013.
- 605 A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017.
- 606 Jue Wang, Binhang Yuan, Luka Rimanic, Yongjun He, Tri Dao, Beidi Chen, Christopher Ré, and  
607 Ce Zhang. Fine-tuning language models over slow networks using activation quantization with  
608 guarantees. *Advances in Neural Information Processing Systems*, 35:19215–19230, 2022.
- 609 Jianqiao Wangni, Jialei Wang, Ji Liu, and Tong Zhang. Gradient sparsification for communication-  
610 efficient distributed optimization. *Advances in Neural Information Processing Systems*, 31, 2018.
- 611 Sunghyeon Woo, Sunwoo Lee, and Dongsuk Jeon. ALAM: Averaged low-precision activation for  
612 memory-efficient training of transformer models. In *The Twelfth International Conference on Learning*  
613 *Representations*, 2024. URL <https://openreview.net/forum?id=OfXqQ5TRwp>.
- 614 Haibo Yang, Jia Liu, and Elizabeth S Bentley. Cfdavg: achieving efficient communication and fast  
615 convergence in non-iid federated learning. In *2021 19th International Symposium on Modeling*  
616 *and Optimization in Mobile, Ad hoc, and Wireless Networks (WiOpt)*, pp. 1–8. IEEE, 2021.
- 617 Daegun Yoon and Sangyoon Oh. Micro: Near-zero cost gradient sparsification for scaling and  
618 accelerating distributed dnn training. In *2023 IEEE 30th International Conference on High*  
619 *Performance Computing, Data, and Analytics (HiPC)*, pp. 87–96. IEEE, 2023.
- 620 Hao Yu, Ke Liang, Dayu Hu, Wenxuan Tu, Chuan Ma, Sihang Zhou, and Xinwang Liu. Chongqing  
621 gzoo: Black-box node injection attack on graph neural networks via zeroth-order optimization.  
622 *IEEE Transactions on Knowledge and Data Engineering*, 2024.
- 623 Yihua Zhang, Pingzhi Li, Junyuan Hong, Jiayang Li, Yimeng Zhang, Wenqing Zheng, Pin-Yu  
624 Chen, Jason D Lee, Wotao Yin, Mingyi Hong, et al. Revisiting zeroth-order optimization for  
625 memory-efficient llm fine-tuning: A benchmark. *arXiv preprint arXiv:2402.11592*, 2024.
- 626 Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min,  
627 Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv*  
628 *preprint arXiv:2303.18223*, 2023.
- 629 Yanjun Zhao, Sizhe Dang, Haishan Ye, Guang Dai, Yi Qian, and Ivor Tsang. Second-order fine-tuning  
630 without pain for LLMs: A hessian informed zeroth-order optimizer. In *The Thirteenth International*  
631 *Conference on Learning Representations*, 2025. URL [https://openreview.net/forum?](https://openreview.net/forum?id=bEqI61iBue)  
632 [id=bEqI61iBue](https://openreview.net/forum?id=bEqI61iBue).
- 633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647

648	APPENDIX	
649		
650	CONTENT	
651		
652		
653	<b>A Statement of LLM Usage</b>	<b>14</b>
654		
655	<b>B Additional Experiment Details and Results</b>	<b>14</b>
656	B.1 Split Layer Selection . . . . .	14
657	B.2 GPU Peak Memory Usage . . . . .	14
658		
659		
660	<b>C sparQ’s Technical Details and Framework Extensions</b>	<b>14</b>
661	C.1 Why Activation Sparsity can be Helpful to Reduce Communication Cost? . . . . .	14
662	C.2 sparQ with Central Difference Method . . . . .	15
663	C.3 sparQ Across Multiple Computing Nodes ( $M > 2$ ) . . . . .	15
664		
665		
666		
667	<b>D Proof</b>	<b>15</b>
668	D.1 Lemmas . . . . .	15
669	D.2 Proof of Theorem 1 . . . . .	18
670	D.3 Proof of Theorem 2 . . . . .	21
671		
672		
673		
674		
675		
676		
677		
678		
679		
680		
681		
682		
683		
684		
685		
686		
687		
688		
689		
690		
691		
692		
693		
694		
695		
696		
697		
698		
699		
700		
701		

## A STATEMENT OF LLM USAGE

During the preparation of this manuscript, we used LLM tools to help with language refinement and stylistic improvements. After each use of the tool, we carefully reviewed and validated the correctness and appropriateness of the generated text to ensure accuracy and alignment with the intended meaning.

## B ADDITIONAL EXPERIMENT DETAILS AND RESULTS

### B.1 SPLIT LAYER SELECTION

In the two-node ( $M = 2$ ) setting, to balance the computation and memory burden across the two nodes, we select the split layer near the middle block where the quantized output exhibits the highest sparsity. Accordingly, in our experiments, we split the OPT models immediately after the activation function in the middle block (e.g., Block 6 for OPT-125M), the Llama-1B model after the activation function in Block 15, and the GPT-2 model after the activation function in Block 7.

### B.2 GPU PEAK MEMORY USAGE

We execute a group of experiments to test GPU peak memory usages by fine-tuning the SST2 dataset across OPT-1.3B, OPT-2.7B, OPT-6.7B and OPT-13B models. Except for different model sizes, all hyperparameter setups are the same. Fig. 5 reveals several key empirical observations:

- 1) The memory overhead for first-order methods is over  $3\times$  more than that for memory-efficient zeroth-order methods (e.g., MeZO and SparQ).
- 2) The memory cost for SparQ is approximately equal to the model size, matching the conclusion in (Malladi et al., 2023) and showing the significant memory reduction compared with first-order methods.

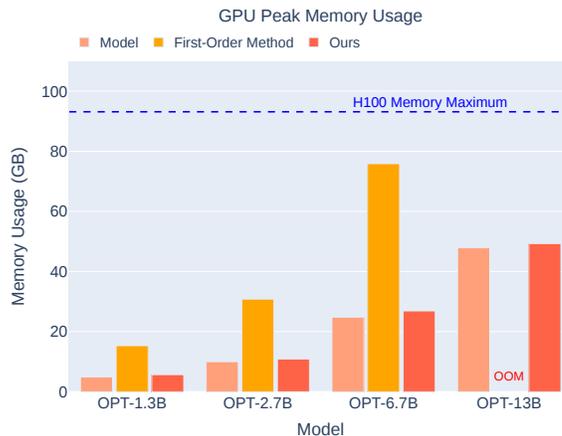


Figure 5: GPU Peak Memory Usage across Multiple LLMs When the Fine-tuning SST-2 Dataset. "OOM" means out of memory. Experiment setup: train batch size=test batch size= 32, momentum= 0.

## C SPARQ'S TECHNICAL DETAILS AND FRAMEWORK EXTENSIONS

### C.1 WHY ACTIVATION SPARSITY CAN BE HELPFUL TO REDUCE COMMUNICATION COST?

The key intuition behind using sparse representations is that activations after quantization often contain a large fraction of zeros. Instead of transmitting the full dense tensors, we can **encode only the non-zero values together with their corresponding indices**. This compressed form eliminates the need to communicate redundant zero entries, thereby drastically reducing the volume

of data transmitted across devices. From a theoretical perspective, if the activation sparsity ratio is  $\rho$  (i.e., only  $\rho$  fraction of entries remain non-zero), then the effective communication cost scales proportionally to  $\rho d$  rather than the full dimensionality  $d$ . In practice, modern LLM activations exhibit high sparsity under low-bit quantization, which makes sparse representations particularly effective. By combining quantization-induced sparsity with index-based encoding, SparQ achieves highly compressed communication while preserving task accuracy.

