

SUBSPACE OPTIMIZATION FOR LARGE LANGUAGE MODELS WITH CONVERGENCE GUARANTEES

Anonymous authors

Paper under double-blind review

ABSTRACT

Subspace optimization algorithms, with GaLore (Zhao et al., 2024) as a representative method, have gained popularity for pre-training or fine-tuning large language models (LLMs) due to their memory efficiency. However, their convergence guarantees remain unclear, particularly in stochastic settings. In this paper, we unexpectedly discover that GaLore does not always converge to the optimal solution and substantiate this finding with an explicit counter-example. We then investigate the conditions under which GaLore can achieve convergence, demonstrating that it does so either in deterministic scenarios or when using a sufficiently large mini-batch size. More significantly, we introduce **GoLore** (**G**radient **r**andom **L**ow-**r**ank **p**rojection), a novel variant of GaLore that provably converges in stochastic settings, even with standard batch sizes. Our convergence analysis can be readily extended to other sparse subspace optimization algorithms. Finally, we conduct numerical experiments to validate our theoretical results and empirically explore the proposed mechanisms.

1 INTRODUCTION

Large Language Models (LLMs) have demonstrated impressive performance across a variety of tasks, including language processing, planning, and coding. However, LLMs require substantial computational resources and memory due to their large model size and the extensive amounts of training data. Consequently, recent advancements in stochastic optimization have focused on developing memory-efficient strategies to pre-train or fine-tune LLMs with significantly reduced computing resources. Most approaches (Vyas et al., 2024; Ramesh et al., 2024; Luo et al., 2023; Liu et al., 2024; Bini et al., 2024; Hao et al., 2024; Zhao et al., 2024; Muhamed et al., 2024; Pan et al., 2024; Loeschcke et al., 2024; Hayou et al., 2024; Lialin et al., 2023; Han et al., 2024; Song et al., 2023) concentrate on reducing the memory of optimizer states, which are critical components of overall training memory consumption. For instance, optimizers such as Adam (Kingma, 2014) and AdamW (Loshchilov, 2017) maintain first and second-order momentum terms for gradients as optimizer states, leading to significant memory overhead for large models.

Among the most popular memory-efficient fine-tuning algorithms is LoRA (Hu et al., 2021), which decreases the number of trainable parameters by employing low-rank model adapters. However, the low-rank constraint on weight updates can result in substantial performance degradation for tasks that require full-rank updates, particularly in the pre-training of LLMs. To address this issue, several LoRA variants have been proposed, including ReLoRA (Lialin et al., 2023) and SLTrain (Han et al., 2024). Recently, GaLore (Zhao et al., 2024) has emerged as an effective solution, significantly reducing optimizer states by projecting full-parameter gradients into periodically recomputed subspaces. By retaining optimizer states in low-rank subspaces, GaLore can reduce memory usage by over 60%, enabling the pre-training of a 7B model on an NVIDIA RTX 4090 with 24GB of memory. In contrast, the vanilla 8-bit Adam without low-rank projection requires over 40GB of memory.

1.1 FUNDAMENTAL OPEN QUESTIONS AND MAIN RESULTS

While GaLore’s memory efficiency has been well established both theoretically and empirically, its convergence guarantees remain unclear. This raises the following fundamental open question:

Q1. Can GaLore converge to stationary solutions, under regular assumptions?

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

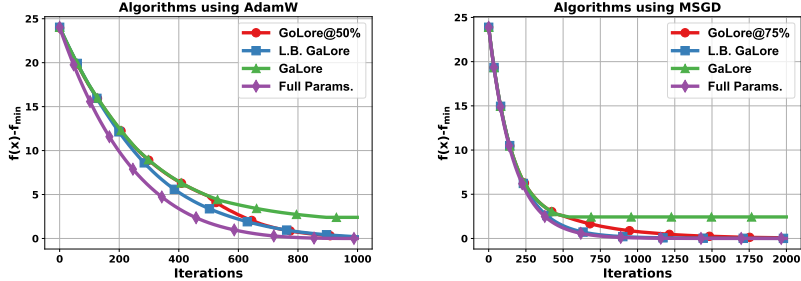


Figure 1: Loss curves of algorithms using AdamW (left) and Momentum SGD (right) on problem (1), where *L.B. GaLore* stands for large-batch GaLore, *GoLore@x%* applies GaLore for the beginning $(100 - x)\%$ iterations and GoLore for the last $x\%$ iterations.

By *stationary solutions*, we refer to first-order stationary points $x \in \mathbb{R}^d$ such that $\nabla f(x) = 0$ for objective function $f : \mathbb{R}^d \rightarrow \mathbb{R}$. By *regular assumptions*, we refer to standard conditions in non-convex smooth optimization, including lower boundedness, L -smoothness and unbiased stochastic gradients with bounded variances, as outlined in Assumptions 1-3 in Sec. 2.

Contrary to expectations, our investigation reveals that GaLore does **NOT** converge to stationary solutions under regular assumptions. The intuition behind this finding is straightforward: GaLore projects the stochastic gradient matrix onto a low-rank subspace spanned by the top r singular vectors obtained via Singular Value Decomposition (SVD), effectively capturing the dominant components of the stochastic gradient matrix. However, the stochastic gradient comprises two components: the true gradient and gradient noise. When the true gradient dominates, the SVD-identified subspace primarily captures the gradient component. *In contrast, as the algorithm approaches a local minimum so that the true gradient diminishes while noise persists, the SVD-derived subspace captures only the noise component, rather than the true gradient, ultimately leading to non-convergence.* To validate this intuition, we construct a counter-example demonstrating that GaLore fails to converge to stationary solutions, see the illustration in Fig. 1. This leads us to a subsequent open question:

Q2. Under what additional assumptions can GaLore converge to stationary solutions?

Based on the preceding discussion, we conclude that the SVD-identified subspace in GaLore aligns well with the descent direction in scenarios where the true gradient component dominates the gradient noise component. This observation naturally leads to two additional assumptions under which GaLore can converge:

- **Noise-Free Assumption.** We theoretically establish that GaLore converges at a rate of $\mathcal{O}(1/T)$ in the deterministic and non-convex setting.
- **Large-Batch Assumption.** We theoretically demonstrate that GaLore converges at a rate of $\mathcal{O}(1/\sqrt{T})$ in the stochastic and non-convex setting, provided that the batch size is extremely large and increases with the number of iterations T , e.g., a batch size of $\Theta(\sqrt{T})$.

However, neither the noise-free assumption nor the large-batch assumption applies to the practical **pre-training** and fine-tuning of LLMs. This leads to another fundamental open question:

Q3. Under what modifications can GaLore provably converge in the LLM setting, in which gradient noise presents and the batch-size cannot be extremely large?

It is evident that SVD-based projections cannot extract meaningful information from noise-dominant matrices. To address this issue, this paper proposes modifying the SVD projection to a **Gradient Random Low-Rank** projection, resulting in the **GoLore** algorithm for pre-training or fine-tuning LLMs. This random projection can effectively capture gradient information even when gradient noise predominates, allowing for convergence in the stochastic and non-convex setting with normal batch sizes. We establish that GoLore converges at a rate of $\mathcal{O}(1/\sqrt{T})$ under standard assumptions.

In our empirical experiments, we implement GaLore during the primary phases of pre-training or fine-tuning LLMs due to its efficacy in capturing the gradient component using SVD-based projection. In contrast, we employ GoLore in the final phase, leveraging its ability to extract the gradient

108 component from noise-dominant stochastic gradients using random projection. This approach en-
 109 hances performance compared to employing GaLore throughout all stages.

110 While our analysis primarily focuses on GaLore, it also has significant connections to other memory-
 111 efficient algorithms. We demonstrate that a ReLoRA-like implementation is equivalent to GaLore,
 112 which is more computational efficient with little additional memory overhead. Furthermore, our
 113 theoretical results can be easily adapted to sparse [subspace descent](#) algorithms with minimal effort.

114 **Contributions.** Our contributions can be summarized as follows:

- 115 • We find that GaLore cannot converge to stationary solutions under regular assumptions. The key
 116 insight is that the SVD-derived subspace primarily captures the noise component rather than the
 117 true gradient in scenarios where gradient noise predominates. We validate the non-convergence
 118 of GaLore by providing an explicit counterexample. This addresses Question Q1.
- 119 • Inspired by the aforementioned insight, we propose [different](#) additional assumptions under
 120 which GaLore can provably converge to stationary solutions. Under the noise-free assump-
 121 tion, we establish that GaLore converges at a rate of $\mathcal{O}(1/T)$. Under the large-batch assumption
 122 or [some additional isotropic noise assumptions](#), we demonstrate that GaLore converges at a rate
 123 of $\mathcal{O}(1/\sqrt{T})$. This addresses Question Q2.
- 124 • In settings where gradient noise persists and the batch size cannot be extremely large, we modify
 125 the SVD projection in GaLore to a random projection, resulting in the GoLore algorithm that
 126 provably converges to stationary solutions at a rate of $\mathcal{O}(1/\sqrt{T})$. This addresses Question Q3.
- 127 • We present an equivalent yet more computationally efficient, ReLoRA-like implementation of
 128 GaLore/GoLore, and extend our analysis to other sparse [subspace descent](#) algorithms.
- 129 • We conduct experiments across various tasks to validate our theoretical findings. In particular,
 130 by alternately using GaLore and GoLore during different phases in LLMs pre-training and fine-
 131 tuning, we achieve enhanced empirical performance.

132 1.2 RELATED WORK

133 **Memory-efficient training.** In LLM training, the primary memory consumption arises not only
 134 from the model parameters but also from activation values and optimizer states. Jiang et al. (2022)
 135 and Yu et al. (2024) have proposed methods to compress activation values into sparse vectors to
 136 alleviate memory usage. Other approaches primarily focus on reducing optimizer states. A notable
 137 work, LoRA (Hu et al., 2021) reparameterizes the weight matrix $\mathbf{W} \in \mathbb{R}^{m \times n}$ as $\mathbf{W} = \mathbf{W}_0 + \mathbf{B}\mathbf{A}$,
 138 where $\mathbf{W}_0 \in \mathbb{R}^{m \times n}$ remains frozen as the pre-trained weights, and $\mathbf{B} \in \mathbb{R}^{m \times r}$ and $\mathbf{A} \in \mathbb{R}^{r \times n}$
 139 are learnable low-rank adapters. Variants of LoRA, such as those proposed by Liu et al. (2024)
 140 and Hayou et al. (2024), aim to enhance training performance. However, constrained to low-rank
 141 updates, LoRA and its variants are primarily effective for fine-tuning tasks and struggle with pre-
 142 training tasks that require high-rank updates. To address this limitation, ReLoRA (Lialin et al.,
 143 2023) enables high-rank updates by accumulating multiple LoRA updates, while LISA (Pan et al.,
 144 2024) learns full-parameter updates on dynamically selected trainable layers. GaLore (Zhao et al.,
 145 2024) and FLORA (Hao et al., 2024) achieve high-rank updates by accumulating low-rank updates
 146 in periodically recomputed subspaces, and SLTrain (Han et al., 2024) employs additional sparse
 147 adapters for high-rank updates. SIFT (Song et al., 2023) also utilizes sparse updates. Although
 148 these algorithms have demonstrated comparable empirical performance to full-parameter training
 149 methods, theoretical guarantees regarding their convergence have not been established. A recent
 150 study by Liang et al. (2024) provides a proof of continuous-time convergence for a class of online
 151 subspace descent algorithms, however, its analysis depends on the availability of true gradients rather
 152 than the stochastic gradients that are more practical in LLM training. To the best of our knowledge,
 153 this work offers the *first* analysis of the discrete-time convergence rate for memory-efficient LLM
 154 training algorithms in stochastic settings.

155 **Convergence for lossy algorithms.** Many optimization algorithms utilize lossy compression on
 156 training dynamics, such as gradients, particularly in the realm of distributed optimization with com-
 157 munication compression. Researchers have established convergence properties for these algorithms
 158 based on either unbiased (Li et al., 2020; Li & Richtárik, 2021; Condat et al., 2024; He et al.,
 159 2024b;a; Mishchenko et al., 2019; Gorbunov et al., 2021; Alistarh et al., 2017; He et al., 2023) or
 160 contractive (Richtárik et al., 2021; Xie et al., 2020; Fatkhullin et al., 2024; He et al., 2023) com-
 161 pressibility. Kozak et al. (2019) provides a convergence analysis for subspace compression under

Polyak-Lojasiewicz or convex conditions, where the subspace compression adheres contractive compressibility at each iteration. Despite these extensive findings, analyzing the convergence properties of [subspace descent](#) algorithms like GaLore remains challenging, as the compressions used can be neither unbiased nor contractive due to the reuse of projection matrices.

2 PRELIMINARIES AND ASSUMPTIONS

Full-parameter training. Training an N_L -layer neural network can be formulated as the following optimization problem:

$$\min_{\mathbf{x}} f(\mathbf{x}) := \mathbb{E}_{\xi \sim \mathcal{D}} F(\mathbf{x}; \xi).$$

Here, $\mathbf{x} = (\text{vec}(\mathbf{X}_1)^\top, \dots, \text{vec}(\mathbf{X}_{N_L})^\top)^\top$ collects all trainable parameters in the model, where N_L is the number of layers, $\mathbf{X}_\ell \in \mathbb{R}^{m_\ell \times n_\ell}$ denotes the weight matrix in the ℓ -th layer, $\ell = 1, \dots, N_L$. $F(\mathbf{x}; \xi)$ computes the loss with respect to data point ξ , \mathcal{D} denotes the training data distribution. In full-parameter training, we directly apply the optimizer to the full-parameter \mathbf{x} :

$$\mathbf{G}_\ell^{(t)} = \nabla_\ell F(\mathbf{x}^{(t)}; \xi^{(t)}), \quad \mathbf{X}_\ell^{(t+1)} = \mathbf{X}_\ell^{(t)} + \rho_\ell^{(t)}(\mathbf{G}_\ell^{(t)}), \quad \ell = 1, \dots, N_L;$$

where ∇_ℓ computes the gradient with respect to the ℓ -th weight matrix \mathbf{X}_ℓ , superscript (t) denotes the variable in the t -th iteration, and $\rho_\ell^{(t)}$ is an entry-wise stateful gradient operator, such as Adam or Momentum SGD (MSGD). Specifically, using MSGD leads to the following $\rho_\ell^{(t)}(\cdot)$:

$$\mathbf{M}_\ell^{(t)} = (1 - \beta_1)\mathbf{M}_\ell^{(t-1)} + \beta_1\mathbf{G}_\ell^{(t)}; \quad \rho_\ell^{(t)}(\mathbf{G}_\ell^{(t)}) = -\eta\mathbf{M}_\ell^{(t)};$$

where η is the learning rate, $\beta_1 \in (0, 1]$ is the momentum coefficient, and $\mathbf{M}_\ell^{(t)}$ is the momentum retained in the optimizer state. In full-parameter pre-training or fine-tuning of LLMs, the memory requirements for storing momentum in MSGD and the additional variance state in Adam are highly demanding. According to Zhao et al. (2024), pre-training a LLaMA 7B model with a single batch size requires 58 GB of memory, with 42 GB allocated to Adam optimizer states and weight gradients.

GaLore algorithm. To address the memory challenge, Zhao et al. (2024) proposes a Gradient Low-Rank Projection (GaLore) approach that allows full-parameter learning but is much more memory-efficient. The key idea is to project each stochastic gradient $\mathbf{G}_\ell \in \mathbb{R}^{m_\ell \times n_\ell}$ onto a low-rank subspace, yielding a low-dimensional gradient approximation. Specifically, GaLore performs SVD on $\mathbf{G}_\ell^{(t)} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$ and obtains rank- r_ℓ projection matrices $\mathbf{P}_\ell^{(t)} = \mathbf{U}[:, : r_\ell] \in \mathbb{R}^{m_\ell \times r_\ell}$ and $\mathbf{Q}_\ell^{(t)} = \mathbf{V}[:, : r_\ell] \in \mathbb{R}^{n_\ell \times r_\ell}$, where $[:, : r]$ denotes the selection of the matrix's first r columns. When $m_\ell \leq n_\ell$, GaLore projects \mathbf{G}_ℓ onto \mathbf{P}_ℓ , yielding a low-rank gradient representation $(\mathbf{P}_\ell^{(t)})^\top \mathbf{G}_\ell^{(t)} \in \mathbb{R}^{r_\ell \times n_\ell}$. Conversely, when $m_\ell > n_\ell$, GaLore projects \mathbf{G}_ℓ onto \mathbf{Q}_ℓ , resulting in $\mathbf{G}_\ell^{(t)} \mathbf{Q}_\ell^{(t)} \in \mathbb{R}^{m_\ell \times r_\ell}$. In either scenarios, the memory cost of optimizer states associated with these low-rank representations can be significantly reduced, leading to memory-efficient LLMs pre-training or fine-tuning:

$$\mathbf{X}_\ell^{(t+1)} = \begin{cases} \mathbf{X}_\ell^{(t)} + \mathbf{P}_\ell^{(t)} \rho_\ell^{(t)}((\mathbf{P}_\ell^{(t)})^\top \mathbf{G}_\ell^{(t)}), & \text{if } m_\ell \leq n_\ell; \\ \mathbf{X}_\ell^{(t)} + \rho_\ell^{(t)}(\mathbf{G}_\ell^{(t)} \mathbf{Q}_\ell^{(t)}) (\mathbf{Q}_\ell^{(t)})^\top, & \text{if } m_\ell > n_\ell. \end{cases}$$

Typically, GaLore selects $\rho_\ell(\cdot)$ as the Adam gradient operator, as illustrated in Alg. 1. However, GaLore can also choose $\rho_\ell(\cdot)$ to be gradient operators in either vanilla SGD or MSGD. Since SVD decomposition is computationally expensive, GaLore updates $\mathbf{P}_\ell^{(t)}$ or $\mathbf{Q}_\ell^{(t)}$ periodically. In other words, GaLore computes $\mathbf{P}_\ell^{(t)}$ or $\mathbf{Q}_\ell^{(t)}$ when iteration step $t \not\equiv 0 \pmod{\tau}$ where $\tau > 0$ is the period, otherwise $\mathbf{P}_\ell^{(t)} = \mathbf{P}_\ell^{(t-1)}$ and $\mathbf{Q}_\ell^{(t)} = \mathbf{Q}_\ell^{(t-1)}$ remain unchanged. Both the gradient subspace projection and periodic switches between different low-rank subspaces pose significant challenges to the convergence analysis for GaLore-like algorithms.

Stiefel manifold. An $m \times r$ Stiefel manifold ($r \leq m$) is defined as

$$\text{St}_{m,r} = \{\mathbf{P} \in \mathbb{R}^{m \times r} \mid \mathbf{P}^\top \mathbf{P} = \mathbf{I}_r\}.$$

Stiefel manifold is the set of low-rank projection matrices to use in subspace optimization. Typically, in GaLore we have $\mathbf{P}_\ell^{(t)} \in \text{St}_{m_\ell, r_\ell}$ and $\mathbf{Q}_\ell^{(t)} \in \text{St}_{n_\ell, r_\ell}$.

Basic assumptions. We introduce the basic assumptions used throughout our theoretical analysis. Each of these assumptions is standard for stochastic optimization.

Assumption 1 (Lower boundedness). *The objective function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ satisfies $\inf_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) > -\infty$, where $d = \sum_{\ell=1}^{N_\ell} m_\ell n_\ell$ is the total number of parameters in the model.*

Assumption 2 (L -smoothness). *The objective function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ satisfies $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \leq L\|\mathbf{x} - \mathbf{y}\|_2$, for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$.*

Assumption 3 (Stochastic gradient). *The gradient oracle (F, \mathcal{D}) satisfies*

$$\mathbb{E}_{\xi \sim \mathcal{D}}[\nabla_\ell F(\mathbf{x}; \xi)] = \nabla_\ell f(\mathbf{x}), \quad \text{and} \quad \mathbb{E}_{\xi \sim \mathcal{D}}[\|\nabla_\ell F(\mathbf{x}; \xi) - \nabla_\ell f(\mathbf{x})\|_F^2] \leq \sigma_\ell^2, \quad \forall \mathbf{x} \in \mathbb{R}^d,$$

where $\sigma_\ell > 0$ is a scalar. Summing all weight matrices we obtain

$$\mathbb{E}_{\xi \sim \mathcal{D}}[\nabla F(\mathbf{x}; \xi)] = \nabla f(\mathbf{x}), \quad \text{and} \quad \mathbb{E}_{\xi \sim \mathcal{D}}[\|\nabla F(\mathbf{x}; \xi) - \nabla f(\mathbf{x})\|_2^2] \leq \sigma^2, \quad \forall \mathbf{x} \in \mathbb{R}^d,$$

where $\sigma = \sqrt{\sum_{\ell=1}^{N_\ell} \sigma_\ell^2}$.