## C.2 SPARQ WITH CENTRAL DIFFERENCE METHOD

In our main paper, we focus on the biased forward difference method to estimate zeroth-order gradients because of its advantage of less activation transmission and comparable precision performance. Here, we introduce another widely used gradient estimation approach - the central difference method, which estimates the ZO gradient as:

$$\hat{G} = \frac{1}{P} \sum_{i=1}^P g_i \cdot \mathbf{u}_i = \frac{1}{P} \sum_{i=1}^P \frac{f(\xi; \mathbf{x} + \mu \mathbf{u}_i) - f(\xi; \mathbf{x} - \mu \mathbf{u}_i)}{2\mu} \cdot \mathbf{u}_i, \quad (7)$$

Additionally, when integrating ZOO using the unbiased central difference method into model parallelism formulation, the expression in Eq. (4) of computing the zeroth-order gradient scalar will be replaced by

$$g_{i,t} = \frac{F(a_{i,t}^+; \mathbf{x}_t^{(2)} + \mu \mathbf{u}_{i,t}^{(2)}) - F(a_{i,t}^-; \mathbf{x}_t^{(2)} - \mu \mathbf{u}_{i,t}^{(2)})}{2\mu}, \quad (8)$$

where activations are computed by

$$a_{i,t}^+ = S_1(\xi_t; \mathbf{x}_t^{(1)} + \mu \mathbf{u}_{i,t}^{(1)}), \quad a_{i,t}^- = S_1(\xi_t; \mathbf{x}_t^{(1)} - \mu \mathbf{u}_{i,t}^{(1)}). \quad (9)$$

In our experiments, the main performance difference of forward and central difference methods lies in communication overhead. Overall, central difference incurs higher communication cost compared to forward difference, as analyzed in Sec. C.3.

## C.3 SPARQ ACROSS MULTIPLE COMPUTING NODES ( $M > 2$ )

Our framework can also be straightforwardly extended to multiple computing node scenarios ( $M > 2$ ). Alg. 2 demonstrates SparQ using the central difference method under the multi-node case ( $M > 2$ ), and Alg. 3 shows SparQ using the forward difference method under the multi-node case ( $M > 2$ ). In Sec. D.3, we provide the convergence theorem and its proof of the  $M > 2$  case. The main impact of LLM fine-tuning on multiple computing nodes is the increased communication overhead attributable to the additional split layers. Nevertheless, compared with AQ-SGD (Wang et al., 2022), SparQ still can achieving  $1 \sim 2\times$  communication savings.

**Communication Cost Analysis ( $M > 2$ ).** Given the number of computing nodes  $M$ , and the number of perturbations  $P$ , the central difference method requires transmitting  $2 \times P \times (M - 1)$  activations and  $P \times (M - 1)$  gradient scalars per training round, whereas the forward difference method requires transmitting only  $(P + 1) \times (M - 1)$  activations and  $P \times (M - 1)$  gradient scalars per round. Therefore, in general, the communication overhead of using the forward difference method is lower than the cost of using the central difference method.

# D PROOF

## D.1 LEMMAS

The following Lemma 2 and Lemma 3 are relative to zeroth-order optimization and have been commonly used in zeroth-order proof. It is worth pointing out that Lemma 3 is dependent on the forward difference method, which can be known by (13). Consequently, we use it to prove the convergence of Alg. 1 ( $M = 2$ ) and Alg. 3 ( $M > 2$ ).

**Lemma 2**  $f \in C_L^{1,1}(\mathbb{R}^d)$  if  $f$  is differentiable and satisfies

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|. \quad (10)$$

**Algorithm 2** SparQ ( $M > 2$ ) with Central Difference Method

---

```

810 1: Initialize: split model immediately after activation functions to get submodels  $S_1, \dots, S_M$ , model
811 parameter  $\mathbf{x}_0 = \text{col}[\mathbf{x}_0^{(1)}, \dots, \mathbf{x}_0^{(M)}]$ , learning rate  $\eta$ , smoothing parameter  $\mu$ , the number of
812 iterations  $T$ .
813
814 2: for  $t = 0, 1, \dots, T - 1$  do
815 3:   On node  $\mathcal{M}_1$ :
816 4:     Sample a random seed  $s$  and a data sample  $\xi_t$ 
817 5:      $\mathbf{x}_t^{(1)} \leftarrow \text{Perturb}(\mathbf{x}_t^{(1)}, \mu, s)$ 
818 6:     Compute  $a_t^+ = S_1(\xi_t; \mathbf{x}_t^{(1)})$ 
819 7:      $\mathbf{x}_t^{(1)} \leftarrow \text{Perturb}(\mathbf{x}_t^{(1)}, -2\mu, s)$ 
820 8:     Compute  $a_t^- = S_1(\xi_t; \mathbf{x}_t^{(1)})$ 
821 9:      $\mathbf{x}_t^{(1)} \leftarrow \text{Perturb}(\mathbf{x}_t^{(1)}, \mu, s)$ 
822 10:    Send  $\mathbb{C}(a_t^+)$  and  $\mathbb{C}(a_t^-)$ , which are sparse representations of quantized  $a_t^+$  and  $a_t^-$ , to  $\mathcal{M}_2$ .
823
824 11:    Here we skip the description of all processes on  $\mathcal{M}_2 \cdots \mathcal{M}_{M-1}$  because they have the
825 similar process.
826
827 12:    On node  $\mathcal{M}_M$ :
828 13:      $\mathbf{x}_t^{(M)} \leftarrow \text{Perturb}(\mathbf{x}_t^{(M)}, \mu, s)$ 
829 14:     Compute  $f_t^+ = F(S_M(\mathbb{C}(a_t^+); \mathbf{x}_t^{(M)}))$ 
830 15:      $\mathbf{x}_t^{(M)} \leftarrow \text{Perturb}(\mathbf{x}_t^{(M)}, -2\mu, s)$ 
831 16:     Compute  $f_t^- = F(S_M(\mathbb{C}(a_t^-); \mathbf{x}_t^{(M)}))$ 
832 17:      $\mathbf{x}_t^{(M)} \leftarrow \text{Perturb}(\mathbf{x}_t^{(M)}, \mu, s)$ 
833 18:      $\hat{g}_t = \frac{1}{2\mu}(f_t^+ - f_t^-)$ 
834 19:      $\mathbf{x}_{t+1}^{(M)} \leftarrow \text{Update}(\mathbf{x}_t^{(M)}, \eta, \hat{g}_t, s)$ 
835 20:     Send  $\hat{g}_t$  back to  $\mathcal{M}_i$ , where  $i \in [1, \dots, M - 1]$ . ► Only a scalar is transmitted
836 21:     On nodes  $\mathcal{M}_i, i \in [1, \dots, M - 1]$ :
837 22:      $\mathbf{x}_{t+1}^{(i)} \leftarrow \text{Update}(\mathbf{x}_t^{(i)}, \eta, \hat{g}_t, s)$ 
838 23:   end for
839
840 24: Function  $\text{Perturb}(\mathbf{x}, \mu, s)$ :
841 25:   Use random seed  $s$  to reset random number generator
842 26:   for  $x_i \in \mathbf{x}$  do
843 27:      $x_i \leftarrow x_i + \mu \cdot u$ , where  $u \sim \mathcal{N}(0, 1)$ .
844 28:   end for
845 29:   return  $\mathbf{x}$ 
846
847 30: Function  $\text{Update}(\mathbf{x}, \eta, \hat{g}, s)$ :
848 31:   Use random seed  $s$  to reset random number generator
849 32:   for  $x_i \in \mathbf{x}$  do
850 33:      $x_i \leftarrow x_i - \eta \cdot \hat{g} \cdot u$ , where  $u \sim \mathcal{N}(0, 1)$ . ► Standard Zeroth-Order SGD
851 34:   end for
852 35:   return  $\mathbf{x}$ 

```

---

Also, we have

$$|f(\mathbf{y}) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle| \leq \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2. \quad (11)$$

**Lemma 3** (Ghadimi & Lan, 2013; Nesterov & Spokoiny, 2017) We define a smooth approximation of objective function  $f$  as  $f_\mu(\cdot)$  that can be formulated as

$$f_\mu(\mathbf{x}) := \frac{1}{(2\pi)^{\frac{d}{2}}} \int f(\mathbf{x} + \mu \mathbf{u}) e^{-\frac{1}{2} \|\mathbf{u}\|^2} d\mathbf{u} = \mathbb{E}[f(\mathbf{x} + \mu \mathbf{u})] \quad (12)$$

Then, for any  $f \in \mathcal{C}_L^{1,1}$ , the following statements hold.

**Algorithm 3** SparQ ( $M > 2$ ) with Forward Difference Method

---

```

864
865
866 1: Initialize: split model immediately after activation functions to get submodels  $S_1, \dots, S_M$ , model
867 parameter  $\mathbf{x}_0 = \text{col}[\mathbf{x}_0^{(1)}, \dots, \mathbf{x}_0^{(M)}]$ , learning rate  $\eta$ , smoothing parameter  $\mu$ , the number of
868 iterations  $T$ .
869
870 2: for  $t = 0, 1, \dots, T - 1$  do
871 3:   On node  $\mathcal{M}_1$ :
872 4:     Sample a random seed  $s$  and a data sample  $\xi_t$ 
873 5:     Compute  $a_t = S_1(\xi_t; \mathbf{x}_t^{(1)})$ 
874 6:      $\mathbf{x}_t^{(1)} \leftarrow \text{Perturb}(\mathbf{x}_t^{(1)}, \mu, s)$ 
875 7:     Compute  $a_t^+ = S_1(\xi_t; \mathbf{x}_t^{(1)})$ 
876 8:      $\mathbf{x}_t^{(1)} \leftarrow \text{Perturb}(\mathbf{x}_t^{(1)}, -\mu, s)$ 
877 9:     Send  $\mathbb{C}(a_t^+)$  and  $\mathbb{C}(a_t)$ , which are sparse representations of quantized  $a_t^+$  and  $a_t$ , to  $\mathcal{M}_2$ .
878 10:    Here we skip the description of all processes on  $\mathcal{M}_2 \cdots \mathcal{M}_{M-1}$  because they have the
879 similar process.
880 11:    On node  $\mathcal{M}_M$ :
881 12:      Compute  $f_t = F(S_M(\mathbb{C}(a_t); \mathbf{x}_t^{(M)}))$ 
882 13:       $\mathbf{x}_t^{(M)} \leftarrow \text{Perturb}(\mathbf{x}_t^{(M)}, \mu, s)$ 
883 14:      Compute  $f_t^+ = F(S_M(\mathbb{C}(a_t^+); \mathbf{x}_t^{(M)}))$ 
884 15:       $\mathbf{x}_t^{(M)} \leftarrow \text{Perturb}(\mathbf{x}_t^{(M)}, -\mu, s)$ 
885 16:       $\hat{g}_t = \frac{1}{\mu}(f_t^+ - f_t)$ 
886 17:       $\mathbf{x}_{t+1}^{(M)} \leftarrow \text{Update}(\mathbf{x}_t^{(M)}, \eta, \hat{g}_t, s)$ 
887 18:      Send  $\hat{g}_t$  back to  $\mathcal{M}_i$ , where  $i \in [1, \dots, M - 1]$ . ► Only a scalar is transmitted
888 19:      On nodes  $\mathcal{M}_i$ ,  $i \in [1, \dots, M - 1]$ :
889 20:         $\mathbf{x}_{t+1}^{(i)} \leftarrow \text{Update}(\mathbf{x}_t^{(i)}, \eta, \hat{g}_t, s)$ 
890 21:    end for
891
892 22: Function  $\text{Perturb}(\mathbf{x}, \mu, s)$ :
893 23:   Use random seed  $s$  to reset random number generator
894 24:   for  $x_i \in \mathbf{x}$  do
895 25:      $x_i \leftarrow x_i + \mu \cdot u$ , where  $u \sim \mathcal{N}(0, 1)$ .
896 26:   end for
897 27:   return  $\mathbf{x}$ 
898
899 28: Function  $\text{Update}(\mathbf{x}, \eta, \hat{g}, s)$ :
900 29:   Use random seed  $s$  to reset random number generator
901 30:   for  $x_i \in \mathbf{x}$  do
902 31:      $x_i \leftarrow x_i - \eta \cdot \hat{g} \cdot u$ , where  $u \sim \mathcal{N}(0, 1)$ . ► Standard Zeroth-Order SGD
903 32:   end for
904 33:   return  $\mathbf{x}$ 

```

---

(a) The gradient of  $f_\mu(\cdot)$  is  $L_\mu$ -Lipschitz continuous where  $L_\mu \leq L$ .  $\nabla f_\mu(\mathbf{x})$  can be shown as

$$\nabla f_\mu(\mathbf{x}) = \frac{1}{(2\pi)^{\frac{d}{2}}} \int \frac{f(\mathbf{x} + \mu \mathbf{u}) - f(\mathbf{x})}{\mu} \mathbf{u} e^{-\frac{1}{2}\|\mathbf{u}\|^2} d\mathbf{u}. \quad (13)$$

(b) For any  $\mathbf{x} \in \mathbb{R}^n$ ,

$$|f_\mu(\mathbf{x}) - f(\mathbf{x})| \leq \frac{\mu^2}{2} Ld \quad (14)$$

$$\|\nabla f_\mu(\mathbf{x}) - \nabla f(\mathbf{x})\| \leq \frac{\mu}{2} L(d+3)^{\frac{3}{2}} \quad (15)$$

(c) For any  $\mathbf{x} \in \mathbb{R}^n$ ,

$$\frac{1}{\mu^2} \mathbb{E}_{\mathbf{u}} \left[ \left( f(\mathbf{x} + \mu \mathbf{u}) - f(\mathbf{x}) \right)^2 \|\mathbf{u}\|^2 \right] \leq \frac{\mu^2}{2} L^2 (d+6)^3 + 2(d+4) \|\nabla f(\mathbf{x})\|^2 \quad (16)$$