3 NON-CONVERGENCE OF GaLore: INTUITION AND COUNTER-EXAMPLE

In this section, we demonstrate why GaLore cannot guarantee exact convergence under Assumptions 1-3. We first illustrate the insight behind the result, then present its formal description.

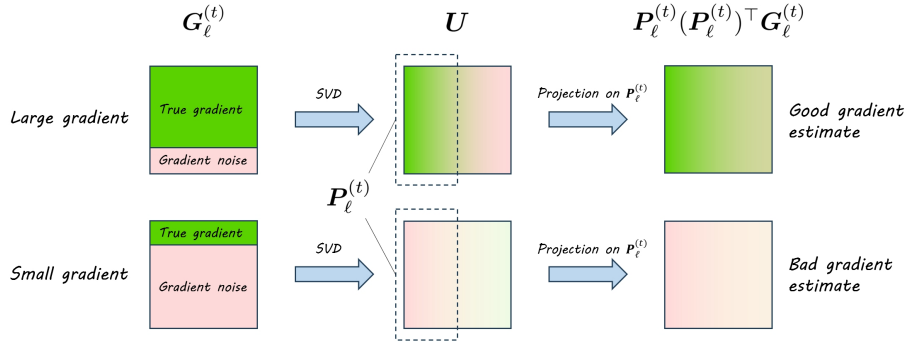


Figure 2: An illustration of the insight on why GaLore fails to converge in small-gradient scenarios. We use color green for true gradient and red for gradient noise.

Insight behind non-convergence. As reviewed in Sec. 2, GaLore performs SVD on stochastic gradient $\mathbf{G} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$ and obtains rank- r projection matrices $\mathbf{P} = \mathbf{U}[:, :r] \in \mathbb{R}^{m \times r}$. GaLore projects \mathbf{G} onto \mathbf{P} , yielding a low-rank gradient representation $\mathbf{P}^\top \mathbf{G} \in \mathbb{R}^{r \times n}$. In other words, GaLore projects the stochastic gradient matrix onto a low-rank subspace spanned by the top r singular vectors, capturing the dominant components of the stochastic gradient matrix. However, the stochastic gradient comprises two components: the true gradient and gradient noise, as shown in Fig. 2. When the true gradient significantly exceeds the gradient noise, typically at the start of training, the low-rank subspace obtained via SVD effectively preserves the true gradient information. As training progresses and the true gradient diminishes to zero, especially near a local minimum, the subspace may become increasingly influenced by gradient noise. In the extreme case, this noise-dominated subspace can become orthogonal to the true gradient subspace, leading to non-convergence.

Counter-Example. We consider the following quadratic problem with gradient noise:

$$f(\mathbf{X}) = \frac{1}{2} \|\mathbf{A}\mathbf{X}\|_F^2 + \langle \mathbf{B}, \mathbf{X} \rangle_F, \quad \nabla F(\mathbf{X}; \xi) = \nabla f(\mathbf{X}) + \xi \sigma \mathbf{C}, \quad (1)$$

where $\mathbf{A} = (\mathbf{I}_{n-r} \quad \mathbf{0}) \in \mathbb{R}^{(n-r) \times n}$, $\mathbf{B} = \begin{pmatrix} \mathbf{D} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \in \mathbb{R}^{n \times n}$ with $\mathbf{D} \in \mathbb{R}^{(n-r) \times (n-r)}$ generated randomly, $\mathbf{C} = \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_r \end{pmatrix} \in \mathbb{R}^{n \times n}$, ξ is a random variable uniformly sampled from $\{1, -1\}$ per iteration, and σ is used to control the gradient noise. It is straightforward to verify that problem (1) satisfies Assumptions 1-3. Moreover, as \mathbf{X} approaches the global minimum of $f(\mathbf{X})$, the true gradient $\nabla f(\mathbf{X}) \rightarrow \mathbf{0}$, while the gradient noise persists with a variance on the order of σ^2 . Fig. 1

illustrates the performance of GaLore when solving problem (1). It is observed that GaLore fails to converge to the optimal solution, regardless of whether the AdamW or MSGD optimizer is used.

Non-convergence of GaLore. Based on the aforementioned insight, we establish the following theorem regarding the non-convergence of GaLore.

Theorem 1 (Non-convergence of GaLore). *There exists an objective function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ satisfying Assumptions 1, 2, a stochastic gradient oracle (F, \mathcal{D}) satisfying Assumption 3, an initial point $\mathbf{x}^{(0)} \in \mathbb{R}^d$, a constant $\epsilon_0 > 0$ such that for any rank $r_\ell < \min\{m_\ell, n_\ell\}$, subspace changing frequency τ , any subspace optimizer ρ inputting subspace gradient of shape $r_\ell \times n_\ell$ and outputting subspace update direction of shape $r_\ell \times n_\ell$ with arbitrary hyperparameters and any $t > 0$, it holds that*

$$\|\nabla f(\mathbf{x}^{(t)})\|_2^2 \geq \epsilon_0.$$

4 CONDITIONS UNDER WHICH GALORE CAN CONVERGE

GaLore provably converges in the noise-free setting. According to the insight presented in Sec. 3, GaLore fails to converge when gradient noise dominates the true gradient in magnitudes. This motivates us to examine the deterministic scenario where the true gradient $\nabla f(\mathbf{x})$ can be accessed without any gradient noise. The GaLore algorithm with noise-free gradients is presented in Alg. 1 (or Alg. 2 in Appendix B.3), where the true gradient oracle is highlighted with the label **(deterministic)**.

Since no gradient noise exists, the projection matrix $\mathbf{P}_\ell^{(t)}$ obtained by SVD can effectively capture the true gradient even when the algorithm approaches a local minimum. For simplicity, we analyze GaLore with MSGD and the following momentum updating mechanism:

$$\mathbf{M}_\ell^{(t)} = \begin{cases} (1 - \beta_1)(\mathbf{P}_\ell^{(t)})^\top \mathbf{P}_\ell^{(t-1)} \mathbf{M}_\ell^{(t-1)} + \beta_1 (\mathbf{P}_\ell^{(t)})^\top \mathbf{G}_\ell^{(t)}, & \text{if } m_\ell \leq n_\ell, \\ (1 - \beta_1) \mathbf{M}_\ell^{(t-1)} (\mathbf{Q}_\ell^{(t-1)})^\top \mathbf{Q}_\ell^{(t)} + \beta_1 \mathbf{G}_\ell^{(t)} \mathbf{Q}_\ell^{(t)}, & \text{if } m_\ell > n_\ell. \end{cases} \quad (2)$$

If the subspace does not change at iteration t , $(\mathbf{P}_\ell^{(t)})^\top \mathbf{P}_\ell^{(t-1)} = (\mathbf{Q}_\ell^{(t-1)})^\top \mathbf{Q}_\ell^{(t)} = \mathbf{I}_{r_\ell}$ and (2) reduces to regular momentum updates. If the subspace changes at iteration t , we inherit $\mathbf{M}_\ell^{(t-1)}$ by first projecting back to the previous space and then to the new subspace. For convenience, we use *momentum projection (MP)* to refer to mechanism (2). When MP is used in the algorithm, we label the corresponding with **(with MP)** in Alg. 1 otherwise **(without MP)**. The following theorem provides convergence guarantees for GaLore using deterministic gradients and MSGD with MP.

Theorem 2 (Convergence rate of deterministic GaLore). *Under Assumptions 1-2, if the number of iterations $T \geq 64/(3\delta)$ and we choose*

$$\beta_1 = 1, \quad \tau = \left\lceil \frac{64}{3\delta\beta_1} \right\rceil, \quad \text{and} \quad \eta = \left(4L + \sqrt{\frac{80L^2}{3\delta\beta_1^2}} + \sqrt{\frac{80\tau^2 L^2}{3\delta}} + \sqrt{\frac{16\tau L^2}{3\beta_1}} \right)^{-1},$$

GaLore using deterministic gradients and MSGD with MP converges as

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}^{(t)})\|_2^2 = \mathcal{O}\left(\frac{L\Delta}{\delta^{5/2}T}\right),$$

where $\Delta = f(\mathbf{x}^{(0)}) - \inf_{\mathbf{x}} f(\mathbf{x})$ and $\delta := \min_\ell \frac{r_\ell}{\min\{m_\ell, n_\ell\}}$.

Remark. In fact, *MSGD* here reduces to momentum gradient descent by using deterministic gradients. Theorem 2 demonstrates that GaLore converges at a rate of $\mathcal{O}(1/T)$ in the deterministic scenario, which is on the same order as full-parameter training. A more detailed result is presented in Theorem 6 in Appendix B.3, where we established convergence for more general hyperparameter choices. However, in deep learning tasks with exceptionally large training datasets, computing the true gradient becomes impractical due to significant computational and memory costs. Therefore, we will next focus on the stochastic setting.

GaLore provably converges with large-batch stochastic gradients. Inspired by the insight presented in Sec. 3, GaLore converges in cases where the true gradient dominates the gradient noise.

Algorithm 1 GaLore / GoLore algorithm using stochastic / deterministic / large-batch gradients with / without momentum projection

Input: Initial point $\mathbf{x}^{(0)}$, data distribution \mathcal{D} , learning rate η , subspace changing frequency τ , rank $\{r_\ell\}_{\ell=1}^{N_L}$, optimizer hyperparameters $\beta_1, \beta_2, \epsilon$, large batch size \mathcal{B} .

Output: $\{\mathbf{x}^{(t)}\}_{t=0}^T$.

Initialize optimizer state $\{M_\ell^{(-1)}\}_{\ell=1}^{N_L}$ and $\{V_\ell^{(-1)}\}_{\ell=1}^{N_L}$ to zero;

for $t = 0, 1, \dots, T - 1$ **do**

for $\ell = 1, 2, \dots, N_L$ **do**

if $t \equiv 0 \pmod{\tau}$ **then**

$G_\ell^{(t)} \leftarrow \nabla_\ell F(\mathbf{x}^{(t)}; \xi^{(t)});$ (stochastic)

$G_\ell^{(t)} \leftarrow \nabla_\ell f(\mathbf{x}^{(t)});$ (deterministic)

$G_\ell^{(t)} \leftarrow \frac{1}{\mathcal{B}} \sum_{b=1}^{\mathcal{B}} \nabla_\ell F(\mathbf{x}^{(t)}; \xi^{(t,b)});$ (large-batch)

$U, \Sigma, V \leftarrow \text{SVD}(G_\ell^{(t)}), P_\ell^{(t)} \leftarrow U[:, :r_\ell], Q_\ell^{(t)} \leftarrow V[:, :r_\ell];$ (GaLore)

 Sample $P_\ell^{(t)} \sim \mathcal{U}(\text{St}_{m_\ell, r_\ell}), Q_\ell^{(t)} \sim \mathcal{U}(\text{St}_{n_\ell, r_\ell});$ (GoLore)

else

$G_\ell^{(t)} \leftarrow \nabla_\ell F(\mathbf{x}^{(t)}; \xi^{(t)});$ (stochastic)

$G_\ell^{(t)} \leftarrow \nabla_\ell f(\mathbf{x}^{(t)});$ (deterministic)

$G_\ell^{(t)} \leftarrow \nabla_\ell F(\mathbf{x}^{(t)}; \xi^{(t)});$ (large-batch)

$P_\ell^{(t)} \leftarrow P_\ell^{(t-1)}, Q_\ell^{(t)} \leftarrow Q_\ell^{(t-1)};$

end if

$R_\ell^{(t)} \leftarrow \begin{cases} (P_\ell^{(t)})^\top G_\ell^{(t)}, & \text{if } m_\ell \leq n_\ell; \\ G_\ell^{(t)} Q_\ell^{(t)}, & \text{if } m_\ell > n_\ell; \end{cases}$

$M_\ell^{(t)} \leftarrow \begin{cases} (1 - \beta_1)(P_\ell^{(t)})^\top P_\ell^{(t-1)} M_\ell^{(t-1)} + \beta_1 R_\ell^{(t)}, & \text{if } m_\ell \leq n_\ell; \\ (1 - \beta_1)M_\ell^{(t-1)}(Q_\ell^{(t-1)})^\top Q_\ell^{(t)} + \beta_1 R_\ell^{(t)}, & \text{if } m_\ell > n_\ell; \end{cases}$ (with MP)

$M_\ell^{(t)} \leftarrow (1 - \beta_1)M_\ell^{(t-1)} + \beta_1 R_\ell^{(t)};$ (without MP)

$V_\ell^{(t)} \leftarrow (1 - \beta_2)V_\ell^{(t-1)} + \beta_2 R_\ell^{(t)} \odot R_\ell^{(t)};$

if using Adam **then**

$M_\ell^{(t)} \leftarrow M_\ell^{(t)} / (1 - \beta_1^t), V_\ell^{(t)} \leftarrow V_\ell^{(t)} / (1 - \beta_2^t), N_\ell^{(t)} \leftarrow M_\ell^{(t)} / (\sqrt{V_\ell^{(t)}} + \epsilon);$

else if using MSGD **then**

$N_\ell^{(t)} \leftarrow M_\ell^{(t)};$

end if

$X_\ell^{(t+1)} \leftarrow \begin{cases} X_\ell^{(t)} - \eta P_\ell^{(t)} N_\ell^{(t)}, & \text{if } m_\ell \leq n_\ell; \\ X_\ell^{(t)} - \eta N_\ell^{(t)} (Q_\ell^{(t)})^\top, & \text{if } m_\ell > n_\ell; \end{cases}$

end for

end for

This convergence can be ensured by reducing the gradient noise through an increased batch size, particularly as the algorithm approaches a local minimum. Specifically, we replace the stochastic gradient $G_\ell^{(t)} = \nabla_\ell F(\mathbf{x}^{(t)}; \xi^{(t)})$ with large-batch gradient $G_\ell^{(t)} = \frac{1}{\mathcal{B}} \sum_{b=1}^{\mathcal{B}} \nabla_\ell F(\mathbf{x}^{(t)}; \xi^{(t,b)})$, which reduces the variance of gradient noise by \mathcal{B} times. The GaLore algorithm with large-batch stochastic gradients is presented in Alg. 1 (or Alg. 3 in Appendix B.4), where the large-batch stochastic gradient oracle is highlighted with the label (large-batch). It is worth noting that the non-convergence of GaLore primarily stems from the erroneous subspace dominated by gradient noise. Therefore, we compute a large-batch gradient only for the SVD step while maintaining a smaller batch size for

other computations, see Alg. 1. As the batch size \mathcal{B} increases with iteration T , GaLore provably converge to stationary solutions, as established in the following theorem:

Theorem 3 (Convergence rate of large-batch GaLore). *Under Assumptions 1-3, if $T \geq 2 + 256/(3\underline{\delta}) + (256\sigma)^2/(9\sqrt{\underline{\delta}}L\Delta)$ and we choose $\tau = \lceil 128/(3\underline{\delta}\beta_1) \rceil$, $\mathcal{B} = \lceil 1/(\underline{\delta}\beta_1) \rceil$,*

$$\beta_1 = \left(1 + \sqrt{\frac{\underline{\delta}^{3/2}\sigma^2 T}{L\Delta}}\right)^{-1}, \quad \text{and} \quad \eta = \left(4L + \sqrt{\frac{80L^2}{3\underline{\delta}\beta_1^2}} + \sqrt{\frac{40\tau^2 L^2}{\underline{\delta}}} + \sqrt{\frac{32\tau L^2}{3\beta_1}}\right)^{-1},$$

GaLore using large-batch gradients and MSGD with MP converges as

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(\mathbf{x}^{(t)})\|_2^2] = \mathcal{O}\left(\frac{L\Delta}{\underline{\delta}^{5/2}T} + \sqrt{\frac{L\Delta\sigma^2}{\underline{\delta}^{7/2}T}}\right),$$

where $\Delta = f(\mathbf{x}^{(0)}) - \inf_{\mathbf{x}} f(\mathbf{x})$ and $\underline{\delta} := \min_{\ell} \frac{r_{\ell}}{\min\{m_{\ell}, n_{\ell}\}}$.

Remark. A more detailed result is presented in Theorem 7 in Appendix B.4, where we established convergence for more general hyperparameter choices. The batch size $\mathcal{B} = \Theta(\sqrt{T})$ in large-batch GaLore grows with iteration T , leading to increased memory overhead, making it less practical than small-batch GaLore. With gradient accumulation, an additional variable is needed to track the gradient, complicating compatibility with per-layer weight updates. Otherwise, larger batch sizes raise the memory required for activation values. Therefore, exploring algorithms that can converge with standard small-batch stochastic gradients becomes essential.

Empirical validation. Fig. 1 illustrates the convergence of large-batch GaLore (blue curve) in solving problem (1). It demonstrates that large-batch GaLore effectively corrects the bias present in small-batch stochastic GaLore (green curve), achieving convergence to stationary solutions.

GaLore provably converges with isotropic noise assumptions. In Appendix G, we further prove that under some additional isotropic noise assumptions, GaLore with small-batch stochastic gradients can also be guaranteed to converge at a rate of $\mathcal{O}(1/\sqrt{T})$.

5 GOLORE: GRADIENT RANDOM LOW-RANK PROJECTION

GoLore algorithm. The main issue with SVD-based projection in GaLore is that it aims to capture the dominant component in the stochastic gradient matrix. Consequently, when gradient noise overshadows the true gradient as the algorithm approaches a local minimum, the SVD-based projection fails to identify valuable gradient information.

To address this, we propose replacing the SVD-based projection with a random projection, which captures components of the stochastic gradient matrix randomly without any preference. This results in the GoLore algorithm presented in Alg. 1 (or Alg. 4 in Appendix B.5). In Alg. 1, the GaLore method highlighted with the label **(GaLore)** samples the projection matrix $\mathbf{P}_{\ell}^{(t)}$ via SVD decomposition. In contrast, the GoLore method highlighted with the label **(GoLore)** samples $\mathbf{P}_{\ell}^{(t)}$ from $\mathcal{U}(\text{St}_{m_{\ell}, r_{\ell}})$, a uniform distribution on the $m_{\ell} \times r_{\ell}$ Stiefel manifold. The following proposition provides a practical strategy to sample from distribution $\mathcal{U}(\text{St}_{m, r})$.

Proposition 1 (Chikuse (2012), Theorem 2.2.1). *A random matrix \mathbf{X} uniformly distributed on $\text{St}_{m, r}$ is expressed as $\mathbf{X} = \mathbf{Z}(\mathbf{Z}^{\top} \mathbf{Z})^{-1/2}$, where the elements of an $m \times r$ random matrix \mathbf{Z} are independent and identically distributed as normal $\mathcal{N}(0, 1)$.*

Convergence guarantee. Unlike SVD used in GaLore, the random sampling strategy in GoLore prevents the subspace from being dominated by gradient noise. The theorem below provides convergence guarantees for GoLore when using small-batch stochastic gradients and MSGD with MP.

Theorem 4 (Convergence rate of GoLore). *Under Assumptions 1-3, for any $T \geq 2 + 128/(3\underline{\delta}) + (128\sigma)^2/(9\sqrt{\underline{\delta}}L\Delta)$, if we choose $\tau = \lceil 64/(3\underline{\delta}\beta_1) \rceil$,*

$$\beta_1 = \left(1 + \sqrt{\frac{\underline{\delta}^{3/2}\sigma^2 T}{L\Delta}}\right)^{-1}, \quad \text{and} \quad \eta = \left(4L + \sqrt{\frac{80L^2}{3\underline{\delta}\beta_1^2}} + \sqrt{\frac{80\tau^2 L^2}{3\underline{\delta}}} + \sqrt{\frac{16\tau L^2}{3\beta_1}}\right)^{-1},$$

Table 1: Memory and computation comparison between GaLore’s original implementation and our ReLoRA-like version, both utilizing MSGD with batch size b . We assume the weight $\mathbf{W} \in \mathbb{R}^{m \times n}$ satisfies $m \leq n$.

GaLore Implementation	Memory	Computation
(Zhao et al., 2024)	$mn + rm + rn + bm$	$6bmn + 4rmn + 2mn + 3rn$
Our ReLoRA-like version	$mn + rm + 2rn + bm + br$	$4bmn + 4brm + 6brn + 5rn$

GoLore using small-batch stochastic gradients and MSGD with MP converges as

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(\mathbf{x}^{(t)})\|_2^2] = \mathcal{O}\left(\frac{L\Delta}{\underline{\delta}^{5/2}T} + \sqrt{\frac{L\Delta\sigma^2}{\underline{\delta}^{7/2}T}}\right),$$

where $\Delta = f(\mathbf{x}^{(0)}) - \inf_{\mathbf{x}} f(\mathbf{x})$ and $\underline{\delta} := \min_{\ell} \frac{r_{\ell}}{\min\{m_{\ell}, n_{\ell}\}}$.

Remark. Theorem 4 demonstrates that GaLore converges at a rate of $\mathcal{O}(1/\sqrt{T})$, which is consistent with the convergence rate of full-parameter pre-training using standard MSGD. A more detailed result is presented in Theorem 8 in Appendix B.5, where we established convergence for more general hyperparameter choices. Unlike deterministic GaLore and low-rank GaLore discussed in Sec. 4, the newly-proposed GoLore algorithm converges in the non-convex stochastic setting with regular batch sizes, making it far more suitable for LLM pre-training and fine-tuning.

Practical application of GoLore in LLMs. While GoLore have theoretical convergence guarantees, directly applying GoLore in LLM tasks may not be ideal. The advantage of using randomly sampled projection matrices becomes evident in the later stages of training, where stochastic gradients are primarily dominated by gradient noise. However, in the early stages, projection matrices derived from SVD retain more gradient information, leading to more effective subspaces. Therefore, we recommend a *hybrid* approach: initially using GaLore to converge toward the neighborhood of the solution, then switching to GoLore for refinement and achieving more accurate results.

Empirical validation. Fig. 1 shows the convergence of the hybrid algorithm (red curve) applied to problem (1), which employs GaLore during the early training phase and switches to GoLore in the later stage. It is observed that the hybrid algorithm successfully converges to stationary solutions.

6 CONNECTION WITH OTHER SUBSPACE OPTIMIZATION METHODS

Connection with ReLoRA. Algorithms like GaLore/GoLore that optimizes in periodically recomputed subspaces can be implemented in an equivalent yet potentially more computational efficient, ReLoRA-like way. Consider a linear layer $\mathbf{y} = \mathbf{W}\mathbf{x}$ with $\mathbf{W} \in \mathbb{R}^{m \times n}$, where $m \leq n$, GaLore first computes the full-parameter gradient $\nabla_{\mathbf{W}}\mathcal{L} = (\nabla_{\mathbf{y}}\mathcal{L})\mathbf{x}^{\top}$ via back propagation and update \mathbf{W} in the subspace as $\mathbf{W} \leftarrow \mathbf{W} + \mathbf{P}\rho(\mathbf{P}^{\top}(\nabla_{\mathbf{W}}\mathcal{L}))$, where $\mathbf{P} \in \mathbb{R}^{m \times r}$ is a low-rank projection matrix. If we use LoRA adaptation $\mathbf{W} = \mathbf{W}_0 + \mathbf{B}\mathbf{A}$ with $\mathbf{B} \in \mathbb{R}^{m \times r}$ and $\mathbf{A} \in \mathbb{R}^{r \times n}$, we compute \mathbf{A} ’s gradient $\nabla_{\mathbf{A}}\mathcal{L} = (\nabla_{\mathbf{z}}\mathcal{L})\mathbf{x}^{\top} = \mathbf{B}^{\top}(\nabla_{\mathbf{y}}\mathcal{L})\mathbf{x}^{\top}$, where $\mathbf{z} = \mathbf{B}\mathbf{x}$ is the additional activation. If we fix $\mathbf{B} = \mathbf{P}$, update $\mathbf{A} \leftarrow \mathbf{A} + \rho(\nabla_{\mathbf{A}}\mathcal{L})$ is equivalent to $\mathbf{W} \leftarrow \mathbf{W} + \mathbf{P}\rho(\mathbf{P}^{\top}(\nabla_{\mathbf{W}}\mathcal{L}))$. The memory and computational costs of the two implementations are compared in Table 1, showing the potential of our ReLoRA-like implementation to reduce computation with little memory overhead. Detailed algorithm descriptions and calculations are in Appendix D.

Connection with FLORA. Aware of the equivalence of the two (GaLore/ReLoRA-like) implementations, the main difference between GoLore and FLORA lies in the choice of projection matrices. Though both algorithms sample $\mathbf{P} \in \mathbb{R}^{m \times r}$ randomly, GoLore uses a uniform distribution on the Stiefel manifold $\mathcal{U}(\text{St}_{m,r})$, while FLORA uses a random Gaussian distribution where each element in \mathbf{P} is independently sampled from $\mathcal{N}(0, 1/r)$, and thus \mathbf{P} may not belongs to $\text{St}_{m,r}$.

Connection with SIFT. SIFT fine-tunes LLMs with sparsified gradients, which can also be viewed as [subspace descent](#). While GaLore projects gradient \mathbf{G} to $\mathbf{P}^{\top}\mathbf{G}$ via a projection matrix \mathbf{P} , SIFT projects gradient \mathbf{G} to $\mathbf{S} \odot \mathbf{G}$ via a sparse mask matrix \mathbf{S} . Our theoretical analysis can be directly transferred to sparse [subspace descent](#) with little effort, implying similar results as in low-rank [subspace descent](#), see Appendix C.

7 EXPERIMENTS

We evaluate GaLore and GoLore on several different tasks, including solving a counter-example problem (1), pre-training and fine-tuning LLMs with real benchmarks. Throughout our experi-

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

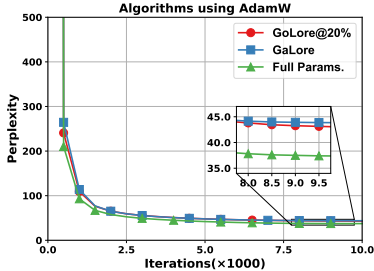


Figure 3: Pre-training curves of various approaches using AdamW with BF16 precision.

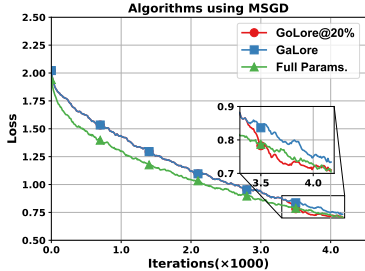


Figure 4: Fine-tuning curves of various approaches using MSGD with BF16 precision.

Table 2: Fine-tuning results on GLUE benchmark using pre-trained RoBERTa-Base.

Algorithm	CoLA	STS-B	MRPC	RTE	SST2	MNLI	QNLI	QQP	Avg
Full Params.	62.07	90.18	92.25	78.34	94.38	87.59	92.46	91.90	86.15
GaLore	61.32	90.24	92.55	77.62	94.61	86.92	92.06	90.84	85.77
FLORA	57.71	89.59	91.96	76.17	94.50	85.42	91.93	90.49	84.72
GoLore@20%	61.66	90.55	92.93	78.34	94.61	87.02	92.20	90.91	86.03

ments, *GoLore@x%* uses GaLore in the first $(100 - x)\%$ iterations and GoLore in the last $x\%$ iterations, *L.B. GaLore* denotes large-batch GaLore, and *Full Params.* denotes full-parameter training. Further results and detailed experimental specifications including the hyperparameter choices and computing resources are deferred to Appendix E.

GoLore’s non-convergence. To validate the non-convergence of GaLore and the convergence properties of GoLore and large-batch GaLore, we compare them with full-parameter training on the constructed quadratic problem defined in (1). Fig. 1 shows that, regardless of whether AdamW or MSGD is employed as the subspace optimizer, GaLore does not converge to the desired solution. In contrast, both GoLore and large-batch GaLore, along with full-parameter training, achieve exact convergence, thereby validating our theoretical results.

Pre-training. To validate the efficiency of GoLore in LLM pre-training tasks, we pre-trained LLaMA-60M on the C4 (Raffel et al., 2020) dataset for 10,000 iterations using various algorithms, including GaLore, GoLore and full-parameter training. All implementations utilized the AdamW optimizer in BF16 format. As illustrated in Fig. 3, there is a noticeable performance gap between GaLore/GoLore and full-parameter training, indicating that the parameters are away from local minima. However, GoLore still demonstrates slightly better training performance compared to GaLore.

Fine-tuning. To validate the efficiency of GoLore in LLM fine-tuning tasks, we fine-tuned pre-trained LLaMA2-7B models (Touvron et al., 2023) on the WinoGrande dataset (Sakaguchi et al., 2021) and pre-trained RoBERTa models (Liu, 2019) on the GLUE benchmark (Wang, 2018) with AdamW optimizers. Fig. 4 displays the fine-tuning loss curves for GaLore and GoLore with rank 1024, while Table 2 presents the task scores for GaLore/GoLore and **FLORA** with rank 4. In both experiments, GoLore outperforms GaLore.

8 CONCLUSION AND LIMITATIONS

This paper investigates subspace optimization approaches for LLM pre-training and fine-tuning. We demonstrate that GaLore fails to converge to the desired solution under regular assumptions, as the SVD-based projection often generates noise-dominated subspaces when the true gradient is relatively small. However, we establish that GaLore can achieve exact convergence when using deterministic or large-batch stochastic gradients. We further introduce GoLore—a variant of GaLore employing randomly sampled projection matrices—and establish its convergence rate even with small-batch stochastic gradients. A limitation of this paper is that convergence guarantees for GoLore are currently provided only when using MSGD as the subspace optimizer. Although GoLore with AdamW performs well empirically, as shown in Table 2, its theoretical convergence guarantees remain unknown and will be addressed in future work.

REFERENCES

- 540
541
542 Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. Qsgd:
543 Communication-efficient sgd via gradient quantization and encoding. *Advances in neural in-*
544 *formation processing systems*, 30, 2017.
- 545 Massimo Bini, Karsten Roth, Zeynep Akata, and Anna Khoreva. Ether: Efficient finetuning of
546 large-scale models with hyperplane reflections. *arXiv preprint arXiv:2405.20271*, 2024.
- 547 Yiming Chen, Yuan Zhang, Liyuan Cao, Kun Yuan, and Zaiwen Wen. Enhancing zeroth-order fine-
548 tuning for language models with low-rank structures. *arXiv preprint arXiv:2410.07698*, 2024.
- 549
550 Yasuko Chikuse. *Statistics on special manifolds*, volume 174. Springer Science & Business Media,
551 2012.
- 552 Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina
553 Toutanova. Boolq: Exploring the surprising difficulty of natural yes/no questions. *arXiv preprint*
554 *arXiv:1905.10044*, 2019.
- 555
556 Laurent Condat, Artavazd Maranjyan, and Peter Richtárik. Locodl: Communication-efficient dis-
557 tributed learning with local training and compression. *arXiv preprint arXiv:2403.04348*, 2024.
- 558 Ilyas Fatkhullin, Alexander Tyurin, and Peter Richtárik. Momentum provably improves error feed-
559 back! *Advances in Neural Information Processing Systems*, 36, 2024.
- 560
561 Eduard Gorbunov, Konstantin P Burlachenko, Zhize Li, and Peter Richtárik. Marina: Faster non-
562 convex distributed learning with compression. In *International Conference on Machine Learning*,
563 pp. 3788–3798. PMLR, 2021.
- 564
565 Andi Han, Jiayang Li, Wei Huang, Mingyi Hong, Akiko Takeda, Pratik Jawanpuria, and Bamdev
566 Mishra. Sltrain: a sparse plus low-rank approach for parameter and memory efficient pretraining.
567 *arXiv preprint arXiv:2406.02214*, 2024.
- 568 Filip Hanzely, Konstantin Mishchenko, and Peter Richtárik. Sega: Variance reduction via gradient
569 sketching. *Advances in Neural Information Processing Systems*, 31, 2018.
- 570
571 Yongchang Hao, Yanshuai Cao, and Lili Mou. Flora: Low-rank adapters are secretly gradient
572 compressors. *arXiv preprint arXiv:2402.03293*, 2024.
- 573 Soufiane Hayou, Nikhil Ghosh, and Bin Yu. Lora+: Efficient low rank adaptation of large models.
574 *arXiv preprint arXiv:2402.12354*, 2024.
- 575
576 Yutong He, Xinmeng Huang, Yiming Chen, Wotao Yin, and Kun Yuan. Lower bounds and accel-
577 erated algorithms in distributed stochastic optimization with communication compression. *arXiv*
578 *preprint arXiv:2305.07612*, 2023.
- 579 Yutong He, Jie Hu, Xinmeng Huang, Songtao Lu, Bin Wang, and Kun Yuan. Distributed bilevel op-
580 timization with communication compression. In *Forty-first International Conference on Machine*
581 *Learning*, 2024a.
- 582
583 Yutong He, Xinmeng Huang, and Kun Yuan. Unbiased compression saves communication in dis-
584 tributed optimization: when and how much? *Advances in Neural Information Processing Systems*,
585 36, 2024b.
- 586 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
587 and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint*
588 *arXiv:2106.09685*, 2021.
- 589
590 Kun Huang and Shi Pu. Cedas: A compressed decentralized stochastic gradient method with im-
591 proved convergence. *arXiv preprint arXiv:2301.05872*, 2023.
- 592
593 Ziyu Jiang, Xuxi Chen, Xueqin Huang, Xianzhi Du, Denny Zhou, and Zhangyang Wang. Back
razor: Memory-efficient transfer learning by self-sparsified backpropagation. *Advances in neural*
information processing systems, 35:29248–29261, 2022.

- 594 Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*,
595 2014.
- 596
- 597 David Kozak, Stephen Becker, Alireza Doostan, and Luis Tenorio. Stochastic subspace descent.
598 *arXiv preprint arXiv:1904.01145*, 2019.
- 599 Zhize Li and Peter Richtárik. Canita: Faster rates for distributed convex optimization with com-
600 munication compression. *Advances in Neural Information Processing Systems*, 34:13770–13781,
601 2021.
- 602
- 603 Zhize Li, Dmitry Kovalev, Xun Qian, and Peter Richtárik. Acceleration for compressed gradient
604 descent in distributed and federated optimization. *arXiv preprint arXiv:2002.11364*, 2020.
- 605 Vladislav Lialin, Sherin Muckatira, Namrata Shivagunde, and Anna Rumshisky. Relora: High-
606 rank training through low-rank updates. In *The Twelfth International Conference on Learning*
607 *Representations*, 2023.
- 608
- 609 Kaizhao Liang, Bo Liu, Lizhang Chen, and Qiang Liu. Memory-efficient llm training with online
610 subspace descent. *arXiv preprint arXiv:2408.12857*, 2024.
- 611 Shih-Yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang-
612 Ting Cheng, and Min-Hung Chen. Dora: Weight-decomposed low-rank adaptation. *arXiv*
613 *preprint arXiv:2402.09353*, 2024.
- 614 Yinhan Liu. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint*
615 *arXiv:1907.11692*, 2019.
- 616
- 617 Sebastian Loeschke, Mads Tofttrup, Michael J Kastoryano, Serge Belongie, and Vésteinn Snæb-
618 jarnarson. Loqt: Low rank adapters for quantized training. *arXiv preprint arXiv:2405.16528*,
619 2024.
- 620 I Loshchilov. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- 621
- 622 Yang Luo, Xiaozhe Ren, Zangwei Zheng, Zhuo Jiang, Xin Jiang, and Yang You. Came: Confidence-
623 guided adaptive memory efficient optimization. *arXiv preprint arXiv:2307.02047*, 2023.
- 624 Sadhika Malladi, Tianyu Gao, Eshaan Nichani, Alex Damian, Jason D Lee, Danqi Chen, and Sanjeev
625 Arora. Fine-tuning language models with just forward passes. *Advances in Neural Information*
626 *Processing Systems*, 36:53038–53075, 2023.
- 627
- 628 Konstantin Mishchenko, Eduard Gorbunov, Martin Takáč, and Peter Richtárik. Distributed learning
629 with compressed gradient differences. *arXiv preprint arXiv:1901.09269*, 2019.
- 630 Aashiq Muhamed, Oscar Li, David Woodruff, Mona Diab, and Virginia Smith. Grass: Compute effi-
631 cient low-memory llm training with structured sparse gradients. *arXiv preprint arXiv:2406.17660*,
632 2024.
- 633
- 634 Rui Pan, Xiang Liu, Shizhe Diao, Renjie Pi, Jipeng Zhang, Chi Han, and Tong Zhang. Lisa: Layer-
635 wise importance sampling for memory-efficient large language model fine-tuning. *arXiv preprint*
636 *arXiv:2403.17919*, 2024.
- 637
- 638 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi
639 Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text
transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- 640 Amrutha Varshini Ramesh, Vignesh Ganapathiraman, Issam H Laradji, and Mark Schmidt. Block-
641 llm: Memory-efficient adaptation of llms by selecting and optimizing the right coordinate blocks.
642 *arXiv preprint arXiv:2406.17296*, 2024.
- 643
- 644 Peter Richtárik, Igor Sokolov, and Ilyas Fatkhullin. Ef21: A new, simpler, theoretically better,
645 and practically faster error feedback. *Advances in Neural Information Processing Systems*, 34:
646 4384–4396, 2021.
- 647
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adver-
sarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021.

- 648 Weixi Song, Zuchao Li, Lefei Zhang, Hai Zhao, and Bo Du. Sparse is enough in fine-tuning pre-
649 trained large language model. *arXiv preprint arXiv:2312.11875*, 2023.
650
- 651 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-
652 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-
653 tion and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 654 Nikhil Vyas, Depen Morwani, and Sham M Kakade. Adamem: Memory efficient momentum for
655 adafactor. In *2nd Workshop on Advancing Neural Network Training: Computational Efficiency,
656 Scalability, and Resource Optimization (WANT@ ICML 2024)*, 2024.
657
- 658 Alex Wang. Glue: A multi-task benchmark and analysis platform for natural language understand-
659 ing. *arXiv preprint arXiv:1804.07461*, 2018.
- 660 Yilong Wang, Haishan Ye, Guang Dai, and Ivor Tsang. Can gaussian sketching converge faster on
661 a preconditioned landscape? In *Forty-first International Conference on Machine Learning*, 2024.
662
- 663 Cong Xie, Shuai Zheng, Sanmi Koyejo, Indranil Gupta, Mu Li, and Haibin Lin. Cser:
664 Communication-efficient sgd with error reset. *Advances in Neural Information Processing Sys-
665 tems*, 33:12593–12603, 2020.
- 666 Zhiyuan Yu, Li Shen, Liang Ding, Xinmei Tian, Yixin Chen, and Dacheng Tao. Sheared back-
667 propagation for fine-tuning foundation models. In *Proceedings of the IEEE/CVF Conference on
668 Computer Vision and Pattern Recognition*, pp. 5883–5892, 2024.
- 669 Liang Zhang, Kiran Koshy Thekumparampil, Sewoong Oh, and Niao He. Dpzero: dimension-
670 independent and differentially private zeroth-order optimization. In *International Workshop on
671 Federated Learning in the Age of Foundation Models in Conjunction with NeurIPS 2023*, 2023.
672
- 673 Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christo-
674 pher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt
675 Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer.
676 Opt: Open pre-trained transformer language models, 2022.
- 677 Jiawei Zhao, Zhenyu Zhang, Beidi Chen, Zhangyang Wang, Anima Anandkumar, and Yuandong
678 Tian. Galore: Memory-efficient llm training by gradient low-rank projection. *arXiv preprint
679 arXiv:2403.03507*, 2024.
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

APPENDIX

A CHALLENGES IN THEORETICAL ANALYSIS

Gradient projection onto a low-rank subspace poses two significant challenges for the convergence analysis of (momentum) stochastic gradient descent:

- **Neither unbiased nor contractive compression.** gradient projection onto this subspace can be viewed as gradient compression. Traditional analyses of optimization algorithms with lossy compression typically rely on either unbiased (Li et al., 2020; Li & Richtárik, 2021; Huang & Pu, 2023; He et al., 2024a;b; Condat et al., 2024) compressibility, *i.e.*, the compressor \mathcal{C} satisfies

$$\mathbb{E}[\mathcal{C}(\mathbf{x})] = \mathbf{x}, \quad \mathbb{E}[\|\mathcal{C}(\mathbf{x}) - \mathbf{x}\|_2^2] \leq \omega \|\mathbf{x}\|_2^2, \quad \forall \mathbf{x} \in \mathbb{R}^d,$$

for some $\omega \geq 0$, or contractive (Richtárik et al., 2021; Xie et al., 2020; Fatkhullin et al., 2024; He et al., 2023) compressibility, *i.e.*,

$$\mathbb{E}[\|\mathcal{C}(\mathbf{x}) - \mathbf{x}\|_2^2] \leq (1 - \delta) \|\mathbf{x}\|_2^2, \quad \forall \mathbf{x} \in \mathbb{R}^d,$$

for some $\delta \in (0, 1]$. However, GaLore’s subspace compression is neither unbiased nor contractive due to the reuse of projection matrices. For example, consider a pre-computed projection matrix $\mathbf{P} \in \mathbb{R}^{m \times r}$. There exists a full-parameter gradient $\mathbf{G} \in \mathbb{R}^{m \times n}$ such that $\mathbf{G} \neq 0$ and $\mathcal{C}(\mathbf{G}) := \mathbf{P}\mathbf{P}^\top \mathbf{G} = 0$, violating both unbiased and contractive compressibility.