Following form (15) and utilizing Jensen's inequality  $\|a\|^2 \leq 2\|a - b\|^2 + 2\|b\|^2$ , we have

$$\|\nabla f_\mu(\mathbf{x})\|^2 \leq 2\|\nabla f(\mathbf{x})\|^2 + \frac{\mu^2}{2}L^2(d+3)^3, \quad (17)$$

$$\|\nabla f(\mathbf{x})\|^2 \leq 2\|\nabla f_\mu(\mathbf{x})\|^2 + \frac{\mu^2}{2}L^2(d+3)^3. \quad (18)$$

Moreover, we denote  $f_\mu^* := \min_{\mathbf{x} \in \mathbb{R}^d} f_\mu(\mathbf{x})$  and conclude  $|f_\mu^* - f^*| \leq \frac{\mu^2 L d}{2}$  from (14). Then, we further conclude that

$$-\mu^2 L d \leq [f_\mu(\mathbf{x}) - f_\mu^*] - [f(\mathbf{x}) - f^*] \leq \mu^2 L d \quad (19)$$

**Lemma 4 (Distance between Gradients of Uncompressed and Compressed Activations,  $M = 2$ )**  
For Alg. 1, the difference between gradients of uncompressed and compressed activations can be bounded by a constant term as follows:

$$\left\| \nabla f(\mathbf{x}_t) - \nabla \hat{f}(\mathbf{x}_t) \right\|^2 \leq (1 + C_{S_1}^2) L_2^2 \kappa^2 L_{S_1}^2.$$

*Proof of Lemma 4:*

$$\begin{aligned} \left\| \nabla f(\mathbf{x}_t) - \nabla \hat{f}(\mathbf{x}_t) \right\|^2 &= \left\| \nabla_{\mathbf{x}^{(1)}} f(\mathbf{x}_t) - \nabla_{\mathbf{x}^{(1)}} \hat{f}(\mathbf{x}_t) \right\|^2 + \left\| \nabla_{\mathbf{x}^{(2)}} f(\mathbf{x}_t) - \nabla_{\mathbf{x}^{(2)}} \hat{f}(\mathbf{x}_t) \right\|^2 \\ &= \left\| \nabla_{\mathbf{x}^{(1)}} (F \circ S_2 \circ S_1) \Big|_{(\mathbf{x}_t^{(1)}, \mathbf{x}_t^{(2)})} - \nabla_{S_1} (F \circ S_2) \Big|_{(\mathbb{C}(S_1(\xi; \mathbf{x}_t^{(1)}), \mathbf{x}_t^{(2)}))} \cdot \nabla_{\mathbf{x}^{(1)}} S_1 \Big|_{\mathbf{x}_t^{(1)}} \right\|^2 \\ &\quad + \left\| \nabla_{\mathbf{x}^{(2)}} (F \circ S_2) \Big|_{(S_1(\xi; \mathbf{x}_t^{(1)}), \mathbf{x}_t^{(2)})} - \nabla_{\mathbf{x}^{(2)}} (F \circ S_2) \Big|_{(\mathbb{C}(S_1(\xi; \mathbf{x}_t^{(1)}), \mathbf{x}_t^{(2)}))} \right\|^2 \\ &\leq C_{S_1}^2 L_2^2 \left\| (\mathbb{C}(S_1(\xi; \mathbf{x}_t^{(1)}), \mathbf{x}_t^{(2)}) - (S_1(\xi; \mathbf{x}_t^{(1)}), \mathbf{x}_t^{(2)})) \right\|^2 \\ &\quad + L_2^2 \left\| (\mathbb{C}(S_1(\xi; \mathbf{x}_t^{(1)}), \mathbf{x}_t^{(2)}) - (S_1(\xi; \mathbf{x}_t^{(1)}), \mathbf{x}_t^{(2)})) \right\|^2 \\ &= (1 + C_{S_1}^2) L_2^2 \left\| (\mathbb{C}(S_1(\xi; \mathbf{x}_t^{(1)}), \mathbf{x}_t^{(2)}) - (S_1(\xi; \mathbf{x}_t^{(1)}), \mathbf{x}_t^{(2)})) \right\|^2 \\ &\leq (1 + C_{S_1}^2) L_2^2 \kappa^2 \left\| S_1(\xi; \mathbf{x}_t^{(1)}) \right\|^2 \\ &\leq (1 + C_{S_1}^2) L_2^2 \kappa^2 L_{S_1}^2, \end{aligned} \quad (20)$$

where we get (20) by assumption 4.  $\blacksquare$

## D.2 PROOF OF THEOREM 1

Before showing the proof of the theorem, we present notations, definitions and update rules used in the following proof.

$$\mathbf{x}_{t+1}^{(1)} = \mathbf{x}_t^{(1)} - \eta \cdot \hat{G}_\mu(\mathbf{x}_t, \xi_t, \mathbf{u}_t) = \mathbf{x}_t^{(1)} - \eta \cdot \hat{g}_\mu(\mathbf{x}_t, \xi_t, \mathbf{u}_t) \cdot \mathbf{u}_t^{(1)} \quad (21)$$

$$\mathbf{x}_{t+1}^{(2)} = \mathbf{x}_t^{(2)} - \eta \cdot \hat{G}_\mu(\mathbf{x}_t, \xi_t, \mathbf{u}_t) = \mathbf{x}_t^{(2)} - \eta \cdot \hat{g}_\mu(\mathbf{x}_t, \xi_t, \mathbf{u}_t) \cdot \mathbf{u}_t^{(2)} \quad (22)$$

Introducing  $\mathbf{x}_t = \text{col}[\mathbf{x}_t^{(1)}, \mathbf{x}_t^{(2)}]$  and  $\mathbf{u}_t = \text{col}[\mathbf{u}_t^{(1)}, \mathbf{u}_t^{(2)}]$ , we can write it into one line

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \cdot g_\mu(\mathbf{x}_t, \xi_t, \mathbf{u}_t) \cdot \mathbf{u}_t \quad (23)$$

Next, we define a new (virtual) zeroth-order gradient scalar without the compression step:

$$\begin{aligned} &g_\mu(\mathbf{x}_t, \xi_t, \mathbf{u}_t) \\ &= \frac{F(S_2(a_t^+; \mathbf{x}_t^{(2)} + \mu \mathbf{u}_t^{(2)})) - F(S_2(a_t^-; \mathbf{x}_t^{(2)} - \mu \mathbf{u}_t^{(2)}))}{\mu} \\ &= \frac{F(S_2(S_1(\xi_t; \mathbf{x}_t^{(1)} + \mu \mathbf{u}_t^{(1)}); \mathbf{x}_t^{(2)} + \mu \mathbf{u}_t^{(2)})) - F(S_2(S_1(\xi_t; \mathbf{x}_t^{(1)} - \mu \mathbf{u}_t^{(1)}); \mathbf{x}_t^{(2)} - \mu \mathbf{u}_t^{(2)}))}{\mu} \end{aligned} \quad (24)$$

It is not hard to see that  $\mathbb{E}_{\xi_t, \mathbf{u}_t} [G_\mu(\mathbf{x}_t, \xi_t, \mathbf{u}_t)] = \nabla f_\mu(\mathbf{x}_t)$  using the Lemma 3 (a). Moreover, we define the difference between these two gradient scalars as:

$$\Delta_\mu(\mathbf{x}_t, \xi_t, \mathbf{u}_t) := g_\mu(\mathbf{x}_t, \xi_t, \mathbf{u}_t) - \hat{g}_\mu(\mathbf{x}_t, \xi_t, \mathbf{u}_t) \quad (25)$$

Now, we arrive at

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta g_\mu(\mathbf{x}_t, \xi_t, \mathbf{u}_t) \cdot \mathbf{u}_t + \eta \Delta_\mu(\mathbf{x}_t, \xi_t, \mathbf{u}_t) \cdot \mathbf{u}_t \quad (26)$$

$$= \mathbf{x}_t - \eta G_\mu(\mathbf{x}_t, \xi_t, \mathbf{u}_t) + \eta \Delta_\mu(\mathbf{x}_t, \xi_t, \mathbf{u}_t) \cdot \mathbf{u}_t \quad (27)$$

When there is no ambiguities, we simply use  $G_\mu$  and  $\Delta_\mu$  to replace  $G_\mu(\mathbf{x}_t, \xi_t, \mathbf{u}_t)$  and  $\Delta_\mu(\mathbf{x}_t, \xi_t, \mathbf{u}_t)$  respectively.

*Proof of Theorem 1:*

We start with the Lipschitz condition:

$$\begin{aligned} & \mathbb{E} [f_\mu(\mathbf{x}_{t+1})] \\ & \leq f_\mu(\mathbf{x}_t) + \mathbb{E} [\langle \nabla f_\mu(\mathbf{x}_t), \mathbf{x}_{t+1} - \mathbf{x}_t \rangle] + \frac{L}{2} \mathbb{E} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \\ & = f_\mu(\mathbf{x}_t) - \eta \langle \nabla f_\mu(\mathbf{x}_t), \mathbb{E} [G_\mu - \Delta_\mu \cdot \mathbf{u}_t] \rangle + \frac{L}{2} \eta^2 \mathbb{E} \|\hat{G}_\mu\|^2 \\ & \leq f_\mu(\mathbf{x}_t) - \eta \mathbb{E} [\langle \nabla f_\mu(\mathbf{x}_t), G_\mu \rangle] + \frac{\eta}{2\epsilon} \mathbb{E} \|\nabla f_\mu(\mathbf{x}_t)\|^2 + \frac{\eta\epsilon}{2} \mathbb{E} \|\Delta_\mu \cdot \mathbf{u}_t\|^2 + L\eta^2 \mathbb{E} \|\hat{G}_\mu\|^2 \quad (28) \\ & = f_\mu(\mathbf{x}_t) - \eta \mathbb{E} [\langle \nabla f_\mu(\mathbf{x}_t), \nabla f_\mu(\mathbf{x}_t) - \nabla f_\mu(\mathbf{x}_t) + G_\mu \rangle] + \frac{\eta}{2\epsilon} \mathbb{E} \|\nabla f_\mu(\mathbf{x}_t)\|^2 \\ & \quad + \frac{\eta\epsilon}{2} \mathbb{E} \|G_\mu - \hat{G}_\mu\|^2 + L\eta^2 \mathbb{E} \|\hat{G}_\mu\|^2 \\ & = f_\mu(\mathbf{x}_t) - \eta \mathbb{E} \|\nabla f_\mu(\mathbf{x}_t)\|^2 - \eta \mathbb{E} [\langle \nabla f_\mu(\mathbf{x}_t), G_\mu - \nabla f_\mu(\mathbf{x}_t) \rangle] + \frac{\eta}{2\epsilon} \mathbb{E} \|\nabla f_\mu(\mathbf{x}_t)\|^2 \\ & \quad + \frac{\eta\epsilon}{2} \mathbb{E} \|\nabla f_\mu(\mathbf{x}_t) - \nabla \hat{f}_\mu(\mathbf{x}_t)\|^2 + L\eta^2 \mathbb{E} \|\hat{G}_\mu\|^2 \\ & = f_\mu(\mathbf{x}_t) + \left(\frac{\eta}{2\epsilon} - \eta\right) \mathbb{E} \|\nabla f_\mu(\mathbf{x}_t)\|^2 + \frac{\eta\epsilon}{2} \mathbb{E} \|\nabla f_\mu(\mathbf{x}_t) - \nabla \hat{f}_\mu(\mathbf{x}_t)\|^2 + L\eta^2 \mathbb{E} \|\hat{G}_\mu\|^2, \quad (29) \end{aligned}$$

where (28) applies the Young's inequality on the cross term, where  $\epsilon$  is an arbitrary positive number, and the Jensen's inequality on the last term.

Re-arranging (29), we get

$$\eta \left(1 - \frac{1}{2\epsilon}\right) \mathbb{E} \|\nabla f_\mu(\mathbf{x}_t)\|^2 \leq f_\mu(\mathbf{x}_t) - \mathbb{E} [f_\mu(\mathbf{x}_{t+1})] + \frac{\eta\epsilon}{2} \mathbb{E} \|\nabla f_\mu(\mathbf{x}_t) - \nabla \hat{f}_\mu(\mathbf{x}_t)\|^2 + L\eta^2 \mathbb{E} \|\hat{G}_\mu\|^2. \quad (30)$$

Then, we sum up all (30) for all  $t = 0, \dots, T-1$  and establish

$$\begin{aligned} & \eta \left(1 - \frac{1}{2\epsilon}\right) \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f_\mu(\mathbf{x}_t)\|^2 \\ & \leq f_\mu(\mathbf{x}_0) - f_\mu(\mathbf{x}_T) + \frac{\eta\epsilon}{2} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f_\mu(\mathbf{x}_t) - \nabla \hat{f}_\mu(\mathbf{x}_t)\|^2 + L\eta^2 \sum_{t=0}^{T-1} \mathbb{E} \|\hat{G}_\mu\|^2 \quad (31) \end{aligned}$$

$$\leq f_\mu(\mathbf{x}_0) - f_\mu(\mathbf{x}^*) + \underbrace{\frac{\eta\epsilon}{2} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f_\mu(\mathbf{x}_t) - \nabla \hat{f}_\mu(\mathbf{x}_t)\|^2}_{A_1} + \underbrace{L\eta^2 \sum_{t=0}^{T-1} \mathbb{E} \|\hat{G}_\mu\|^2}_{A_2}, \quad (32)$$

where the last inequality follows from  $f_\mu(\mathbf{x}^*) \leq f_\mu(\mathbf{x}_T)$ .

Next, we process  $A_1$  as follows. We apply (15) on the first term and the third term of (33). For the second term of (33), we obtain it by Lemma 4.

Now, we deal with  $A_1$ :

$$A_1 = \mathbb{E} \|\nabla f_\mu(\mathbf{x}_t) - \nabla \hat{f}_\mu(\mathbf{x}_t)\|^2$$

$$\begin{aligned}
&= \left\| \nabla f_\mu(\mathbf{x}_t) - \nabla f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t) - \nabla \hat{f}(\mathbf{x}_t) + \nabla \hat{f}(\mathbf{x}_t) - \nabla \hat{f}_\mu(\mathbf{x}_t) \right\|^2 \\
&\leq 3 \left\| \nabla f_\mu(\mathbf{x}_t) - \nabla f(\mathbf{x}_t) \right\|^2 + 3 \left\| \nabla f(\mathbf{x}_t) - \nabla \hat{f}(\mathbf{x}_t) \right\|^2 + 3 \left\| \nabla \hat{f}(\mathbf{x}_t) - \nabla \hat{f}_\mu(\mathbf{x}_t) \right\|^2 \quad (33)
\end{aligned}$$

$$\leq \frac{3}{4} \mu^2 L^2 (d+3)^3 + 3(1 + C_{S_1}^2) L_2^2 \kappa^2 L_{S_1}^2 + \frac{3}{4} \mu^2 L^2 (d+3)^3 \quad (34)$$

$$= \frac{3}{2} \mu^2 L^2 (d+3)^3 + 3(1 + C_{S_1}^2) L_2^2 \kappa^2 L_{S_1}^2, \quad (35)$$

where we obtain (34) by Lemma 4.