- **Periodically projected optimizer states.** When GaLore changes the subspace, the retained momentum terms must be adjusted to track the gradients in the new subspace. Since these momentum terms were initially aligned with the gradients in the original subspace, such adjustments inevitably introduce additional errors, especially when the two subspaces differ significantly. In the extreme case where the two subspaces are entirely orthogonal, the momentum from the previous subspace becomes largely irrelevant for optimization in the new one.

B THEORETICAL PROOFS

B.1 NOTATIONS AND USEFUL LEMMAS

We assume the model parameters consist of N_L weight matrices. We use $\mathbf{X}_\ell \in \mathbb{R}^{m_\ell \times n_\ell}$ to denote the ℓ -th weight matrix and $\mathbf{x} \in \mathbb{R}^d = (\text{vec}(\mathbf{X}_1)^\top, \dots, \text{vec}(\mathbf{X}_{N_L})^\top)^\top$ to denote the vector collecting all the parameters, $d = \sum_{\ell=1}^{N_L} m_\ell n_\ell$. We assume GaLore/GoLore applies rank- r_ℓ projection to the ℓ -th weight matrix and denote

$$\delta_\ell = \frac{r_\ell}{\min\{m_\ell, n_\ell\}}, \quad \underline{\delta} = \min_{1 \leq \ell \leq N_L} \delta_\ell, \quad \bar{\delta} = \max_{1 \leq \ell \leq N_L} \delta_\ell.$$

We define $\tilde{\mathbf{M}}_\ell^{(t)}$ as

$$\tilde{\mathbf{M}}_\ell^{(t)} = \begin{cases} \mathbf{P}_\ell^{(t)} \mathbf{M}_\ell^{(t)}, & \text{if } m_\ell \leq n_\ell, \\ \mathbf{M}_\ell^{(t)} (\mathbf{Q}_\ell^{(t)})^\top, & \text{if } m_\ell > n_\ell, \end{cases}$$

and $\tilde{\mathbf{m}} = (\text{vec}(\tilde{\mathbf{M}}_1)^\top, \dots, \text{vec}(\tilde{\mathbf{M}}_{N_L})^\top)^\top$. While using Alg. 1 with MSGD and MP, it holds for $m_\ell \leq n_\ell$ that

$$\tilde{\mathbf{M}}_\ell^{(t)} = \begin{cases} \beta_1 \mathbf{P}_\ell^{(0)} (\mathbf{P}_\ell^{(0)})^\top \mathbf{G}_\ell^{(0)}, & t = 0; \\ \mathbf{P}_\ell^{(t)} (\mathbf{P}_\ell^{(t)})^\top \left((1 - \beta_1) \tilde{\mathbf{M}}_\ell^{(t-1)} + \beta_1 \mathbf{G}_\ell^{(t)} \right), & t = k\tau, k \in \mathbb{N}^*; \\ (1 - \beta_1) \tilde{\mathbf{M}}_\ell^{(t-1)} + \beta_1 \mathbf{P}_\ell^{(t)} (\mathbf{P}_\ell^{(t)})^\top \mathbf{G}_\ell^{(t)}, & t = k\tau + r, k \in \mathbb{N}, 1 \leq r < \tau; \end{cases}$$

for $m_\ell > n_\ell$ that

$$\tilde{\mathbf{M}}_\ell^{(t)} = \begin{cases} \beta_1 \mathbf{G}_\ell^{(0)} \mathbf{Q}_\ell^{(0)} (\mathbf{Q}_\ell^{(0)})^\top, & t = 0; \\ \left((1 - \beta_1) \tilde{\mathbf{M}}_\ell^{(t-1)} + \beta_1 \mathbf{G}_\ell^{(t)} \right) \mathbf{Q}_\ell^{(t)} (\mathbf{Q}_\ell^{(t)})^\top, & t = k\tau, k \in \mathbb{N}^*; \\ (1 - \beta_1) \tilde{\mathbf{M}}_\ell^{(t-1)} + \beta_1 \mathbf{G}_\ell^{(t)} \mathbf{Q}_\ell^{(t)} (\mathbf{Q}_\ell^{(t)})^\top, & t = k\tau + r, k \in \mathbb{N}, 1 \leq r < \tau; \end{cases}$$

and for both cases that

$$\mathbf{X}_\ell^{(t+1)} = \mathbf{X}_\ell^{(t)} - \eta \tilde{\mathbf{M}}_\ell^{(t)}.$$

Lemma 1 (Error of GaLore’s projection). *Let $\mathbf{G} = \mathbf{U}\Sigma\mathbf{V}^\top$ be the SVD of $\mathbf{G} \in \mathbb{R}^{m \times n}$, projection matrix $\mathbf{P} = \mathbf{U}[:, :r]$, $\mathbf{Q} = \mathbf{V}[:, :r]$, $r < \min\{m, n\}$. It holds for $m \leq n$ that*

$$\|\mathbf{P}\mathbf{P}^\top\mathbf{G} - \mathbf{G}\|_F^2 \leq \left(1 - \frac{r}{m}\right) \|\mathbf{G}\|_F^2,$$

and for $m > n$ that

$$\|\mathbf{G}\mathbf{Q}\mathbf{Q}^\top - \mathbf{G}\|_F^2 \leq \left(1 - \frac{r}{n}\right) \|\mathbf{G}\|_F^2.$$

Proof. Without loss of generality assume $m \leq n$ (the other case can be proved similarly). Let $\mathbf{Q} = \mathbf{U}[:, (r+1):]$, It holds that $\mathbf{I} = \mathbf{U}\mathbf{U}^\top = \mathbf{P}\mathbf{P}^\top + \mathbf{Q}\mathbf{Q}^\top$. Thus,

$$\begin{aligned} \|\mathbf{P}\mathbf{P}^\top\mathbf{G} - \mathbf{G}\|_F^2 &= \|(\mathbf{I} - \mathbf{P}\mathbf{P}^\top)\mathbf{U}\Sigma\mathbf{V}^\top\|_F^2 \\ &= \text{tr}(\mathbf{V}\Sigma^\top\mathbf{U}^\top(\mathbf{I} - \mathbf{P}\mathbf{P}^\top)^2\mathbf{U}\Sigma\mathbf{V}^\top) \\ &= \text{tr}(\Sigma^\top\mathbf{U}^\top\mathbf{Q}\mathbf{Q}^\top\mathbf{U}\Sigma), \end{aligned} \quad (3)$$

where the second equation uses $\|\mathbf{X}\|_F^2 = \text{tr}(\mathbf{X}^\top\mathbf{X})$ and the last equation uses $\text{tr}(\mathbf{A}\mathbf{B}) = \text{tr}(\mathbf{B}\mathbf{A})$, $\mathbf{V}^\top\mathbf{V} = \mathbf{I}$ and $\mathbf{Q}^\top\mathbf{Q} = \mathbf{I}$. By $\mathbf{Q}^\top\mathbf{P} = 0$ and $\mathbf{P}^\top\mathbf{Q} = 0$, we have

$$\mathbf{U}^\top\mathbf{Q}\mathbf{Q}^\top\mathbf{U} = \begin{pmatrix} \mathbf{P}^\top \\ \mathbf{Q}^\top \end{pmatrix} \mathbf{Q}\mathbf{Q}^\top \begin{pmatrix} \mathbf{P} & \mathbf{Q} \end{pmatrix} = \begin{pmatrix} \mathbf{0}_{r \times r} & \mathbf{0}_{r \times (m-r)} \\ \mathbf{0}_{(m-r) \times r} & \mathbf{I}_{m-r} \end{pmatrix}. \quad (4)$$

Let $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_m \geq 0$ denote the eigenvalues of \mathbf{G} , (4) implies

$$\Sigma^\top\mathbf{U}^\top\mathbf{Q}\mathbf{Q}^\top\mathbf{U}\Sigma = \begin{pmatrix} \mathbf{0}_{r \times r} & \mathbf{0}_{r \times (m-r)} & \mathbf{0}_{r \times (n-m)} \\ \mathbf{0}_{(m-r) \times r} & \text{diag}(\sigma_{r+1}, \dots, \sigma_m) & \mathbf{0}_{(m-r) \times (n-m)} \\ \mathbf{0}_{(n-m) \times r} & \mathbf{0}_{(n-m) \times (m-r)} & \mathbf{0}_{(n-m) \times (n-m)} \end{pmatrix}. \quad (5)$$

Applying (5) to (3) yields

$$\|\mathbf{P}\mathbf{P}^\top\mathbf{G} - \mathbf{G}\|_F^2 = \text{tr}(\Sigma^\top\mathbf{U}^\top\mathbf{Q}\mathbf{Q}^\top\mathbf{U}\Sigma) = \sum_{i=r+1}^m \sigma_i^2 \leq \frac{m-r}{m} \|\mathbf{G}\|_F^2,$$

where the inequality uses $\|\mathbf{G}\|_F^2 = \text{tr}(\mathbf{G}^\top\mathbf{G}) = \text{tr}(\Sigma^\top\Sigma) = \sum_{i=1}^m \sigma_i^2$. \square

Lemma 2 (Gradient connections). *It holds for any $t, \tau > 0$ that*

$$\|\nabla_\ell f(\mathbf{x}^{(0)})\|_F^2 \leq \frac{2}{\tau} \sum_{r=0}^{\tau-1} \|\nabla_\ell f(\mathbf{x}^{(t+r)})\|_F^2 + (\tau-1) \sum_{r=0}^{\tau-2} \|\nabla_\ell f(\mathbf{x}^{(t+r+1)}) - \nabla_\ell f(\mathbf{x}^{(t+r)})\|_F^2. \quad (6)$$

Proof. For any $r = 1, \dots, \tau-1$, it holds that

$$\begin{aligned} \|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2 &= \|\nabla_\ell f(\mathbf{x}^{(t+r)}) - (\nabla_\ell f(\mathbf{x}^{(t+r)}) - \nabla_\ell f(\mathbf{x}^{(t)}))\|_F^2 \\ &\leq 2\|\nabla_\ell f(\mathbf{x}^{(t+r)})\|_F^2 + 2\|\nabla_\ell f(\mathbf{x}^{(t+r)}) - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2. \end{aligned} \quad (7)$$

For any $r = 2, \dots, \tau-1$, it holds that

$$\begin{aligned} \|\nabla_\ell f(\mathbf{x}^{(t+r)}) - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2 &= \left\| \sum_{i=1}^r \nabla_\ell f(\mathbf{x}^{(t+i)}) - \nabla_\ell f(\mathbf{x}^{(t+i-1)}) \right\|_F^2 \\ &\leq r \sum_{i=1}^r \|\nabla_\ell f(\mathbf{x}^{(t+i)}) - \nabla_\ell f(\mathbf{x}^{(t+i-1)})\|_F^2, \end{aligned} \quad (8)$$

where the inequality uses Cauchy's inequality. Summing (7) from $r = 1$ to $\tau - 1$ and applying (8) yields

$$\begin{aligned} \tau \|\nabla_{\ell} f(\mathbf{x}^{(t)})\|_F^2 &\leq 2 \sum_{r=0}^{\tau-1} \|\nabla_{\ell} f(\mathbf{x}^{(t+r)})\|_F^2 + 2 \sum_{i=1}^{\tau-1} \sum_{j=1}^i i \|\nabla_{\ell} f(\mathbf{x}^{(t+j)}) - \nabla_{\ell} f(\mathbf{x}^{(t+j-1)})\|_F^2 \\ &\leq 2 \sum_{r=0}^{\tau-1} \|\nabla_{\ell} f(\mathbf{x}^{(t+r)})\|_F^2 + 2 \sum_{j=1}^{\tau-1} \sum_{i=1}^{\tau-1} i \|\nabla_{\ell} f(\mathbf{x}^{(t+j)}) - \nabla_{\ell} f(\mathbf{x}^{(t+j-1)})\|_F^2 \\ &= 2 \sum_{r=0}^{\tau-1} \|\nabla_{\ell} f(\mathbf{x}^{(t+r)})\|_F^2 + \tau(\tau-1) \sum_{j=1}^{\tau-1} \|\nabla_{\ell} f(\mathbf{x}^{(t+j)}) - \nabla_{\ell} f(\mathbf{x}^{(t+j-1)})\|_F^2, \end{aligned}$$

which is exactly (6). \square

Lemma 3 (Projection orthogonality). *If $P \in \text{St}_{m,r}$, it holds for any $A, B \in \mathbb{R}^{m \times n}$ that*

$$\|PP^{\top}A + (I - PP^{\top})B\|_F^2 = \|PP^{\top}A\|_F^2 + \|(I - PP^{\top})B\|_F^2. \quad (9)$$

Proof. By definition we have $P^{\top}P = I$. It suffices to note that

$$\langle PP^{\top}A, (I - PP^{\top})B \rangle_F = \text{tr}(A^{\top}PP^{\top}(I - PP^{\top})B) = \text{tr}(0) = 0.$$

Lemma 4 (Descent lemma). *Under Assumption 2, for update*

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \eta \tilde{\mathbf{m}}^{(t)},$$

it holds that

$$\begin{aligned} f(\mathbf{x}^{(t+1)}) &\leq f(\mathbf{x}^{(t)}) - \left(\frac{1}{2\eta} - \frac{L}{2}\right) \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_2^2 + \frac{\eta}{2} \|\tilde{\mathbf{m}}^{(t)} - \nabla f(\mathbf{x}^{(t)})\|_2^2 \\ &\quad - \frac{\eta}{2} \|\nabla f(\mathbf{x}^{(t)})\|_2^2. \end{aligned} \quad (10)$$

Proof. By L -smoothness of f (Assumption 2) we have

$$\begin{aligned} &f(\mathbf{x}^{(t+1)}) - f(\mathbf{x}^{(t)}) \\ &\leq \langle \nabla f(\mathbf{x}^{(t)}), \mathbf{x}^{(t+1)} - \mathbf{x}^{(t)} \rangle + \frac{L}{2} \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_2^2 \\ &= \left\langle \frac{\tilde{\mathbf{m}}^{(t)}}{2}, \mathbf{x}^{(t+1)} - \mathbf{x}^{(t)} \right\rangle + \left\langle \nabla f(\mathbf{x}^{(t)}) - \frac{\tilde{\mathbf{m}}^{(t)}}{2}, \mathbf{x}^{(t+1)} - \mathbf{x}^{(t)} \right\rangle + \frac{L}{2} \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_2^2 \\ &= - \left(\frac{1}{2\eta} - \frac{L}{2}\right) \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_2^2 + \frac{\eta}{2} \|\nabla f(\mathbf{x}^{(t)}) - \tilde{\mathbf{m}}^{(t)}\|_2^2 - \frac{\eta}{2} \|\nabla f(\mathbf{x}^{(t)})\|_2^2, \end{aligned}$$

which is exactly (10). \square

Lemma 5 (Error of GoLore's projection). *Let $P \sim \mathcal{U}(\text{St}_{m,r})$, $Q \sim \mathcal{U}(\text{St}_{n,r})$, it holds for all $G \in \mathbb{R}^{m \times n}$ that*

$$\mathbb{E}[PP^{\top}] = \frac{r}{m} \cdot I, \quad \mathbb{E}[QQ^{\top}] = \frac{r}{n} \cdot I, \quad (11)$$

and

$$\mathbb{E}[\|PP^{\top}G - G\|_F^2] = \left(1 - \frac{r}{m}\right) \|G\|_F^2, \quad \mathbb{E}[\|GQQ^{\top} - G\|_F^2] = \left(1 - \frac{r}{n}\right) \|G\|_F^2. \quad (12)$$

Proof. We refer the proof of (11) to Theorem 2.2.2 in Chikuse (2012). By $P^{\top}P = I$, we have

$$\begin{aligned} \mathbb{E}[\|PP^{\top}G - G\|_F^2] &= \mathbb{E}[\text{tr}(G^{\top}(I - PP^{\top})^2G)] \\ &= \mathbb{E}[\text{tr}(G^{\top}(I - PP^{\top})G)] \\ &= \text{tr}(G^{\top}(I - \mathbb{E}[PP^{\top}])G). \end{aligned} \quad (13)$$

864 Applying (11) to (13) yields

$$\begin{aligned}
865 \mathbb{E}[\|\mathbf{P}\mathbf{P}^\top \mathbf{G} - \mathbf{G}\|_F^2] &= \text{tr} \left(\mathbf{G}^\top \left(\mathbf{I} - \frac{r}{m} \mathbf{I} \right) \mathbf{G} \right) \\
866 &= \left(1 - \frac{r}{m} \right) \text{tr}(\mathbf{G}^\top \mathbf{G}) \\
867 &= \left(1 - \frac{r}{m} \right) \|\mathbf{G}\|_F^2.
\end{aligned}$$

871 The other part of (12) can be proved similarly. \square

873 B.2 NON-CONVERGENCE OF GALORE

875 In this subsection, we present the proof for Theorem 1. We first restate Theorem 1 as follows:

876 **Theorem 5** (Non-convergence of GaLore). *There exists an objective function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ satisfying*
877 *Assumptions 1, 2, a stochastic gradient oracle (F, \mathcal{D}) satisfying Assumption 3, an initial point $\mathbf{x}^{(0)} \in$*
878 *\mathbb{R}^d , a constant $\epsilon_0 > 0$ such that for GaLore with any rank $r_\ell < \min\{m_\ell, n_\ell\}$, subspace changing*
879 *frequency τ , any subspace optimizer ρ with arbitrary hyperparameters and any $t > 0$, it holds that*

$$880 \|\nabla f(\mathbf{x}^{(t)})\|_2^2 \geq \epsilon_0.$$

882 *Proof.* Consider target function $f(\mathbf{X}) = \frac{L}{2} \text{tr}(\mathbf{X}^\top \mathbf{p}\mathbf{p}^\top \mathbf{X})$ where $L > 0$, $\mathbf{X} \in \mathbb{R}^{n \times n}$ with $n > 1$
883 and $\mathbf{p} = (1, 0, \dots, 0)^\top \in \mathbb{R}^n$. It holds that

$$884 f(\mathbf{X}) = \frac{L}{2} \|\mathbf{p}^\top \mathbf{X}\|_2^2 \geq 0,$$

885 thus f satisfies Assumption 1. Since $\nabla f(\mathbf{X}) = L\mathbf{p}\mathbf{p}^\top \mathbf{X}$, it holds that

$$886 \|\nabla f(\mathbf{X}) - \nabla f(\mathbf{Y})\|_F = L\|\mathbf{p}\mathbf{p}^\top (\mathbf{X} - \mathbf{Y})\|_F \leq L\|\mathbf{p}\mathbf{p}^\top\|_2 \|\mathbf{X} - \mathbf{Y}\|_F = L\|\mathbf{X} - \mathbf{Y}\|_F,$$

887 thus f satisfies Assumption 2.

888 Consider the following stochastic gradient oracle:

$$889 F(\mathbf{X}; \xi) = f(\mathbf{X}) + \xi \tilde{\sigma} \cdot \text{tr}(\mathbf{Q}\mathbf{Q}^\top \mathbf{X}), \quad \text{and} \quad \mathbb{P}_{\xi \sim \mathcal{D}}[\xi = 1] = \mathbb{P}_{\xi \sim \mathcal{D}}[\xi = -1] = 0.5,$$

890 where $\tilde{\sigma} = \sigma / \sqrt{(n-1)n/2}$ and

$$891 \mathbf{Q} = \begin{pmatrix} 0 \\ \text{diag}(1, \sqrt[4]{2}, \dots, \sqrt[4]{n-1}) \end{pmatrix} \in \mathbb{R}^{n \times (n-1)}.$$

892 Note that $\nabla F(\mathbf{X}; \xi) = \nabla f(\mathbf{X}) + \xi \tilde{\sigma} \mathbf{Q}\mathbf{Q}^\top$, it holds for any $\mathbf{X} \in \mathbb{R}^{n \times n}$ that

$$893 \mathbb{E}_{\xi \sim \mathcal{D}}[\nabla F(\mathbf{X}; \xi)] = \nabla f(\mathbf{X})$$

$$894 \mathbb{E}_{\xi \sim \mathcal{D}}[\|\nabla F(\mathbf{X}; \xi) - \nabla f(\mathbf{X})\|_F^2] = \tilde{\sigma}^2 \|\mathbf{Q}\mathbf{Q}^\top\|_F^2 = \frac{\sigma^2}{(n-1)n/2} \cdot \sum_{i=1}^{n-1} i = \sigma^2,$$

895 thus oracle (F, \mathcal{D}) satisfies Assumption 3.

896 Consider the following initial point:

$$897 \mathbf{X}^{(0)} = \begin{pmatrix} \lambda \mathbf{p}^\top \\ \mathbf{\Lambda} \end{pmatrix},$$

898 where $0 < \lambda < \tilde{\sigma}/L$ is a scalar and $\mathbf{\Lambda} \in \mathbb{R}^{(n-1) \times n}$ is an arbitrary matrix. We show that GaLore
899 with the above objective function f , stochastic gradient oracle (F, \mathcal{D}) , initial point $\mathbf{X}^{(0)}$, arbitrary
900 rank $0 < r < n$, arbitrary subspace changing frequency τ and arbitrary subspace optimizer ρ , can
901 only output points $\mathbf{X}^{(t)}$ with $\|\nabla f(\mathbf{X}^{(t)})\|_F^2 \geq \epsilon_0$ for $\epsilon_0 = L^2 \lambda^2 > 0$.

902 When $\tau \mid t$, GaLore recomputes the subspace projection matrix at iteration t . If the first row of $\mathbf{X}^{(t)}$
903 equals $\lambda \mathbf{p}^\top$, i.e., $\mathbf{X}^{(t)}[1, :] = \lambda \mathbf{p}^\top$, the stochastic gradient is given by

$$904 \mathbf{G}^{(t)} = L\mathbf{p}\mathbf{p}^\top \mathbf{X} + \xi^{(t)} \tilde{\sigma} \mathbf{Q}\mathbf{Q}^\top = \text{diag} \left(L\lambda, \xi^{(t)} \tilde{\sigma}, \sqrt{2} \xi^{(t)} \tilde{\sigma}, \dots, \sqrt{n-1} \xi^{(t)} \tilde{\sigma} \right).$$

918 since $L\lambda < \tilde{\sigma}$, computing SVD yields
 919

$$\begin{aligned}
 \mathbf{G}^{(t)} &= \begin{pmatrix} L\lambda & 0 & \cdots & 0 \\ 0 & \xi^{(t)}\tilde{\sigma} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sqrt{n-1}\xi^{(t)}\tilde{\sigma} \end{pmatrix} \\
 &= \underbrace{\begin{pmatrix} 0 & \cdots & 0 & \zeta_1 \\ 0 & \cdots & \zeta_2 & 0 \\ \vdots & \ddots & \vdots & \vdots \\ \zeta_n & \cdots & 0 & 0 \end{pmatrix}}_{:=\mathbf{U}} \underbrace{\begin{pmatrix} \sqrt{n-1}\tilde{\sigma} & \cdots & 0 & 0 \\ \vdots & \ddots & \vdots & \vdots \\ 0 & \cdots & \tilde{\sigma} & 0 \\ 0 & \cdots & 0 & L\lambda \end{pmatrix}}_{:=\mathbf{\Sigma}} \underbrace{\begin{pmatrix} 0 & 0 & \cdots & \zeta_n\xi^{(t)} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \zeta_2\xi^{(t)} & \cdots & 0 \\ \zeta_1 & 0 & \cdots & 0 \end{pmatrix}}_{:=\mathbf{V}^\top},
 \end{aligned}$$

931 where $\zeta_1, \dots, \zeta_n \in \{-1, 1\}$. For any rank $r < n$, the projection matrix is thus
 932

$$\mathbf{P}^{(t)} = \begin{pmatrix} 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 \\ 0 & 0 & \cdots & \zeta_{n-r+1} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \zeta_{n-1} & \cdots & 0 \\ \zeta_n & 0 & \cdots & 0 \end{pmatrix} \in \mathbb{R}^{n \times r}.$$

942 Using this projection matrix, the subspace updates in the following τ iterations is as
 943

$$\mathbf{X}^{(t+\Delta_t)} = \mathbf{X}^{(t)} + \mathbf{P}^{(t)} \sum_{s=0}^{\Delta_t-1} \rho^{(t+s)} ((\mathbf{P}^{(t)})^\top \mathbf{G}^{(t)}) \Rightarrow \mathbf{X}^{(t+\Delta_t)}[1, :] = \mathbf{X}^{(t)}[1, :] = \lambda \mathbf{p}^\top,$$

947 for $\Delta_t = 1, 2, \dots, \tau$. Since $\mathbf{X}^{(0)}[1, :] = \lambda \mathbf{p}^\top$, it holds for all $t > 0$ that $\mathbf{X}^{(t)}[1, :] = \lambda \mathbf{p}^\top$ and thus
 948

$$\|\nabla f(\mathbf{X}^{(t)})\|_F^2 = L^2 \lambda^2 = \epsilon_0.$$

951 \square

953 **Remark.** When setting $\mathbf{B} = 0$ in the quadratic problem setting (Sec. 7), the quadratic problem is
 954 equivalent to the counter-example we construct in the proof of Theorem 5. The illustration in Fig. 5
 955 displays the loss curves for this problem.
 956

957 B.3 CONVERGENCE OF DETERMINISTIC GALORE

958 In this subsection, we present the proof for Theorem 2. GaLore using deterministic gradients and
 959 MSGD with MP is specified as Alg. 2.

962 **Lemma 6** (Momentum contraction). *In deterministic GaLore using MSGD with MP (Alg. 2), if*
 963 *$0 < \beta_1 \leq 1$, term $\tilde{\mathbf{M}}_\ell^{(t)}$ has the following contraction properties:*
 964

- 965 • When $t = 0$, it holds that

$$\begin{aligned}
 \|\tilde{\mathbf{M}}_\ell^{(0)} - \nabla_\ell f(\mathbf{X}^{(0)})\|_F^2 &\leq (\tau - 1)(1 - \delta_\ell \beta_1) \sum_{r=0}^{\tau-2} \|\nabla_\ell f(\mathbf{x}^{(r+1)}) - \nabla_\ell f(\mathbf{x}^{(r)})\|_F^2 \\
 &\quad + \frac{2(1 - \delta_\ell \beta_1)}{\tau} \sum_{r=0}^{\tau-1} \|\nabla_\ell f(\mathbf{x}^{(r)})\|_F^2; \tag{14}
 \end{aligned}$$

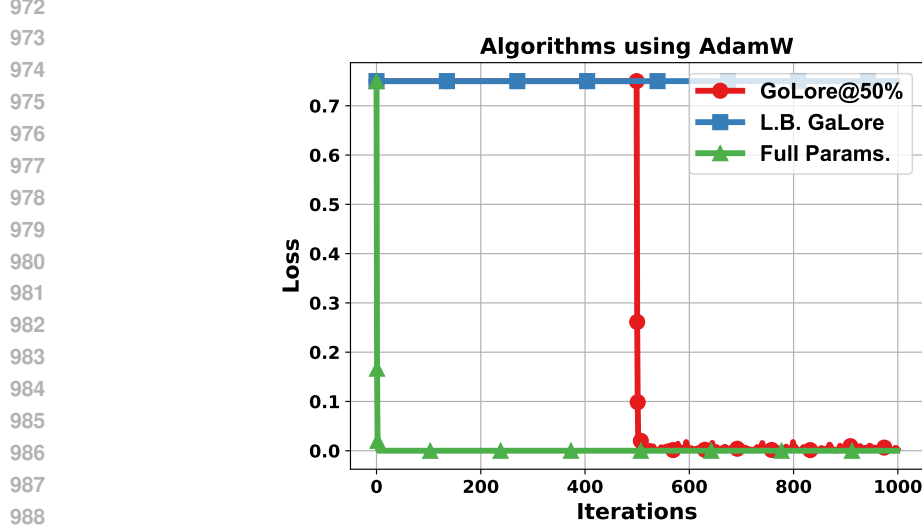


Figure 5: Loss curves of algorithms using AdamW. *GoLore@50%* uses GaLore in the first half and shifts to GoLore in the last half, *Full Params.* denotes full-parameter training.

- When $t = k\tau$, $k \in \mathbb{N}^*$, it holds that

$$\begin{aligned}
& \|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2 - \left(1 - \left(1 - \frac{\delta_\ell}{4}\right)\beta_1\right) \|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2 \\
& \leq \frac{2(1 - \delta_\ell)}{\tau} \sum_{r=0}^{\tau-1} \|\nabla_\ell f(\mathbf{x}^{(k\tau+r)})\|_F^2 + \frac{5(1 - \beta_1)}{\delta_\ell \beta_1} \|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2 \\
& \quad + (\tau - 1)(1 - \delta_\ell) \sum_{r=0}^{\tau-2} \|\nabla_\ell f(\mathbf{x}^{(k\tau+r+1)}) - \nabla_\ell f(\mathbf{x}^{(k\tau+r)})\|_F^2; \tag{15}
\end{aligned}$$

- When $t = k\tau + r$, $k \in \mathbb{N}$, $1 \leq r < \tau$, it holds that

$$\begin{aligned}
& \|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2 - \left(1 - \left(1 - \frac{\delta_\ell}{4}\right)\beta_1\right) \|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2 \\
& \leq \left(1 - \frac{\delta_\ell}{2}\right)\beta_1 \|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2 + \frac{5(1 - \beta_1)}{\delta_\ell \beta_1} \|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2 \\
& \quad + \frac{10r\beta_1}{\delta_\ell} \sum_{i=1}^r \|\nabla_\ell f(\mathbf{x}^{(k\tau+i)}) - \nabla_\ell f(\mathbf{x}^{(k\tau+i-1)})\|_F^2. \tag{16}
\end{aligned}$$

Proof. Without loss of generality assume $m_\ell \leq n_\ell$ (the other case can be proved similarly). When $t = 0$, we have

$$\begin{aligned}
\|\tilde{\mathbf{M}}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2 &= \|\beta_1(\mathbf{P}_\ell^{(0)}(\mathbf{P}_\ell^{(0)})^\top - \mathbf{I})\nabla_\ell f(\mathbf{x}^{(0)}) - (1 - \beta_1)\nabla_\ell f(\mathbf{x}^{(0)})\|_F^2 \\
&\leq \beta_1(1 - \delta_\ell)\|\nabla_\ell f(\mathbf{x}^{(0)})\|_F^2 + (1 - \beta_1)\|\nabla_\ell f(\mathbf{x}^{(0)})\|_F^2 \\
&= (1 - \delta_\ell\beta_1)\|\nabla_\ell f(\mathbf{x}^{(0)})\|_F^2, \tag{17}
\end{aligned}$$

where the inequality uses Lemma 1 and Jensen's inequality. Applying Lemma 2 to (17) yields (14).

When $t = k\tau$, $k \in \mathbb{N}^*$, we have

$$\begin{aligned}
& \|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2 \\
&= \|\mathbf{P}_\ell^{(t)}(\mathbf{P}_\ell^{(t)})^\top [(1 - \beta_1)\tilde{\mathbf{M}}_\ell^{(t-1)} + \beta_1\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})] - (\mathbf{I} - \mathbf{P}_\ell^{(t)}(\mathbf{P}_\ell^{(t)})^\top)\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2 \\
&= \|\mathbf{P}_\ell^{(t)}(\mathbf{P}_\ell^{(t)})^\top [(1 - \beta_1)(\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)}))]\|_F^2 + \|(\mathbf{I} - \mathbf{P}_\ell^{(t)}(\mathbf{P}_\ell^{(t)})^\top)\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2 \\
&\leq \|(1 - \beta_1)(\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)}))\|_F^2 + (1 - \delta_\ell)\|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2, \tag{18}
\end{aligned}$$

Algorithm 2 GaLore using deterministic gradients and MSGD with MP

Input: Initial point $\mathbf{x}^{(0)}$, learning rate η , subspace changing frequency τ , rank $\{r_\ell\}_{\ell=1}^{N_L}$, momentum parameter β_1 .

Output: $\{\mathbf{x}^{(t)}\}_{t=0}^T$.

Initialize optimizer state $\{M_\ell^{(-1)}\}_{\ell=1}^{N_L}$ to zero;

for $t = 0, 1, \dots, T - 1$ **do**

for $\ell = 1, 2, \dots, N_L$ **do**

$\mathbf{G}_\ell^{(t)} \leftarrow \nabla_\ell f(\mathbf{x}^{(t)});$

if $t \equiv 0 \pmod{\tau}$ **then**

$\mathbf{U}, \Sigma, \mathbf{V} \leftarrow \text{SVD}(\mathbf{G}_\ell^{(t)});$

if $m_\ell \leq n_\ell$ **then**

$\mathbf{P}_\ell^{(t)} \leftarrow \mathbf{U}[:, :r_\ell];$

$\mathbf{M}_\ell^{(t)} \leftarrow (1 - \beta_1)(\mathbf{P}_\ell^{(t)})^\top \mathbf{P}_\ell^{(t-1)} \mathbf{M}_\ell^{(t-1)} + \beta_1(\mathbf{P}_\ell^{(t)})^\top \mathbf{G}_\ell^{(t)};$

$\mathbf{X}_\ell^{(t+1)} \leftarrow \mathbf{X}_\ell^{(t)} - \eta \mathbf{P}_\ell^{(t)} \mathbf{M}_\ell^{(t)};$

else

$\mathbf{Q}_\ell^{(t)} \leftarrow \mathbf{V}[:, :r_\ell];$

$\mathbf{M}_\ell^{(t)} \leftarrow (1 - \beta_1)\mathbf{M}_\ell^{(t-1)}(\mathbf{Q}_\ell^{(t-1)})^\top \mathbf{Q}_\ell^{(t)} + \beta_1\mathbf{G}_\ell^{(t)}\mathbf{Q}_\ell^{(t)};$

$\mathbf{X}_\ell^{(t+1)} \leftarrow \mathbf{X}_\ell^{(t)} - \eta \mathbf{M}_\ell^{(t)}(\mathbf{Q}_\ell^{(t)})^\top;$

end if

else

if $m_\ell \leq n_\ell$ **then**

$\mathbf{P}_\ell^{(t)} \leftarrow \mathbf{P}_\ell^{(t-1)};$

$\mathbf{M}_\ell^{(t)} \leftarrow (1 - \beta_1)\mathbf{M}_\ell^{(t-1)} + \beta_1(\mathbf{P}_\ell^{(t)})^\top \mathbf{G}_\ell^{(t)};$

$\mathbf{X}_\ell^{(t+1)} \leftarrow \mathbf{X}_\ell^{(t)} - \eta \mathbf{P}_\ell^{(t)} \mathbf{M}_\ell^{(t)};$

else

$\mathbf{Q}_\ell^{(t)} \leftarrow \mathbf{Q}_\ell^{(t-1)};$

$\mathbf{M}_\ell^{(t)} \leftarrow (1 - \beta_1)\mathbf{M}_\ell^{(t-1)} + \beta_1\mathbf{G}_\ell^{(t)}\mathbf{Q}_\ell^{(t)};$

$\mathbf{X}_\ell^{(t+1)} \leftarrow \mathbf{X}_\ell^{(t)} - \eta \mathbf{M}_\ell^{(t)}(\mathbf{Q}_\ell^{(t)})^\top;$

end if

end if

end for

end for

where the second equality uses Lemma 3 and $\mathbf{G}_\ell^{(t)} = \nabla_\ell f(\mathbf{x}^{(t)})$, the inequality uses Lemma 1 and $\|\mathbf{P}_\ell^{(t)}(\mathbf{P}_\ell^{(t)})^\top\|_2 = 1$. By Young's inequality, we have

$$\begin{aligned} & \|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2 \\ &= \|(\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t-1)})) - (\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(t-1)}))\|_F^2 \\ &\leq \left(1 + \frac{\delta_\ell \beta_1}{4}\right) \|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2 + \left(1 + \frac{4}{\delta_\ell \beta_1}\right) \|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2. \end{aligned} \quad (19)$$

Applying Lemma 2 and (19) to (18) yields (15).

When $t = k\tau + r$, $k \in \mathbb{N}$, $1 \leq r < \tau$, we have

$$\begin{aligned} & \|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2 \\ &= \|(1 - \beta_1)(\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)})) + \beta_1(\mathbf{P}_\ell^{(t)}(\mathbf{P}_\ell^{(t)})^\top - \mathbf{I})\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2 \\ &\leq (1 - \beta_1)\|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2 + \beta_1\|(\mathbf{I} - \mathbf{P}_\ell^{(k\tau)}(\mathbf{P}_\ell^{(k\tau)})^\top)\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2, \end{aligned} \quad (20)$$

where the inequality uses Jensen's inequality and $\mathbf{P}_\ell^{(t)} = \mathbf{P}_\ell^{(t-1)} = \dots = \mathbf{P}_\ell^{(k\tau)}$. The first term can be similarly upper bounded as (19). For the second term, we have

$$\begin{aligned}
& \|(\mathbf{I} - \mathbf{P}_\ell^{(k\tau)})(\mathbf{P}_\ell^{(k\tau)})^\top \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2 \\
& \leq \left(1 + \frac{\delta_\ell}{4}\right) \|(\mathbf{I} - \mathbf{P}_\ell^{(k\tau)})(\mathbf{P}_\ell^{(k\tau)})^\top \nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2 \\
& \quad + \left(1 + \frac{4}{\delta_\ell}\right) \|(\mathbf{I} - \mathbf{P}_\ell^{(k\tau)})(\mathbf{P}_\ell^{(k\tau)})^\top (\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(k\tau)}))\|_F^2 \\
& \leq \left(1 + \frac{\delta_\ell}{4}\right) (1 - \delta_\ell) \|\nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2 + \frac{5}{\delta_\ell} \|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2, \tag{21}
\end{aligned}$$

where the first inequality uses Young's inequality and the second inequality uses Lemma 1. By Young's inequality, we have

$$\|\nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2 \leq \left(1 + \frac{\delta_\ell}{4}\right) \|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2 + \left(1 + \frac{4}{\delta_\ell}\right) \|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2. \tag{22}$$

Note that $t = k\tau + r$, we further have

$$\begin{aligned}
\|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2 &= \left\| \sum_{i=1}^r \nabla_\ell f(\mathbf{x}^{(k\tau+i)}) - \nabla_\ell f(\mathbf{x}^{(k\tau+i-1)}) \right\|_F^2 \\
&\leq r \sum_{i=1}^r \|\nabla_\ell f(\mathbf{x}^{(k\tau+i)}) - \nabla_\ell f(\mathbf{x}^{(k\tau+i-1)})\|_F^2, \tag{23}
\end{aligned}$$

where the inequality uses Cauchy's inequality. Applying (22)(23) to (21) yields

$$\begin{aligned}
& \|(\mathbf{I} - \mathbf{P}_\ell^{(k\tau)})(\mathbf{P}_\ell^{(k\tau)})^\top \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2 \\
& \leq \left(1 - \frac{\delta_\ell}{2}\right) \|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2 + \frac{10r}{\delta_\ell} \sum_{i=1}^r \|\nabla_\ell f(\mathbf{x}^{(k\tau+i)}) - \nabla_\ell f(\mathbf{x}^{(k\tau+i-1)})\|_F^2. \tag{24}
\end{aligned}$$

Applying (19)(24) to (20) yields (16). \square

Lemma 7 (Momentum error). *Under Assumption 2, if $0 < \beta_1 \leq 1$ in deterministic GaLore using MSGD and MP (Alg. 2), it holds for any $K \geq 1$ that*

$$\begin{aligned}
& \sum_{t=0}^{K\tau-1} \|\tilde{\mathbf{m}}^{(t)} - \nabla f(\mathbf{x}^{(t)})\|_2^2 \\
& \leq \left(\frac{5(1-\beta_1)}{(1-\underline{\delta}/4)\underline{\delta}\beta_1^2} + \frac{5\tau(\tau-1)}{(1-\underline{\delta}/4)\underline{\delta}} + \frac{\tau-1}{(1-\underline{\delta}/4)\beta_1} \right) L^2 \sum_{t=0}^{K\tau-2} \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_2^2 \\
& \quad + \left(\frac{1-\underline{\delta}/2}{1-\underline{\delta}/4} + \frac{2}{(1-\underline{\delta}/4)\tau\beta_1} \right) \sum_{t=0}^{K\tau-1} \|\nabla f(\mathbf{x}^{(t)})\|_2^2. \tag{25}
\end{aligned}$$

Proof. By Lemma 6 we have

$$\begin{aligned}
& \sum_{t=0}^{K\tau-1} \|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2 - \left(1 - \left(1 - \frac{\delta_\ell}{4}\right)\beta_1\right) \sum_{t=0}^{K\tau-2} \|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2 \\
& \leq \left(\frac{5(1-\beta_1)}{\delta_\ell\beta_1} + \frac{5\tau(\tau-1)\beta_1}{\delta_\ell} + (\tau-1) \right) \sum_{t=0}^{K\tau-2} \|\nabla_\ell f(\mathbf{x}^{(t+1)}) - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2 \\
& \quad + \left(\frac{2}{\tau} + \left(1 - \frac{\delta_\ell}{2}\right)\beta_1 \right) \sum_{t=0}^{K\tau-1} \|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2,
\end{aligned}$$