Then, we start to deal with  $A_2$ .

$$\begin{aligned}
A_2 &= L\eta^2 \sum_{t=0}^{T-1} \mathbb{E} \left\| \hat{G}_\mu \right\|^2 \\
&\leq L\eta^2 \sum_{t=1}^T \left( \frac{\mu^2}{2} L^2 (d+6)^3 + 2(d+4) \mathbb{E} \left\| \hat{G} \right\|^2 \right) \\
&\leq L\eta^2 \sum_{t=0}^{T-1} \left( \frac{\mu^2}{2} L^2 (d+6)^3 + 2(d+4) \left( \mathbb{E} \left\| \nabla \hat{f}(\mathbf{x}_t) \right\|^2 + \sigma^2 \right) \right) \\
&= L\eta^2 \sum_{t=0}^{T-1} \left( \frac{\mu^2}{2} L^2 (d+6)^3 + 2(d+4) \left( \mathbb{E} \left\| \nabla \hat{f}(\mathbf{x}_t) - \nabla f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t) \right\|^2 + \sigma^2 \right) \right) \\
&\leq L\eta^2 \sum_{t=0}^{T-1} \left( \frac{\mu^2}{2} L^2 (d+6)^3 + 2(d+4) \left( 2\mathbb{E} \left\| \nabla \hat{f}(\mathbf{x}_t) - \nabla f(\mathbf{x}_t) \right\|^2 + 2\mathbb{E} \left\| \nabla f(\mathbf{x}_t) \right\|^2 + \sigma^2 \right) \right) \quad (36) \\
&\leq L\eta^2 \sum_{t=0}^{T-1} \left( \frac{\mu^2}{2} L^2 (d+6)^3 + 2(d+4) \left( 2(1 + C_{S_1}^2) L_2^2 \kappa^2 L_{S_1}^2 + \sigma^2 \right) + 4(d+4) \mathbb{E} \left\| \nabla f(\mathbf{x}_t) \right\|^2 \right) \\
&= L\eta^2 \sum_{t=0}^{T-1} \left( \frac{\mu^2}{2} L^2 (d+6)^3 + 2(d+4) \left( 2(1 + C_{S_1}^2) L_2^2 \kappa^2 L_{S_1}^2 + \sigma^2 \right) \right) + 4(d+4) L\eta^2 \sum_{t=0}^{T-1} \mathbb{E} \left\| \nabla f(\mathbf{x}_t) \right\|^2, \quad (37)
\end{aligned}$$

where we apply Lemma 4 on the term  $\left\| \nabla \hat{f}(\mathbf{x}_t) - \nabla f(\mathbf{x}_t) \right\|^2$  of (36).

Putting all pieces above together, we establish:

$$\begin{aligned}
\left( \eta - \frac{\eta}{2\epsilon} \right) \sum_{t=0}^{T-1} \mathbb{E} \left\| \nabla f_\mu(\mathbf{x}_t) \right\|^2 &\leq f_\mu(\mathbf{x}_0) - f_\mu(\mathbf{x}^*) + \frac{\eta\epsilon}{2} \sum_{t=0}^{T-1} \left( \frac{3}{2} \mu^2 L^2 (d+3)^3 + 3(1 + C_{S_1}^2) L_2^2 \kappa^2 L_{S_1}^2 \right) \\
&\quad + L\eta^2 \sum_{t=0}^{T-1} \left( \frac{\mu^2}{2} L^2 (d+6)^3 + 2(d+4) \left( 2(1 + C_{S_1}^2) L_2^2 \kappa^2 L_{S_1}^2 + \sigma^2 \right) \right) \\
&\quad + 4(d+4) L\eta^2 \sum_{t=0}^{T-1} \mathbb{E} \left\| \nabla f(\mathbf{x}_t) \right\|^2 \quad (38)
\end{aligned}$$

By using (19) and re-arranging (38), we get

$$\begin{aligned}
&\left( \frac{\eta}{2} \left( 1 - \frac{1}{2\epsilon} \right) - 4(d+4) L\eta^2 \right) \sum_{t=0}^{T-1} \mathbb{E} \left\| \nabla f(\mathbf{x}_t) \right\|^2 \\
&\leq f(\mathbf{x}_0) - f(\mathbf{x}^*) + \mu^2 Ld + \frac{\eta\epsilon}{2} \sum_{t=0}^{T-1} \left( \frac{3}{2} \mu^2 L^2 (d+3)^3 + 3(1 + C_{S_1}^2) L_2^2 \kappa^2 L_{S_1}^2 \right)
\end{aligned}$$

1079

$$\begin{aligned}
& + L\eta^2 \sum_{t=0}^{T-1} \left( \frac{\mu^2}{2} L^2 (d+6)^3 + 2(d+4) \left( 2(1+C_{S_1}^2) L_2^2 \kappa^2 L_{S_1}^2 + \sigma^2 \right) \right) \\
& + \eta \left( 1 - \frac{1}{2\epsilon} \right) \sum_{t=0}^{T-1} \frac{\mu^2}{4} L^2 (d+3)^3
\end{aligned} \tag{39}$$

We select  $\epsilon = 1$  and further simplify

$$\begin{aligned}
& \left( \eta - 16(d+4)L\eta^2 \right) \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(\mathbf{x}_t)\|^2 \leq 4(f(\mathbf{x}_0) - f(\mathbf{x}^*)) + 4\mu^2 Ld + \eta(3\mu^2 L^2 (d+3)^3 \\
& \quad + L\eta^2 T \left( 2\mu^2 L^2 (d+6)^3 + 8(d+4) \left( 2(1+C_{S_1}^2) L_2^2 \kappa^2 L_{S_1}^2 + \sigma^2 \right) \right) \\
& \quad + 6(1+C_{S_1}^2) L_2^2 \kappa^2 L_{S_1}^2 T + \frac{1}{2} \eta \mu^2 L^2 (d+3)^3 T
\end{aligned} \tag{40}$$

We assume that  $\frac{1}{2}\eta \leq (\eta - 16(d+4)L\eta^2)$  and then get  $\eta \leq \frac{1}{32(d+4)L}$ . Then, we can further simplify the inequality above.

$$\begin{aligned}
\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(\mathbf{x}_t)\|^2 & \leq \frac{8(f(\mathbf{x}_0) - f(\mathbf{x}^*))}{T\eta} + \frac{8\mu^2 Ln}{T\eta} + 7\mu^2 L^2 (d+3)^3 + 12(1+C_{S_1}^2) L_2^2 \kappa^2 L_{S_1}^2 \\
& \quad + 4L\eta \left( \mu^2 L^2 (d+6)^3 + 4(d+4) \left( 2(1+C_{S_1}^2) L_2^2 \kappa^2 L_{S_1}^2 + \sigma^2 \right) \right)
\end{aligned} \tag{41}$$