1134 which implies

$$\begin{aligned}
1135 & \sum_{t=0}^{K\tau-1} \|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2 \\
1136 & \leq \left(\frac{5(1-\beta_1)}{(1-\delta_\ell/4)\delta_\ell\beta_1^2} + \frac{5\tau(\tau-1)}{(1-\delta_\ell/4)\delta_\ell} + \frac{\tau-1}{(1-\delta_\ell/4)\beta_1} \right) \sum_{t=0}^{K\tau-2} \|\nabla_\ell f(\mathbf{x}^{(t+1)}) - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2 \\
1137 & \quad + \left(\frac{1-\delta_\ell/2}{1-\delta_\ell/4} + \frac{2}{(1-\delta_\ell/4)\tau\beta_1} \right) \sum_{t=0}^{K\tau-1} \|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2. \tag{26}
\end{aligned}$$

1144 Summing (26) for $\ell = 1, \dots, N_L$ and applying Assumption 2 yields (25). \square

1146 Now we are ready to prove the convergence of Alg. 2.

1147 **Theorem 6** (Convergence of deterministic GaLore). *Under Assumptions 1-2, if hyperparameters*

$$1148 \quad 0 < \beta_1 \leq 1, \quad \tau \geq \frac{64}{3\beta_1\delta}, \quad 0 < \eta \leq \min \left\{ \frac{1}{4L}, \sqrt{\frac{3\delta\beta_1^2}{80L^2}}, \sqrt{\frac{3\delta}{80\tau^2L^2}}, \sqrt{\frac{3\beta_1}{16\tau L^2}} \right\}, \tag{27}$$

1151 *GaLore using deterministic gradients and MSGD with MP (Alg. 2) converges as*

$$1152 \quad \frac{1}{K\tau} \sum_{t=0}^{K\tau-1} \|\nabla f(\mathbf{x}^{(t)})\|_2^2 \leq \frac{16\Delta}{\delta\eta K\tau} \tag{28}$$

1155 *for any $K \geq 1$, where $\Delta = f(\mathbf{x}^{(0)}) - \inf_{\mathbf{x}} f(\mathbf{x})$.*

1157 *Proof.* By Lemma 4 we have

$$\begin{aligned}
1158 \quad \sum_{t=0}^{K\tau-1} \|\nabla f(\mathbf{x}^{(t)})\|_2^2 & \leq \frac{2[f(\mathbf{x}^{(0)}) - f(\mathbf{x}^{(K\tau)})]}{\eta} + \sum_{t=0}^{K\tau-1} \|\tilde{\mathbf{m}}^{(t)} - \nabla f(\mathbf{x}^{(t)})\|_2^2 \\
1159 & \quad - \left(\frac{1}{\eta^2} - \frac{L}{\eta} \right) \sum_{t=0}^{K\tau-1} \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_2^2. \tag{29}
\end{aligned}$$

1164 Applying Lemma 7 to (29) and using $\delta \leq \bar{\delta} < 1$ yields

$$\begin{aligned}
1165 \quad & \left(\frac{\delta}{4} - \frac{8}{3\tau\beta_1} \right) \sum_{t=0}^{K\tau-1} \|\nabla f(\mathbf{x}^{(t)})\|_2^2 \\
1166 & \leq \frac{2}{\eta} f(\mathbf{x}^{(0)}) - f(\mathbf{x}^{(K\tau)}) \\
1167 & \quad - \left(\frac{1}{\eta^2} - \frac{L}{\eta} - \frac{20(1-\beta_1)L^2}{3\delta\beta_1^2} - \frac{20\tau(\tau-1)L^2}{3\delta} - \frac{4(\tau-1)L^2}{3\beta_1} \right) \sum_{t=0}^{K\tau-1} \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_2^2. \tag{30}
\end{aligned}$$

1173 By (27) we have

$$1174 \quad \frac{\delta}{4} - \frac{8}{3\tau\beta_1} \geq \frac{\delta}{8}, \quad \text{and} \quad \frac{1}{4\eta^2} \geq \max \left\{ \frac{L}{\eta}, \frac{20(1-\beta_1)L^2}{3\delta\beta_1^2}, \frac{20\tau(\tau-1)L^2}{3\delta}, \frac{4(\tau-1)L^2}{3\beta_1} \right\}. \tag{31}$$

1176 Applying (31) to (30) yields (28). \square

1178 We now prove Theorem 2, which is restated as follows.

1179 **Corollary 1** (Convergence complexity of deterministic GaLore). *Under Assumptions 1-2, if $T \geq 64/(3\delta)$ and we choose*

$$\begin{aligned}
1182 \quad & \beta_1 = 1 \\
1183 \quad & \tau = \left\lceil \frac{64}{3\delta\beta_1} \right\rceil \\
1184 & \eta = \left(4L + \sqrt{\frac{80L^2}{3\delta\beta_1^2}} + \sqrt{\frac{80\tau^2L^2}{3\delta}} + \sqrt{\frac{16\tau L^2}{3\beta_1}} \right)^{-1}, \tag{32}
\end{aligned}$$

1188 *GaLore using deterministic gradients and MSGD with MP (Alg. 2) converges as*
 1189

$$1190 \frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}^{(t)})\|_2^2 = \mathcal{O}\left(\frac{L\Delta}{\underline{\delta}^{5/2}T}\right), \quad (32)$$

1193 where $\Delta = f(\mathbf{x}^{(0)}) - \inf_{\mathbf{x}} f(\mathbf{x})$. Consequently, the computation complexity to reach an ε -accurate
 1194 solution \mathbf{x} such that $\|\nabla f(\mathbf{x})\|_2^2 \leq \varepsilon$ is $\mathcal{O}\left(\frac{L\Delta}{\underline{\delta}^{5/2}\varepsilon} + \frac{1}{\underline{\delta}}\right)$.
 1195
 1196
 1197

1198 *Proof.* $T \geq 1 + 64/(3\underline{\delta})$ guarantees $T \geq \tau$. Let $T = K\tau + r$, where $K \in \mathbb{N}^*$ and $0 \leq r < \tau$. If
 1199 $r = 0$, (32) is a direct result of Theorem 6. If $r > 0$, applying Theorem 6 to $\tilde{K} := K + 1$ yields
 1200

$$1201 \frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}^{(t)})\|_2^2 \leq \frac{\tilde{K}\tau}{T} \cdot \frac{1}{\tilde{K}\tau} \sum_{t=0}^{\tilde{K}\tau-1} \|\nabla f(\mathbf{x}^{(t)})\|_2^2 = \mathcal{O}\left(\frac{L\Delta}{\underline{\delta}^{5/2}T}\right).$$

1205 □

1207 B.4 CONVERGENCE OF LARGE-BATCH GALORE

1208
 1209 In this subsection, we present the proof for Theorem 3. GaLore using large-batch stochastic gradi-
 1210 ents and MSGD with MP is specified as Alg. 3.

1211 **Lemma 8** (Momentum contraction). *Under Assumption 3, in large-batch GaLore using MSGD with*
 1212 *MP (Alg. 3), if $0 < \beta_1 \leq 1$, term $\tilde{\mathbf{M}}_\ell^{(t)}$ has the following contraction properties:*
 1213

- 1214 • When $t = 0$, it holds that

$$1215 \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(0)} - \nabla_\ell f(\mathbf{X}^{(0)})\|_F^2] \leq 2(\tau - 1)(1 - \delta_\ell \beta_1) \sum_{r=0}^{\tau-2} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(r+1)}) - \nabla_\ell f(\mathbf{x}^{(r)})\|_F^2]$$

$$1216 + \frac{4(1 - \delta_\ell \beta_1)}{\tau} \sum_{r=0}^{\tau-1} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(r)})\|_F^2] + \frac{4\beta_1 \sigma_\ell^2}{\mathcal{B}}; \quad (33)$$

- 1222 • When $t = k\tau$, $k \in \mathbb{N}^*$, it holds that

$$1223 \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] - \left(1 - \left(1 - \frac{\delta_\ell}{4}\right)\beta_1\right) \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2]$$

$$1224 \leq \frac{4(1 - \delta_\ell)}{\tau} \sum_{r=0}^{\tau-1} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(k\tau+r)})\|_F^2] + \frac{5(1 - \beta_1)}{\delta_\ell \beta_1} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2]$$

$$1225 + 2(\tau - 1)(1 - \delta_\ell) \sum_{r=0}^{\tau-2} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(k\tau+r+1)}) - \nabla_\ell f(\mathbf{x}^{(k\tau+r)})\|_F^2] + \frac{5\sigma_\ell^2}{\mathcal{B}}; \quad (34)$$

- 1233 • When $t = k\tau + r$, $k \in \mathbb{N}$, $1 \leq r < \tau$, it holds that

$$1234 \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] - \left(1 - \left(1 - \frac{\delta_\ell}{4}\right)\beta_1\right) \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2]$$

$$1235 \leq \left(1 - \frac{\delta_\ell}{2}\right)\beta_1 \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + \frac{5(1 - \beta_1)}{\delta_\ell \beta_1} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2]$$

$$1236 + \frac{15r\beta_1}{\delta_\ell} \sum_{i=1}^r \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(k\tau+i)}) - \nabla_\ell f(\mathbf{x}^{(k\tau+i-1)})\|_F^2] + \left(\frac{11\beta_1}{\delta_\ell \mathcal{B}} + \beta_1^2\right) \sigma_\ell^2. \quad (35)$$

Algorithm 3 GaLore using large-batch stochastic gradients and MSGD with MP

Input: Initial point $\mathbf{x}^{(0)}$, data distribution \mathcal{D} , learning rate η , subspace changing frequency τ , rank $\{r_\ell\}_{\ell=1}^{N_L}$, momentum parameter β_1 , large batch size \mathcal{B} .

Output: $\{\mathbf{x}^{(t)}\}_{t=0}^T$.

Initialize optimizer state $\{\mathbf{M}_\ell^{(-1)}\}_{\ell=1}^{N_L}$ to zero;

for $t = 0, 1, \dots, T - 1$ **do**

if $t \equiv 0 \pmod{\tau}$ **then**

 Sample $\{\xi^{(t,b)}\}_{b=1}^{\mathcal{B}} \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}$;

else

 Sample $\xi^{(t)} \sim \mathcal{D}$;

end if

for $\ell = 1, 2, \dots, N_L$ **do**

if $t \equiv 0 \pmod{\tau}$ **then**

$\mathbf{G}_\ell^{(t)} = \frac{1}{\mathcal{B}} \sum_{b=1}^{\mathcal{B}} \nabla_\ell F(\mathbf{x}^{(t)}; \xi^{(t,b)})$;

$\mathbf{U}, \Sigma, \mathbf{V} \leftarrow \text{SVD}(\mathbf{G}_\ell^{(t)})$;

if $m_\ell \leq n_\ell$ **then**

$\mathbf{P}_\ell^{(t)} \leftarrow \mathbf{U}[:, :r_\ell]$;

$\mathbf{M}_\ell^{(t)} \leftarrow (1 - \beta_1)(\mathbf{P}_\ell^{(t)})^\top \mathbf{P}_\ell^{(t-1)} \mathbf{M}_\ell^{(t-1)} + \beta_1(\mathbf{P}_\ell^{(t)})^\top \mathbf{G}_\ell^{(t)}$;

$\mathbf{X}_\ell^{(t+1)} \leftarrow \mathbf{X}_\ell^{(t)} - \eta \mathbf{P}_\ell^{(t)} \mathbf{M}_\ell^{(t)}$;

else

$\mathbf{Q}_\ell^{(t)} \leftarrow \mathbf{V}[:, :r_\ell]$;

$\mathbf{M}_\ell^{(t)} \leftarrow (1 - \beta_1)\mathbf{M}_\ell^{(t-1)}(\mathbf{Q}_\ell^{(t-1)})^\top \mathbf{Q}_\ell^{(t)} + \beta_1 \mathbf{G}_\ell^{(t)} \mathbf{Q}_\ell^{(t)}$;

$\mathbf{X}_\ell^{(t+1)} \leftarrow \mathbf{X}_\ell^{(t)} - \eta \mathbf{M}_\ell^{(t)}(\mathbf{Q}_\ell^{(t)})^\top$;

end if

else

$\mathbf{G}_\ell^{(t)} = \nabla_\ell F(\mathbf{x}^{(t)}; \xi^{(t)})$;

if $m_\ell \leq n_\ell$ **then**

$\mathbf{P}_\ell^{(t)} \leftarrow \mathbf{P}_\ell^{(t-1)}$;

$\mathbf{M}_\ell^{(t)} \leftarrow (1 - \beta_1)\mathbf{M}_\ell^{(t-1)} + \beta_1(\mathbf{P}_\ell^{(t)})^\top \mathbf{G}_\ell^{(t)}$;

$\mathbf{X}_\ell^{(t+1)} \leftarrow \mathbf{X}_\ell^{(t)} - \eta \mathbf{P}_\ell^{(t)} \mathbf{M}_\ell^{(t)}$;

else

$\mathbf{Q}_\ell^{(t)} \leftarrow \mathbf{Q}_\ell^{(t-1)}$;

$\mathbf{M}_\ell^{(t)} \leftarrow (1 - \beta_1)\mathbf{M}_\ell^{(t-1)} + \beta_1 \mathbf{G}_\ell^{(t)} \mathbf{Q}_\ell^{(t)}$;

$\mathbf{X}_\ell^{(t+1)} \leftarrow \mathbf{X}_\ell^{(t)} - \eta \mathbf{M}_\ell^{(t)}(\mathbf{Q}_\ell^{(t)})^\top$;

end if

end if

end for

end for

Proof. Without loss of generality assume $m_\ell \leq n_\ell$ (the other case can be proved similarly). When $t = 0$, we have

$$\begin{aligned}
& \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2] \\
&= \mathbb{E}[\|\beta_1 \mathbf{P}_\ell^{(0)} (\mathbf{P}_\ell^{(0)})^\top \mathbf{G}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2] \\
&= \mathbb{E}[\|\beta_1 (\mathbf{P}_\ell^{(0)} (\mathbf{P}_\ell^{(0)})^\top - \mathbf{I}) \mathbf{G}_\ell^{(0)} + \beta_1 (\mathbf{G}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)})) - (1 - \beta_1) \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2] \\
&\leq \beta_1 \mathbb{E}[\|(\mathbf{P}_\ell^{(0)} (\mathbf{P}_\ell^{(0)})^\top - \mathbf{I}) \mathbf{G}_\ell^{(0)} + \mathbf{G}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2] + (1 - \beta_1) \|\nabla_\ell f(\mathbf{x}^{(0)})\|_F^2, \quad (36)
\end{aligned}$$

where the inequality uses Jensen's inequality. For the first term we have

$$\begin{aligned}
& \mathbb{E}[\|(\mathbf{P}_\ell^{(0)}(\mathbf{P}_\ell^{(0)})^\top - \mathbf{I})\mathbf{G}_\ell^{(0)} + \mathbf{G}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2] \\
& \leq 2\mathbb{E}[\|(\mathbf{I} - \mathbf{P}_\ell^{(0)}(\mathbf{P}_\ell^{(0)})^\top)\mathbf{G}_\ell^{(0)}\|_F^2] + 2\mathbb{E}[\|\mathbf{G}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2] \\
& \leq 2(1 - \delta_\ell)\mathbb{E}[\|\mathbf{G}_\ell^{(0)}\|_F^2] + 2\mathbb{E}[\|\mathbf{G}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2] \\
& \leq 2(1 - \delta_\ell)\|\nabla_\ell f(\mathbf{x}^{(0)})\|_F^2 + \frac{(4 - 2\delta_\ell)\sigma_\ell^2}{\mathcal{B}}, \tag{37}
\end{aligned}$$

where the first inequality uses Cauchy's inequality, the second inequality uses Lemma 1, the third inequality uses $\mathbb{E}[\|\mathbf{G}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2] \leq \sigma_\ell^2/\mathcal{B}$ (Assumption 3). Applying (37) and Lemma 2 to (36) yields (33).

When $t = k\tau$, $k \in \mathbb{N}^*$, we have

$$\begin{aligned}
& \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
& = \mathbb{E}[\|\mathbf{P}_\ell^{(t)}(\mathbf{P}_\ell^{(t)})^\top [(1 - \beta_1)\tilde{\mathbf{M}}_\ell^{(t-1)} + \beta_1\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})] - (\mathbf{I} - \mathbf{P}_\ell^{(t)}(\mathbf{P}_\ell^{(t)})^\top)\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
& = \mathbb{E}[\|\mathbf{P}_\ell^{(t)}(\mathbf{P}_\ell^{(t)})^\top [(1 - \beta_1)\tilde{\mathbf{M}}_\ell^{(t-1)} + \beta_1\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})]\|_F^2] \\
& \quad + \mathbb{E}[\|(\mathbf{I} - \mathbf{P}_\ell^{(t)}(\mathbf{P}_\ell^{(t)})^\top)\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2], \tag{38}
\end{aligned}$$

where the second equality uses Lemma 3. By $\|\mathbf{P}_\ell^{(t)}(\mathbf{P}_\ell^{(t)})^\top\|_2 = 1$, we have

$$\begin{aligned}
& \mathbb{E}[\|\mathbf{P}_\ell^{(t)}(\mathbf{P}_\ell^{(t)})^\top [(1 - \beta_1)\tilde{\mathbf{M}}_\ell^{(t-1)} + \beta_1\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})]\|_F^2] \\
& \leq \mathbb{E}[\|(1 - \beta_1)\tilde{\mathbf{M}}_\ell^{(t-1)} + \beta_1\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
& = \mathbb{E}[\|(1 - \beta_1)(\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)})) + \beta_1(\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)}))\|_F^2] \\
& \leq \mathbb{E}[\|(1 - \beta_1)(\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)}))\|_F^2] + \beta_1^2\mathbb{E}[\|\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2], \tag{39}
\end{aligned}$$

where the last inequality uses the unbiasedness of $\mathbf{G}_\ell^{(t)}$ (Assumption 3). By Young's inequality, we have

$$\begin{aligned}
& \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
& = \mathbb{E}[\|(\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t-1)})) - (\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(t-1)}))\|_F^2] \\
& \leq \left(1 + \frac{\delta_\ell\beta_1}{4}\right)\mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2] + \left(1 + \frac{4}{\delta_\ell\beta_1}\right)\mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2]. \tag{40}
\end{aligned}$$

Applying (40) to (39) yields

$$\begin{aligned}
& \mathbb{E}[\|\mathbf{P}_\ell^{(t)}(\mathbf{P}_\ell^{(t)})^\top [(1 - \beta_1)\tilde{\mathbf{M}}_\ell^{(t-1)} + \beta_1\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})]\|_F^2] \\
& \leq \left(1 - \left(1 - \frac{\delta_\ell}{4}\right)\beta_1\right)\mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2] + \frac{\beta_1^2\sigma^2}{\mathcal{B}} \\
& \quad + \frac{5(1 - \beta_1)}{\delta_\ell\beta_1}\mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2]. \tag{41}
\end{aligned}$$

For the second term in (38), we have

$$\begin{aligned}
& \mathbb{E}[\|(\mathbf{I} - \mathbf{P}_\ell^{(t)}(\mathbf{P}_\ell^{(t)})^\top)\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
& \leq 2\mathbb{E}[\|(\mathbf{I} - \mathbf{P}_\ell^{(t)}(\mathbf{P}_\ell^{(t)})^\top)\mathbf{G}_\ell^{(t)}\|_F^2] + 2\mathbb{E}[\|(\mathbf{I} - \mathbf{P}_\ell^{(t)}(\mathbf{P}_\ell^{(t)})^\top)(\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)}))\|_F^2] \\
& \leq 2(1 - \delta_\ell)\mathbb{E}[\|\mathbf{G}_\ell^{(t)}\|_F^2] + 2\mathbb{E}[\|\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
& \leq 2(1 - \delta_\ell)\mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + \frac{4\sigma_\ell^2}{\mathcal{B}}, \tag{42}
\end{aligned}$$

where the first inequality uses Cauchy's inequality, the second inequality uses Lemma 1 and $\|\mathbf{I} - \mathbf{P}_\ell^{(t)}(\mathbf{P}_\ell^{(t)})^\top\|_2 = 1$, the third inequality uses Assumption 3. Applying (41)(42) to (38) and using Lemma 2 yields (34).

When $t = k\tau + r$, $k \in \mathbb{N}$, $1 \leq r < \tau$, we have

$$\begin{aligned}
& \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
&= \mathbb{E}[\|(1 - \beta_1)(\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)})) + \beta_1(\mathbf{P}_\ell^{(t)}(\mathbf{P}_\ell^{(t)})^\top \mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)}))\|_F^2] \\
&= \mathbb{E}[\|(1 - \beta_1)(\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)})) + \beta_1(\mathbf{P}_\ell^{(t)}(\mathbf{P}_\ell^{(t)})^\top - \mathbf{I})\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
&\quad + \beta_1^2 \mathbb{E}[\|\mathbf{P}_\ell^{(t)}(\mathbf{P}_\ell^{(t)})^\top (\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)}))\|_F^2] \\
&\leq (1 - \beta_1) \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + \beta_1 \mathbb{E}[\|(\mathbf{I} - \mathbf{P}_\ell^{(t)}(\mathbf{P}_\ell^{(t)})^\top)\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
&\quad + \beta_1^2 \mathbb{E}[\|\mathbf{P}_\ell^{(t)}(\mathbf{P}_\ell^{(t)})^\top (\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)}))\|_F^2], \tag{43}
\end{aligned}$$

where the second equality uses the unbiasedness of $\mathbf{G}_\ell^{(t)}$ and the independence implied by $\mathbf{P}_\ell^{(t)} = \mathbf{P}_\ell^{(t-1)}$, the inequality uses Jensen's inequality. The first term is similarly bounded as (40). For the second term, we have

$$\begin{aligned}
& \mathbb{E}[\|(\mathbf{I} - \mathbf{P}_\ell^{(k\tau)}(\mathbf{P}_\ell^{(k\tau)})^\top)\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
&\leq \left(1 + \frac{\delta_\ell}{4}\right) \mathbb{E}[\|(\mathbf{I} - \mathbf{P}_\ell^{(k\tau)}(\mathbf{P}_\ell^{(k\tau)})^\top)\mathbf{G}_\ell^{(k\tau)}\|_F^2] \\
&\quad + \left(1 + \frac{4}{\delta_\ell}\right) \mathbb{E}[\|(\mathbf{I} - \mathbf{P}_\ell^{(k\tau)}(\mathbf{P}_\ell^{(k\tau)})^\top)(\nabla_\ell f(\mathbf{x}^{(t)}) - \mathbf{G}_\ell^{(k\tau)})\|_F^2] \\
&\leq \left(1 - \frac{3\delta_\ell}{4}\right) \mathbb{E}[\|\mathbf{G}_\ell^{(k\tau)}\|_F^2] + 2 \left(1 + \frac{4}{\delta_\ell}\right) \mathbb{E}[\|\mathbf{G}_\ell^{(k\tau)} - \nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2] \\
&\quad + 2 \left(1 + \frac{4}{\delta_\ell}\right) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2], \tag{44}
\end{aligned}$$

where the first inequality uses Young's inequality, the second inequality uses Lemma 1 and Cauchy's inequality. We further have

$$\begin{aligned}
& \left(1 - \frac{3\delta_\ell}{4}\right) \mathbb{E}[\|\mathbf{G}_\ell^{(k\tau)}\|_F^2] + 2 \left(1 + \frac{4}{\delta_\ell}\right) \mathbb{E}[\|\mathbf{G}_\ell^{(k\tau)} - \nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2] \\
&\leq \left(1 - \frac{3\delta_\ell}{4}\right) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2] + \frac{11}{\delta_\ell} \mathbb{E}[\|\mathbf{G}_\ell^{(k\tau)} - \nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2] \\
&\leq \left(1 - \frac{3\delta_\ell}{4}\right) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2] + \frac{11\sigma_\ell^2}{\delta_\ell \mathcal{B}} \\
&\leq \left(1 - \frac{\delta_\ell}{2}\right) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + \left(1 + \frac{4}{\delta_\ell}\right) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2] + \frac{11\sigma_\ell^2}{\delta_\ell \mathcal{B}}, \tag{45}
\end{aligned}$$

where the first inequality uses unbiasedness of $\mathbf{G}_\ell^{(k\tau)}$, the second inequality uses Assumption 3, the third inequality uses Young's inequality.

Applying (45) to (44) and applying Cauchy's inequality yields

$$\begin{aligned}
& \mathbb{E}[\|(\mathbf{I} - \mathbf{P}_\ell^{(k\tau)}(\mathbf{P}_\ell^{(k\tau)})^\top)\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
&\leq \left(1 - \frac{\delta_\ell}{2}\right) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + \frac{11\sigma_\ell^2}{\delta_\ell \mathcal{B}} + \frac{15r}{\delta_\ell} \sum_{i=1}^r \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(k\tau+i)}) - \nabla_\ell f(\mathbf{x}^{(k\tau+i-1)})\|_F^2]. \tag{46}
\end{aligned}$$

For the third term, we have

$$\mathbb{E}[\|\mathbf{P}_\ell^{(k\tau)}(\mathbf{P}_\ell^{(k\tau)})^\top (\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)}))\|_F^2] \leq \mathbb{E}[\|\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \leq \sigma_\ell^2, \tag{47}$$

where the first inequality uses $\|\mathbf{P}_\ell^{(k\tau)}(\mathbf{P}_\ell^{(k\tau)})^\top\|_2 = 1$, the second inequality uses Assumption 3.

Applying (40)(46)(47) to (43) yields (35). \square

Lemma 9 (Momentum error). *Under Assumption 2-3, if $0 < \beta_1 \leq 1$ in large-batch GaLore using MSGD and MP (Alg. 3), it holds for any $K \geq 1$ that*

$$\begin{aligned}
& \sum_{t=0}^{K\tau-1} \mathbb{E}[\|\tilde{\mathbf{m}}^{(t)} - \nabla f(\mathbf{x}^{(t)})\|_2^2] \\
& \leq \left(\frac{5(1-\beta_1)}{(1-\underline{\delta}/4)\underline{\delta}\beta_1^2} + \frac{15\tau(\tau-1)}{2(1-\underline{\delta}/4)\underline{\delta}} + \frac{2(\tau-1)}{(1-\bar{\delta}/4)\beta_1} \right) L^2 \sum_{t=0}^{K\tau-2} \mathbb{E}[\|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_2^2] \\
& \quad + \left(\frac{1-\underline{\delta}/2}{1-\underline{\delta}/4} + \frac{4}{(1-\bar{\delta}/4)\tau\beta_1} \right) \sum_{t=0}^{K\tau-1} \mathbb{E}[\|\nabla f(\mathbf{x}^{(t)})\|_2^2] \\
& \quad + \left(\frac{5K}{(1-\bar{\delta}/4)\beta_1\mathcal{B}} + \frac{11K\tau}{(1-\underline{\delta}/4)\underline{\delta}\mathcal{B}} + \frac{K\tau\beta_1}{1-\bar{\delta}/4} \right) \sigma^2. \tag{48}
\end{aligned}$$

Proof. By Lemma 8 we have

$$\begin{aligned}
& \sum_{t=0}^{K\tau-1} \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] - \left(1 - \left(1 - \frac{\delta_\ell}{4}\right)\beta_1\right) \sum_{t=0}^{K\tau-2} \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
& \leq \left(\frac{5(1-\beta_1)}{\delta_\ell\beta_1} + \frac{15\tau(\tau-1)\beta_1}{2\delta_\ell} + 2(\tau-1) \right) \sum_{t=0}^{K\tau-2} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t+1)}) - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
& \quad + \left(\frac{4}{\tau} + \left(1 - \frac{\delta_\ell}{2}\right)\beta_1 \right) \sum_{t=0}^{K\tau-1} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + \left(\frac{5K}{\mathcal{B}} + \frac{11K\tau\beta_1}{\delta_\ell\mathcal{B}} + K\tau\beta_1^2 \right) \sigma_\ell^2,
\end{aligned}$$

which implies

$$\begin{aligned}
& \sum_{t=0}^{K\tau-1} \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
& \leq \left(\frac{5(1-\beta_1)}{(1-\delta_\ell/4)\delta_\ell\beta_1^2} + \frac{15\tau(\tau-1)}{2(1-\delta_\ell/4)\delta_\ell} + \frac{2(\tau-1)}{(1-\delta_\ell/4)\beta_1} \right) \sum_{t=0}^{K\tau-2} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t+1)}) - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
& \quad + \left(\frac{1-\delta_\ell/2}{1-\delta_\ell/4} + \frac{4}{(1-\delta_\ell/4)\tau\beta_1} \right) \sum_{t=0}^{K\tau-1} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
& \quad + \left(\frac{5K}{(1-\delta_\ell/4)\beta_1\mathcal{B}} + \frac{11K\tau}{(1-\delta_\ell/4)\delta_\ell\mathcal{B}} + \frac{K\tau\beta_1}{1-\delta_\ell/4} \right) \sigma_\ell^2. \tag{49}
\end{aligned}$$

Summing (49) for $\ell = 1, \dots, N_L$ and applying Assumption 2-3 yields (48). \square

Now we are ready to prove the convergence of Alg. 3.

Theorem 7 (Convergence of large-batch GaLore). *Under Assumptions 1-3, if hyperparameters*

$$0 < \beta_1 \leq 1, \quad \tau \geq \frac{128}{3\beta_1\underline{\delta}}, \quad 0 < \eta \leq \min \left\{ \frac{1}{4L}, \sqrt{\frac{3\underline{\delta}\beta_1^2}{80L^2}}, \sqrt{\frac{\underline{\delta}}{40\tau^2L^2}}, \sqrt{\frac{3\beta_1}{32\tau L^2}} \right\}, \tag{50}$$

GaLore using large-batch stochastic gradients and MSGD with MP (Alg. 3) converges as

$$\frac{1}{K\tau} \sum_{t=0}^{K\tau-1} \mathbb{E}[\|\nabla f(\mathbf{x}^{(t)})\|_2^2] \leq \frac{16\Delta}{\underline{\delta}\eta K\tau} + \left(\frac{160}{3\beta_1\underline{\delta}\tau\mathcal{B}} + \frac{352}{3\underline{\delta}^2\mathcal{B}} + \frac{32\beta_1}{3\underline{\delta}} \right) \sigma^2 \tag{51}$$

for any $K \geq 1$, where $\Delta = f(\mathbf{x}^{(0)}) - \inf_{\mathbf{x}} f(\mathbf{x})$.

1458 *Proof.* By Lemma 4 we have

$$1459 \sum_{t=0}^{K\tau-1} \mathbb{E}[\|\nabla f(\mathbf{x}^{(t)})\|_2^2] \leq \frac{2[f(\mathbf{x}^{(0)}) - \mathbb{E}[f(\mathbf{x}^{(K\tau)})]]}{\eta} + \sum_{t=0}^{K\tau-1} \mathbb{E}[\|\tilde{\mathbf{m}}^{(t)} - \nabla f(\mathbf{x}^{(t)})\|_2^2] \\ 1461 \\ 1462 - \left(\frac{1}{\eta^2} - \frac{L}{\eta}\right) \sum_{t=0}^{K\tau-1} \mathbb{E}[\|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_2^2]. \quad (52)$$

1463 Applying Lemma 9 to (52) and using $\underline{\delta} \leq \bar{\delta} < 1$ yields

$$1464 \left(\frac{\underline{\delta}}{4} - \frac{16}{3\tau\beta_1}\right) \sum_{t=0}^{K\tau-1} \mathbb{E}[\|\nabla f(\mathbf{x}^{(t)})\|_2^2] \\ 1465 \\ 1466 \leq \frac{2}{\eta} \mathbb{E}[f(\mathbf{x}^{(0)}) - f(\mathbf{x}^{(K\tau)})] + \left(\frac{20K}{3\beta_1\mathcal{B}} + \frac{44K\tau}{3\underline{\delta}\mathcal{B}} + \frac{4K\tau\beta_1}{3}\right) \sigma^2 \\ 1467 \\ 1468 - \left(\frac{1}{\eta^2} - \frac{L}{\eta} - \frac{20(1-\beta_1)L^2}{3\underline{\delta}\beta_1^2} - \frac{10\tau(\tau-1)L^2}{\underline{\delta}} - \frac{8(\tau-1)L^2}{3\beta_1}\right) \sum_{t=0}^{K\tau-1} \mathbb{E}[\|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_2^2]. \\ 1469 \\ 1470 \\ 1471 \\ 1472 \\ 1473 \\ 1474 \\ 1475 \\ 1476 \\ 1477 \quad (53)$$

By (50) we have

$$1478 \frac{\underline{\delta}}{4} - \frac{16}{3\tau\beta_1} \geq \frac{\underline{\delta}}{8}, \quad \text{and} \quad \frac{1}{4\eta^2} \geq \max \left\{ \frac{L}{\eta}, \frac{20(1-\beta_1)L^2}{3\underline{\delta}\beta_1^2}, \frac{10\tau(\tau-1)L^2}{\underline{\delta}}, \frac{8(\tau-1)L^2}{3\beta_1} \right\}. \quad (54)$$

1479 Applying (54) to (53) yields (51). \square

1480 We now prove Theorem 3, which is restated as follows.

1481 **Corollary 2** (Convergence complexity of large-batch GaLore). *Under Assumptions 1-3, if $T \geq 2 + 256/(3\underline{\delta}) + (256\sigma)^2/(9\sqrt{\underline{\delta}}L\Delta)$ and we choose*

$$1482 \beta_1 = \left(1 + \sqrt{\frac{\underline{\delta}^{3/2}\sigma^2 T}{L\Delta}}\right)^{-1}, \\ 1483 \\ 1484 \tau = \left\lceil \frac{128}{3\underline{\delta}\beta_1} \right\rceil, \\ 1485 \\ 1486 \eta = \left(4L + \sqrt{\frac{80L^2}{3\underline{\delta}\beta_1^2}} + \sqrt{\frac{40\tau^2 L^2}{\underline{\delta}}} + \sqrt{\frac{32\tau L^2}{3\beta_1}}\right)^{-1}, \\ 1487 \\ 1488 \mathcal{B} = \left\lceil \frac{1}{\underline{\delta}\beta_1} \right\rceil, \\ 1489 \\ 1490 \\ 1491 \\ 1492 \\ 1493 \\ 1494 \\ 1495 \\ 1496$$

1497 *GaLore using large-batch stochastic gradients and MSGD with MP (Alg. 3) converges as*

$$1498 \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(\mathbf{x}^{(t)})\|_2^2] = \mathcal{O} \left(\frac{L\Delta}{\underline{\delta}^{5/2}T} + \sqrt{\frac{L\Delta\sigma^2}{\underline{\delta}^{7/2}T}} \right), \quad (55)$$

1499 where $\Delta = f(\mathbf{x}^{(0)}) - \inf_{\mathbf{x}} f(\mathbf{x})$. Consequently, the computation complexity to reach an ε -accurate solution \mathbf{x} such that $\|\nabla f(\mathbf{x})\|_2^2 \leq \varepsilon$ is $\mathcal{O} \left(\frac{L\Delta\sigma^2}{\underline{\delta}^{7/2}\varepsilon^2} + \frac{L\Delta}{\underline{\delta}^{5/2}\varepsilon} + \frac{\sigma^2}{\underline{\delta}^{1/2}L\Delta} + \frac{1}{\underline{\delta}} \right)$.

1500 *Proof.* $T \geq 2 + 128/(3\underline{\delta}) + (128\sigma)^2/(9\sqrt{\underline{\delta}}L\Delta)$ guarantees $T \geq \tau$. Let $T = K\tau + r$, where $K \in \mathbb{N}^*$ and $0 \leq r < \tau$. If $r = 0$, (55) is a direct result of Theorem 7. If $r > 0$, applying Theorem 7 to $\tilde{K} := K + 1$ yields

$$1501 \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(\mathbf{x}^{(t)})\|_2^2] \leq \frac{\tilde{K}\tau}{T} \cdot \frac{1}{\tilde{K}\tau} \sum_{t=0}^{\tilde{K}\tau-1} \mathbb{E}[\|\nabla f(\mathbf{x}^{(t)})\|_2^2] = \mathcal{O} \left(\frac{L\Delta}{\underline{\delta}^{5/2}T} + \sqrt{\frac{L\Delta\sigma^2}{\underline{\delta}^{7/2}T}} \right).$$

1502 \square

Algorithm 4 GoLore using small-batch stochastic gradients and MSGD with MP

Input: Initial point $\mathbf{x}^{(0)}$, data distribution \mathcal{D} , learning rate η , subspace changing frequency τ , rank $\{r_\ell\}_{\ell=1}^{N_L}$, momentum parameter β_1 .

Output: $\{\mathbf{x}^{(t)}\}_{t=0}^T$.

Initialize optimizer state $\{M_\ell^{(-1)}\}_{\ell=1}^{N_L}$ to zero;

for $t = 0, 1, \dots, T - 1$ **do**

 Sample $\xi^{(t)} \sim \mathcal{D}$;

$\mathbf{G}_\ell^{(t)} = \nabla_\ell F(\mathbf{x}^{(t)}; \xi^{(t)})$;

for $\ell = 1, 2, \dots, N_L$ **do**

if $t \equiv 0 \pmod{\tau}$ **then**

if $m_\ell \leq n_\ell$ **then**

 Sample $\mathbf{P}_\ell^{(t)} \sim \mathcal{U}(\text{St}_{m_\ell, r_\ell})$;

$\mathbf{M}_\ell^{(t)} \leftarrow (1 - \beta_1)(\mathbf{P}_\ell^{(t)})^\top \mathbf{P}_\ell^{(t-1)} \mathbf{M}_\ell^{(t-1)} + \beta_1(\mathbf{P}_\ell^{(t)})^\top \mathbf{G}_\ell^{(t)}$;

$\mathbf{X}_\ell^{(t+1)} \leftarrow \mathbf{X}_\ell^{(t)} - \eta \mathbf{P}_\ell^{(t)} \mathbf{M}_\ell^{(t)}$;

else

 Sample $\mathbf{Q}_\ell^{(t)} \sim \mathcal{U}(\text{St}_{n_\ell, r_\ell})$;

$\mathbf{M}_\ell^{(t)} \leftarrow (1 - \beta_1)\mathbf{M}_\ell^{(t-1)}(\mathbf{Q}_\ell^{(t-1)})^\top \mathbf{Q}_\ell^{(t)} + \beta_1\mathbf{G}_\ell^{(t)}\mathbf{Q}_\ell^{(t)}$;

$\mathbf{X}_\ell^{(t+1)} \leftarrow \mathbf{X}_\ell^{(t)} - \eta \mathbf{M}_\ell^{(t)}(\mathbf{Q}_\ell^{(t)})^\top$;

end if

else

if $m_\ell \leq n_\ell$ **then**

$\mathbf{P}_\ell^{(t)} \leftarrow \mathbf{P}_\ell^{(t-1)}$;

$\mathbf{M}_\ell^{(t)} \leftarrow (1 - \beta_1)\mathbf{M}_\ell^{(t-1)} + \beta_1(\mathbf{P}_\ell^{(t)})^\top \mathbf{G}_\ell^{(t)}$;

$\mathbf{X}_\ell^{(t+1)} \leftarrow \mathbf{X}_\ell^{(t)} - \eta \mathbf{P}_\ell^{(t)} \mathbf{M}_\ell^{(t)}$;

else

$\mathbf{Q}_\ell^{(t)} \leftarrow \mathbf{Q}_\ell^{(t-1)}$;

$\mathbf{M}_\ell^{(t)} \leftarrow (1 - \beta_1)\mathbf{M}_\ell^{(t-1)} + \beta_1\mathbf{G}_\ell^{(t)}\mathbf{Q}_\ell^{(t)}$;

$\mathbf{X}_\ell^{(t+1)} \leftarrow \mathbf{X}_\ell^{(t)} - \eta \mathbf{M}_\ell^{(t)}(\mathbf{Q}_\ell^{(t)})^\top$;

end if

end if

end for

end for

B.5 CONVERGENCE OF GOLORE

In this subsection, we present the proof for Theorem 4. GoLore using small-batch stochastic gradients and MSGD with MP is specified as Alg. 4.