Assuming that  $\eta = \mathcal{O}\left(\frac{1}{\sqrt{Td}}\right)$ , then we obtain

$$\begin{aligned}
\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(\mathbf{x}_t)\|^2 & = \mathcal{O}\left(\frac{D\sqrt{d}}{\sqrt{T}}\right) + \mathcal{O}\left(\frac{\mu^2 d^{\frac{3}{2}}}{\sqrt{T}}\right) + \mathcal{O}\left(\mu^2 (d+3)^3\right) + \mathcal{O}\left((1+C_{S_1}^2) L_2^2 \kappa^2 L_{S_1}^2\right) \\
& \quad + \mathcal{O}\left(\frac{\mu^2 (d+6)^3}{\sqrt{Td}}\right) + \mathcal{O}\left(\frac{(d+4)(1+C_{S_1}^2) L_2^2 \kappa^2 L_{S_1}^2}{\sqrt{Td}}\right) + \mathcal{O}\left(\frac{(d+4)\sigma^2}{\sqrt{Td}}\right),
\end{aligned} \tag{42}$$

where  $D = f(\mathbf{x}_0) - f(\mathbf{x}^*)$ .

Further, we set  $\mu \leq \frac{1}{(d+6)\sqrt{T}}$ , then we can get

$$\begin{aligned}
\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(\mathbf{x}_t)\|^2 & = \mathcal{O}\left(\frac{D\sqrt{d}}{\sqrt{T}}\right) + \mathcal{O}\left(\frac{1}{T\sqrt{d}}\right) + \mathcal{O}\left(\frac{d}{T}\right) + \mathcal{O}\left((1+C_{S_1}^2) L_2^2 \kappa^2 L_{S_1}^2\right) \\
& \quad + \mathcal{O}\left(\frac{\sqrt{d}}{T^{\frac{3}{2}}}\right) + \mathcal{O}\left(\frac{\sqrt{d}}{\sqrt{T}}(1+C_{S_1}^2) L_2^2 \kappa^2 L_{S_1}^2\right) + \mathcal{O}\left(\frac{\sqrt{d}}{\sqrt{T}}\sigma^2\right)
\end{aligned}$$

Further, we simplify it, then we can get

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(\mathbf{x}_t)\|^2 = \mathcal{O}\left(\frac{D\sqrt{d}}{\sqrt{T}}\right) + \mathcal{O}\left(\frac{d}{T}\right) + \mathcal{O}\left((1+C_{S_1}^2) L_2^2 \kappa^2 L_{S_1}^2 \left(\frac{\sqrt{d}}{\sqrt{T}} + 1\right)\right) + \mathcal{O}\left(\frac{\sqrt{d}}{\sqrt{T}}\sigma^2\right)$$

■

### D.3 PROOF OF THEOREM 2

**Assumption 6 (Lipschitz Gradients)** For multi-node cases ( $M > 2$ ), we suppose that

- $f(\cdot)$  has  $L$ -Lipschitz gradient,
- $f \circ S_M \circ \dots \circ S_{i+1}$  has a  $L_{f \circ S_M \circ \dots \circ S_{i+1}}$ -Lipschitz gradient, and its gradient is bounded by  $C_{f \circ S_M \circ \dots \circ S_{i+1}}$  for all  $i = 1, \dots, M-1$ ,
- The submodel  $S_i$  is  $L_{S_i}$ -Lipschitz, and its gradient is bounded by  $C_{S_i}$ , for all  $i = 1, \dots, M$ .

**Lemma 5 (Distance between Gradients of Uncompressed and Compressed Activations,  $M > 2$ )**  
 For Alg. 3, the difference between gradients of uncompressed and compressed activations can be bounded by a constant term as follows:

$$\left\| \nabla f(\mathbf{x}_t) - \nabla \hat{f}(\mathbf{x}_t) \right\|^2 \leq 2M\kappa^2 L_S^2 \Psi$$

*Proof of lemma 5:*

For simplicity in the following proof, we denote that

$$\bar{S}_i = S_i(\cdots(S_2(S_1(\xi, \mathbf{x}_t^{(1)}); \mathbf{x}_t^{(2)}); \cdots); \mathbf{x}_t^{(i)}) \quad \text{and} \quad \bar{m}_i = S_i(\cdots(\mathcal{C}(S_1(\xi; \mathbf{x}_t^{(1)})); \cdots); \mathbf{x}_t^{(i)}).$$

Then, we bound the gradient difference as follows:

$$\begin{aligned} & \left\| \nabla f(\mathbf{x}_t) - \nabla \hat{f}(\mathbf{x}_t) \right\|^2 \\ &= \left\| \nabla_{\mathbf{x}^{(1)}} f(\mathbf{x}_t) - \nabla_{\mathbf{x}^{(1)}} \hat{f}(\mathbf{x}_t) \right\|^2 + \cdots + \left\| \nabla_{\mathbf{x}^{(M)}} f(\mathbf{x}_t) - \nabla_{\mathbf{x}^{(M)}} \hat{f}(\mathbf{x}_t) \right\|^2 \\ &= \sum_{i=1}^M \left\| \nabla_{S_i}(F \circ S_M \circ \cdots \circ S_{i+1}) \Big|_{(\bar{S}_i, \dots, \mathbf{x}_t^{(M)})} \cdot \nabla_{\mathbf{x}^{(2)}} S_i \Big|_{(\bar{S}_{i-1}, \mathbf{x}_t^{(i)})} \right. \\ & \quad \left. - \nabla_{S_i}(F \circ S_M \circ \cdots \circ S_{i+1}) \Big|_{(\bar{m}_i, \mathbf{x}_t^{(i+1)}, \dots, \mathbf{x}_t^{(M)})} \cdot \nabla_{\mathbf{x}^{(i)}} S_i \Big|_{(\bar{m}_{i-1}, \mathbf{x}_t^{(i)})} \right\|^2 \\ &\leq (1 + 2C_{S_{M-1}}^2) L_{f \circ S_M}^2 \kappa^2 L_{S_{M-1}}^2 + 2\kappa^2 \sum_{i=1}^{M-2} (C_{S_{M-2}}^2 L_{f \circ S_M \circ \cdots \circ S_{i+1}}^2 + C_{f \circ S_M \circ \cdots \circ S_{i+2}}^2 L_{S_{i+1}}^2) L_{S_i}^2 \\ &\leq 2M\kappa^2 L_S^2 \max \left\{ (1 + 2C_{S_{M-1}}^2) L_{f \circ S_M}^2 \kappa^2, \max_{i \in [M-2]} \{ C_{S_{M-2}}^2 L_{f \circ S_M \circ \cdots \circ S_{i+1}}^2 + C_{f \circ S_M \circ \cdots \circ S_{i+2}}^2 L_{S_{i+1}}^2 \} \right\}, \end{aligned}$$

where  $L_S^2 = \max\{L_{S_1}^2, \dots, L_{S_M}^2\}$ . ■

**Theorem 2 (Convergence of SparQ under Non-Convexity,  $M > 2$ )** For Alg. 3, under assumptions 1, 2, 3, 5 and 6, if the number of computing nodes  $M > 2$  and learning rate  $\eta \leq 1/32(d+4)L$ , then the sequence  $\{\mathbf{x}_t\}$  generated by SparQ satisfies

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(\mathbf{x}_t)\|^2 &\leq \frac{8(f(\mathbf{x}_0) - f(\mathbf{x}^*))}{T\eta} + \frac{8\mu^2 Ln}{T\eta} + 7\mu^2 L^2 (d+3)^3 + 12M\kappa^2 L_S^2 \Psi \\ &\quad + 4L\eta \left( \mu^2 L^2 (d+6)^3 + 4(d+4)(4M\kappa^2 L_S^2 \Psi + \sigma^2) \right), \end{aligned}$$

where the definitions of  $L_S$  and  $\Psi$  can be found in the following proof.

*Proof of Theorem 2:*

Here, we only provide a rough proof because we use the same proof framework as the proof of Theorem 1. The key differences are using Lemma 5 in inequalities (33) and (36).

Hence,  $A_1$  and  $A_2$  will be modified as follows:

$$\begin{aligned} A_1 &= \left\| \nabla f_\mu(\mathbf{x}_t) - \nabla \hat{f}_\mu(\mathbf{x}_t) \right\|^2 \\ &= \left\| \nabla f_\mu(\mathbf{x}_t) - \nabla f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t) - \nabla \hat{f}(\mathbf{x}_t) + \nabla \hat{f}(\mathbf{x}_t) - \nabla \hat{f}_\mu(\mathbf{x}_t) \right\|^2 \\ &\leq 3 \left\| \nabla f_\mu(\mathbf{x}_t) - \nabla f(\mathbf{x}_t) \right\|^2 + 3 \left\| \nabla f(\mathbf{x}_t) - \nabla \hat{f}(\mathbf{x}_t) \right\|^2 + 3 \left\| \nabla \hat{f}(\mathbf{x}_t) - \nabla \hat{f}_\mu(\mathbf{x}_t) \right\|^2 \\ &\leq \frac{3}{4} \mu^2 L^2 (d+3)^3 + 6M\kappa^2 L_S^2 \Psi + \frac{3}{4} \mu^2 L^2 (d+3)^3 \tag{43} \\ &= \frac{3}{2} \mu^2 L^2 (d+3)^3 + 6M\kappa^2 L_S^2 \Psi, \end{aligned}$$

1188 where we obtain (43) by Lemma 5.

$$\begin{aligned}
1189 \quad A_2 &= L\eta^2 \sum_{t=0}^{T-1} \mathbb{E} \left\| \hat{G}_\mu \right\|^2 \\
1190 \quad &\leq L\eta^2 \sum_{t=1}^T \left( \frac{\mu^2}{2} L^2 (d+6)^3 + 2(d+4) \mathbb{E} \left\| \hat{G} \right\|^2 \right) \\
1191 \quad &\leq L\eta^2 \sum_{t=0}^{T-1} \left( \frac{\mu^2}{2} L^2 (d+6)^3 + 2(d+4) \left( \mathbb{E} \left\| \nabla \hat{f}(\mathbf{x}_t) \right\|^2 + \sigma^2 \right) \right) \\
1192 \quad &= L\eta^2 \sum_{t=0}^{T-1} \left( \frac{\mu^2}{2} L^2 (d+6)^3 + 2(d+4) \left( \mathbb{E} \left\| \nabla \hat{f}(\mathbf{x}_t) - \nabla f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t) \right\|^2 + \sigma^2 \right) \right) \\
1193 \quad &\leq L\eta^2 \sum_{t=0}^{T-1} \left( \frac{\mu^2}{2} L^2 (d+6)^3 + 2(d+4) \left( 2\mathbb{E} \left\| \nabla \hat{f}(\mathbf{x}_t) - \nabla f(\mathbf{x}_t) \right\|^2 + 2\mathbb{E} \left\| \nabla f(\mathbf{x}_t) \right\|^2 + \sigma^2 \right) \right) \\
1194 \quad &\leq L\eta^2 \sum_{t=0}^{T-1} \left( \frac{\mu^2}{2} L^2 (d+6)^3 + 2(d+4) \left( 4M\kappa^2 L_S^2 \Psi + \sigma^2 \right) + 4(d+4) \mathbb{E} \left\| \nabla f(\mathbf{x}_t) \right\|^2 \right) \\
1195 \quad &= L\eta^2 \sum_{t=0}^{T-1} \left( \frac{\mu^2}{2} L^2 (d+6)^3 + 2(d+4) (4M\kappa^2 L_S^2 \Psi + \sigma^2) \right) + 4(d+4) L\eta^2 \sum_{t=0}^{T-1} \mathbb{E} \left\| \nabla f(\mathbf{x}_t) \right\|^2, \\
1196 \quad & \\
1197 \quad & \\
1198 \quad & \\
1199 \quad & \\
1200 \quad & \\
1201 \quad & \\
1202 \quad & \\
1203 \quad & \\
1204 \quad & \\
1205 \quad & \\
1206 \quad & \\
1207 \quad & \\
1208 \quad & \\
1209 \quad &
\end{aligned}$$

1210 Then, the subsequent steps are the same as the proof of Theorem 1, so we skip them and directly  
1211 arrive at

$$\begin{aligned}
1212 \quad &\left( \eta - 16(d+4)L\eta^2 \right) \sum_{t=0}^{T-1} \mathbb{E} \left\| \nabla f(\mathbf{x}_t) \right\|^2 \leq 4(f(\mathbf{x}_0) - f(\mathbf{x}^*)) + 4\mu^2 Ld + \eta(3\mu^2 L^2 (d+3)^3 \\
1213 \quad &+ L\eta^2 T (2\mu^2 L^2 (d+6)^3 + 8(d+4)(4M\kappa^2 L_S^2 \Psi + \sigma^2)) + 6M\kappa^2 \Psi T + \frac{1}{2} \eta \mu^2 L^2 (d+3)^3 T \\
1214 \quad &
\end{aligned}$$

1215 We assume that  $\frac{1}{2} \eta \leq (\eta - 16(d+4)L\eta^2)$  and then get  $\eta \leq \frac{1}{32(d+4)L}$ . Then, we can further simplify  
1216 the inequality above and obtain

$$\begin{aligned}
1217 \quad &\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \left\| \nabla f(\mathbf{x}_t) \right\|^2 \leq \frac{8(f(\mathbf{x}_0) - f(\mathbf{x}^*))}{T\eta} + \frac{8\mu^2 Ln}{T\eta} + 7\mu^2 L^2 (d+3)^3 + 12M\kappa^2 L_S^2 \Psi \\
1218 \quad &+ 4L\eta \left( \mu^2 L^2 (d+6)^3 + 4(d+4)(4M\kappa^2 L_S^2 \Psi + \sigma^2) \right) \\
1219 \quad & \\
1220 \quad & \\
1221 \quad & \\
1222 \quad & \\
1223 \quad & \\
1224 \quad & \blacksquare \\
1225 \quad & \\
1226 \quad & \\
1227 \quad & \\
1228 \quad & \\
1229 \quad & \\
1230 \quad & \\
1231 \quad & \\
1232 \quad & \\
1233 \quad & \\
1234 \quad & \\
1235 \quad & \\
1236 \quad & \\
1237 \quad & \\
1238 \quad & \\
1239 \quad & \\
1240 \quad & \\
1241 \quad &
\end{aligned}$$