Lemma 10 (Momentum contraction). *Under Assumption 3, in large-batch GoLore using MSGD with MP (Alg. 4), if $0 < \beta_1 \leq 1$, term $\tilde{\mathbf{M}}_\ell^{(t)}$ has the following contraction properties:*

- When $t = 0$, it holds that

$$\begin{aligned} \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(0)} - \nabla_\ell f(\mathbf{X}^{(0)})\|_F^2] &\leq (\tau - 1)(1 - \delta_\ell \beta_1) \sum_{r=0}^{\tau-2} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(r+1)}) - \nabla_\ell f(\mathbf{x}^{(r)})\|_F^2] \\ &\quad + \frac{2(1 - \delta_\ell \beta_1)}{\tau} \sum_{r=0}^{\tau-1} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(r)})\|_F^2] + \delta_\ell \beta_1^2 \sigma_\ell^2; \end{aligned} \quad (56)$$

1566

- When $t = k\tau$, $k \in \mathbb{N}^*$, it holds that

1567

1568

1569

1570

1571

1572

1573

1574

1575

$$\begin{aligned} & \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] - \delta_\ell \left(1 - \left(1 - \frac{\delta_\ell}{4}\right)\beta_1\right) \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2] \\ & \leq \frac{2(1 - \delta_\ell)}{\tau} \sum_{r=0}^{\tau-1} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(k\tau+r)})\|_F^2] + \frac{5(1 - \beta_1)}{\beta_1} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2] \\ & \quad + (\tau - 1)(1 - \delta_\ell) \sum_{r=0}^{\tau-2} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(k\tau+r+1)}) - \nabla_\ell f(\mathbf{x}^{(k\tau+r)})\|_F^2] + \delta_\ell \beta_1^2 \sigma_\ell^2; \end{aligned} \quad (57)$$

1576

1577

- When $t = k\tau + r$, $k \in \mathbb{N}$, $1 \leq r < \tau$, it holds that

1578

1579

1580

1581

1582

1583

1584

1585

$$\begin{aligned} & \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] - \left(1 - \left(1 - \frac{\delta_\ell}{4}\right)\beta_1\right) \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2] \\ & \leq \left(1 - \frac{\delta_\ell}{2}\right) \beta_1 \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + \frac{5(1 - \beta_1)}{\delta_\ell \beta_1} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2] \\ & \quad + \frac{10r\beta_1}{\delta_\ell} \sum_{i=1}^r \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(k\tau+i)}) - \nabla_\ell f(\mathbf{x}^{(k\tau+i-1)})\|_F^2] + \beta_1^2 \sigma_\ell^2. \end{aligned} \quad (58)$$

1586

Proof. Without loss of generality assume $m_\ell \leq n_\ell$ (the other case can be proved similarly). When $t = 0$, we have

1588

1589

1590

1591

1592

1593

1594

1595

$$\begin{aligned} & \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2] \\ & = \mathbb{E}[\|\beta_1 \mathbf{P}_\ell^{(0)} (\mathbf{P}_\ell^{(0)})^\top \mathbf{G}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2] \\ & = \mathbb{E}[\|(\beta_1 \mathbf{P}_\ell^{(0)} (\mathbf{P}_\ell^{(0)})^\top - \mathbf{I}) \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2] + \beta_1^2 \mathbb{E}[\|\mathbf{P}_\ell^{(0)} (\mathbf{P}_\ell^{(0)})^\top (\mathbf{G}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)}))\|_F^2] \\ & = \text{tr}((\nabla_\ell f(\mathbf{x}^{(0)}))^\top \mathbb{E}[(\beta_1 \mathbf{P}_\ell^{(0)} (\mathbf{P}_\ell^{(0)})^\top - \mathbf{I})^2] \nabla_\ell f(\mathbf{x}^{(0)})) \\ & \quad + \beta_1^2 \text{tr}(\mathbb{E}_{\xi^{(0)} \sim \mathcal{D}}[(\mathbf{G}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)}))^\top \mathbb{E}_{\mathbf{P} \sim \mathcal{U}(\text{St}_{m_\ell, r_\ell})}[(\mathbf{P}\mathbf{P}^\top)^2] (\mathbf{G}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)}))]), \end{aligned} \quad (59)$$

1596

where the second equality uses unbiasedness of $\mathbf{G}_\ell^{(0)}$. By Lemma 5 we have

1597

1598

1599

1600

$$\begin{aligned} \mathbb{E}[(\beta_1 \mathbf{P}_\ell^{(0)} (\mathbf{P}_\ell^{(0)})^\top - \mathbf{I})^2] & = \mathbf{I} - (2\beta_1 - \beta_1^2) \mathbb{E}[\mathbf{P}_\ell^{(0)} (\mathbf{P}_\ell^{(0)})^\top] \\ & = \mathbf{I} - (2\beta_1 - \beta_1^2) \delta_\ell \mathbf{I}, \end{aligned}$$

1601

thus

1602

1603

1604

$$\begin{aligned} \text{tr}((\nabla_\ell f(\mathbf{x}^{(0)}))^\top \mathbb{E}[(\beta_1 \mathbf{P}_\ell^{(0)} (\mathbf{P}_\ell^{(0)})^\top - \mathbf{I})^2] \nabla_\ell f(\mathbf{x}^{(0)})) & = (1 - \delta_\ell (2\beta_1 - \beta_1^2)) \|\nabla_\ell f(\mathbf{x}^{(0)})\|_F^2 \\ & \leq (1 - \delta_\ell \beta_1) \|\nabla_\ell f(\mathbf{x}^{(0)})\|_F^2. \end{aligned} \quad (60)$$

1605

Similarly, by Lemma 5 we have

1606

1607

1608

1609

1610

1611

1612

1613

1614

$$\begin{aligned} & \text{tr}(\mathbb{E}_{\xi^{(0)} \sim \mathcal{D}}[(\mathbf{G}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)}))^\top \mathbb{E}_{\mathbf{P} \sim \mathcal{U}(\text{St}_{m_\ell, r_\ell})}[(\mathbf{P}\mathbf{P}^\top)^2] (\mathbf{G}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)}))]) \\ & = \text{tr} \left(\mathbb{E} \left[(\mathbf{G}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)}))^\top \left(\frac{r_\ell}{m_\ell} \cdot \mathbf{I} \right) (\mathbf{G}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)})) \right] \right) \\ & = \delta_\ell \mathbb{E}[\|\mathbf{G}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2] \\ & \leq \delta_\ell \sigma_\ell^2, \end{aligned} \quad (61)$$

where the inequality uses Assumption 3. Applying (60)(61) and Lemma 2 to (59) yields (56).

1615

When $t = k\tau$, $k \in \mathbb{N}^*$, we have

1616

1617

1618

1619

$$\begin{aligned} & \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\ & = \mathbb{E}[\|\mathbf{P}_\ell^{(t)} (\mathbf{P}_\ell^{(t)})^\top [(1 - \beta_1) \tilde{\mathbf{M}}_\ell^{(t-1)} + \beta_1 \mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})] - (\mathbf{I} - \mathbf{P}_\ell^{(t)} (\mathbf{P}_\ell^{(t)})^\top) \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\ & = \delta_\ell \mathbb{E}[\|(1 - \beta_1) \tilde{\mathbf{M}}_\ell^{(t-1)} + \beta_1 \mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + (1 - \delta_\ell) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2], \end{aligned} \quad (62)$$

where the second equality uses Lemma 3 and Lemma 5. For the first term, we have

$$\begin{aligned}
& \mathbb{E}[\|(1 - \beta_1)\tilde{\mathbf{M}}_\ell^{(t-1)} + \beta_1\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
&= \mathbb{E}[\|(1 - \beta_1)(\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)})) + \beta_1(\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)}))\|_F^2] \\
&\leq \mathbb{E}[\|(1 - \beta_1)(\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)}))\|_F^2] + \beta_1^2 \mathbb{E}[\|\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
&\leq (1 - \beta_1) \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + \beta_1^2 \sigma_\ell^2,
\end{aligned} \tag{63}$$

where both inequalities use Assumption 3. By Young's inequality, we have

$$\begin{aligned}
& \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
&= \mathbb{E}[\|(\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t-1)})) - (\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(t-1)}))\|_F^2] \\
&\leq \left(1 + \frac{\delta_\ell \beta_1}{4}\right) \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2] + \left(1 + \frac{4}{\delta_\ell \beta_1}\right) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2].
\end{aligned} \tag{64}$$

Applying (63)(64) and Lemma 2 to (62) yields (57).

When $t = k\tau + r$, $k \in \mathbb{N}$, $1 \leq r < \tau$, we have

$$\begin{aligned}
& \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
&= \mathbb{E}[\|(1 - \beta_1)(\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)})) + \beta_1(\mathbf{P}_\ell^{(t)}(\mathbf{P}_\ell^{(t)})^\top \mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)}))\|_F^2] \\
&= \mathbb{E}[\|(1 - \beta_1)(\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)})) + \beta_1(\mathbf{P}_\ell^{(t)}(\mathbf{P}_\ell^{(t)})^\top - \mathbf{I})\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
&\quad + \beta_1^2 \mathbb{E}[\|\mathbf{P}_\ell^{(t)}(\mathbf{P}_\ell^{(t)})^\top (\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)}))\|_F^2] \\
&\leq (1 - \beta_1) \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + \beta_1 \mathbb{E}[\|(\mathbf{I} - \mathbf{P}_\ell^{(t)}(\mathbf{P}_\ell^{(t)})^\top)\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
&\quad + \beta_1^2 \mathbb{E}[\|\mathbf{P}_\ell^{(t)}(\mathbf{P}_\ell^{(t)})^\top (\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)}))\|_F^2],
\end{aligned} \tag{65}$$

where the second equality uses the unbiasedness of $\mathbf{G}_\ell^{(t)}$ and the independence implied by $\mathbf{P}_\ell^{(t)} = \mathbf{P}_\ell^{(t-1)}$, the inequality uses Jensen's inequality. The first term is similarly bounded as (64). For the second term, we have

$$\begin{aligned}
& \mathbb{E}[\|(\mathbf{I} - \mathbf{P}_\ell^{(k\tau)}(\mathbf{P}_\ell^{(k\tau)})^\top)\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
&\leq \left(1 + \frac{\delta_\ell}{4}\right) \mathbb{E}[\|(\mathbf{I} - \mathbf{P}_\ell^{(k\tau)}(\mathbf{P}_\ell^{(k\tau)})^\top)\nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2] \\
&\quad + \left(1 + \frac{4}{\delta_\ell}\right) \mathbb{E}[\|(\mathbf{I} - \mathbf{P}_\ell^{(k\tau)}(\mathbf{P}_\ell^{(k\tau)})^\top)(\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(k\tau)}))\|_F^2] \\
&\leq \left(1 - \frac{3\delta_\ell}{4}\right) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2] + \left(1 + \frac{4}{\delta_\ell}\right) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2],
\end{aligned} \tag{66}$$

where the first inequality uses Young's inequality, the second inequality uses Lemma 5 and $\|\mathbf{I} - \mathbf{P}_\ell^{(k\tau)}(\mathbf{P}_\ell^{(k\tau)})^\top\|_2 = 1$. By Young's inequality, we have

$$\mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2] \leq \left(1 + \frac{\delta_\ell}{4}\right) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + \left(1 + \frac{4}{\delta_\ell}\right) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2]. \tag{67}$$

Applying (67) to (66) and applying Cauchy's inequality yields

$$\begin{aligned}
& \mathbb{E}[\|(\mathbf{I} - \mathbf{P}_\ell^{(k\tau)}(\mathbf{P}_\ell^{(k\tau)})^\top)\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
&\leq \left(1 - \frac{\delta_\ell}{2}\right) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + \frac{10r}{\delta_\ell} \sum_{i=1}^r \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(k\tau+i)}) - \nabla_\ell f(\mathbf{x}^{(k\tau+i-1)})\|_F^2].
\end{aligned} \tag{68}$$

For the third term, we have

$$\mathbb{E}[\|\mathbf{P}_\ell^{(k\tau)}(\mathbf{P}_\ell^{(k\tau)})^\top (\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)}))\|_F^2] \leq \mathbb{E}[\|\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \leq \sigma_\ell^2, \tag{69}$$

where the first inequality uses $\|\mathbf{P}_\ell^{(k\tau)}(\mathbf{P}_\ell^{(k\tau)})^\top\|_2 = 1$, the second inequality uses Assumption 3.

Applying (64)(68)(69) to (65) yields (58). \square

Lemma 11 (Momentum error). *Under Assumption 2-3, if $0 < \beta_1 \leq 1$ in GoLore using MSGD and MP (Alg. 4), it holds for any $K \geq 1$ that*

$$\begin{aligned}
& \sum_{t=0}^{K\tau-1} \mathbb{E}[\|\tilde{\mathbf{m}}^{(t)} - \nabla f(\mathbf{x}^{(t)})\|_2^2] \\
& \leq \left(\frac{5(1-\beta_1)}{(1-\bar{\delta}/4)\bar{\delta}\beta_1^2} + \frac{5\tau(\tau-1)}{(1-\bar{\delta}/4)\bar{\delta}} + \frac{\tau-1}{(1-\bar{\delta}/4)\beta_1} \right) L^2 \sum_{t=0}^{K\tau-2} \mathbb{E}[\|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_2^2] \\
& \quad + \left(\frac{1-\bar{\delta}/2}{1-\bar{\delta}/4} + \frac{2}{(1-\bar{\delta}/4)\tau\beta_1} \right) \sum_{t=0}^{K\tau-1} \mathbb{E}[\|\nabla f(\mathbf{x}^{(t)})\|_2^2] + \frac{K\tau\beta_1\sigma^2}{1-\bar{\delta}/4}. \tag{70}
\end{aligned}$$

Proof. By Lemma 10 we have

$$\begin{aligned}
& \sum_{t=0}^{K\tau-1} \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] - \left(1 - \left(1 - \frac{\delta_\ell}{4}\right)\beta_1\right) \sum_{t=0}^{K\tau-2} \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
& \leq \left(\frac{5(1-\beta_1)}{\delta_\ell\beta_1} + \frac{5\tau(\tau-1)\beta_1}{\delta_\ell} + (\tau-1) \right) \sum_{t=0}^{K\tau-2} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t+1)}) - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
& \quad + \left(\frac{2}{\tau} + \left(1 - \frac{\delta_\ell}{2}\right)\beta_1 \right) \sum_{t=0}^{K\tau-1} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + K\tau\beta_1^2\sigma_\ell^2,
\end{aligned}$$

which implies

$$\begin{aligned}
& \sum_{t=0}^{K\tau-1} \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
& \leq \left(\frac{5(1-\beta_1)}{(1-\delta_\ell/4)\delta_\ell\beta_1^2} + \frac{5\tau(\tau-1)}{(1-\delta_\ell/4)\delta_\ell} + \frac{\tau-1}{(1-\delta_\ell/4)\beta_1} \right) \sum_{t=0}^{K\tau-2} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t+1)}) - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
& \quad + \left(\frac{1-\delta_\ell/2}{1-\delta_\ell/4} + \frac{2}{(1-\delta_\ell/4)\tau\beta_1} \right) \sum_{t=0}^{K\tau-1} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + \frac{K\tau\beta_1\sigma_\ell^2}{1-\delta_\ell/4}. \tag{71}
\end{aligned}$$

Summing (71) for $\ell = 1, \dots, N_L$ and applying Assumption 2-3 yields (70). \square

Now we are ready to prove the convergence of Alg. 4.

Theorem 8 (Convergence of Golore). *Under Assumptions 1-3, if hyperparameters*

$$0 < \beta_1 \leq 1, \quad \tau \geq \frac{64}{3\beta_1\bar{\delta}}, \quad 0 < \eta \leq \min \left\{ \frac{1}{4L}, \sqrt{\frac{3\bar{\delta}\beta_1^2}{80L^2}}, \sqrt{\frac{3\bar{\delta}}{80\tau^2L^2}}, \sqrt{\frac{3\beta_1}{16\tau L^2}} \right\}, \tag{72}$$

GoLore using small-batch stochastic gradients and MSGD with MP (Alg. 4) converges as

$$\frac{1}{K\tau} \sum_{t=0}^{K\tau-1} \mathbb{E}[\|\nabla f(\mathbf{x}^{(t)})\|_2^2] \leq \frac{16\Delta}{\bar{\delta}\eta K\tau} + \frac{32\beta_1\sigma^2}{3\bar{\delta}} \tag{73}$$

for any $K \geq 1$, where $\Delta = f(\mathbf{x}^{(0)}) - \inf_{\mathbf{x}} f(\mathbf{x})$.

Proof. By Lemma 4 we have

$$\begin{aligned}
\sum_{t=0}^{K\tau-1} \mathbb{E}[\|\nabla f(\mathbf{x}^{(t)})\|_2^2] & \leq \frac{2[f(\mathbf{x}^{(0)}) - \mathbb{E}[f(\mathbf{x}^{(K\tau)})]]}{\eta} + \sum_{t=0}^{K\tau-1} \mathbb{E}[\|\tilde{\mathbf{m}}^{(t)} - \nabla f(\mathbf{x}^{(t)})\|_2^2] \\
& \quad - \left(\frac{1}{\eta^2} - \frac{L}{\eta} \right) \sum_{t=0}^{K\tau-1} \mathbb{E}[\|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_2^2]. \tag{74}
\end{aligned}$$

Applying Lemma 11 to (74) and using $\underline{\delta} \leq \bar{\delta} < 1$ yields

$$\begin{aligned} & \left(\frac{\underline{\delta}}{4} - \frac{8}{3\tau\beta_1} \right) \sum_{t=0}^{K\tau-1} \mathbb{E}[\|\nabla f(\mathbf{x}^{(t)})\|_2^2] \\ & \leq \frac{2}{\eta} \mathbb{E}[f(\mathbf{x}^{(0)}) - f(\mathbf{x}^{(K\tau)})] + \frac{4K\tau\beta_1\sigma^2}{3} \\ & \quad - \left(\frac{1}{\eta^2} - \frac{L}{\eta} - \frac{20(1-\beta_1)L^2}{3\underline{\delta}\beta_1^2} - \frac{20\tau(\tau-1)L^2}{3\underline{\delta}} - \frac{4(\tau-1)L^2}{3\beta_1} \right) \sum_{t=0}^{K\tau-1} \mathbb{E}[\|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_2^2]. \end{aligned} \quad (75)$$

By (72) we have

$$\frac{\underline{\delta}}{4} - \frac{8}{3\tau\beta_1} \geq \frac{\underline{\delta}}{8}, \quad \text{and} \quad \frac{1}{4\eta^2} \geq \max \left\{ \frac{L}{\eta}, \frac{20(1-\beta_1)L^2}{3\underline{\delta}\beta_1^2}, \frac{20\tau(\tau-1)L^2}{3\underline{\delta}}, \frac{4(\tau-1)L^2}{3\beta_1} \right\}. \quad (76)$$

Applying (76) to (75) yields (73). \square

We now prove Theorem 4, which is restated as follows.

Corollary 3 (Convergence complexity of GoLore). *Under Assumptions 1-3, if $T \geq 2 + 128/(3\underline{\delta}) + (128\sigma)^2/(9\sqrt{\underline{\delta}}L\Delta)$ and we choose*

$$\begin{aligned} \beta_1 &= \left(1 + \sqrt{\frac{\underline{\delta}^{3/2}\sigma^2 T}{L\Delta}} \right)^{-1}, \\ \tau &= \left\lceil \frac{64}{3\underline{\delta}\beta_1} \right\rceil, \\ \eta &= \left(4L + \sqrt{\frac{80L^2}{3\underline{\delta}\beta_1^2}} + \sqrt{\frac{80\tau^2 L^2}{3\underline{\delta}}} + \sqrt{\frac{16\tau L^2}{3\beta_1}} \right)^{-1}, \end{aligned}$$

GoLore using small-batch stochastic gradients and MSGD with MP (Alg. 4) converges as

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(\mathbf{x}^{(t)})\|_2^2] = \mathcal{O} \left(\frac{L\Delta}{\underline{\delta}^{5/2}T} + \sqrt{\frac{L\Delta\sigma^2}{\underline{\delta}^{7/2}T}} \right), \quad (77)$$

where $\Delta = f(\mathbf{x}^{(0)}) - \inf_{\mathbf{x}} f(\mathbf{x})$. Consequently, the computation complexity to reach an ε -accurate solution \mathbf{x} such that $\|\nabla f(\mathbf{x})\|_2^2 \leq \varepsilon$ is $\mathcal{O} \left(\frac{L\Delta\sigma^2}{\underline{\delta}^{7/2}\varepsilon^2} + \frac{L\Delta}{\underline{\delta}^{5/2}\varepsilon} + \frac{\sigma^2}{\underline{\delta}^{1/2}L\Delta} + \frac{1}{\underline{\delta}} \right)$.

Proof. $T \geq 2 + 128/(3\underline{\delta}) + (128\sigma)^2/(9\sqrt{\underline{\delta}}L\Delta)$ guarantees $T \geq \tau$. Let $T = K\tau + r$, where $K \in \mathbb{N}^*$ and $0 \leq r < \tau$. If $r = 0$, (77) is a direct result of Theorem 8. If $r > 0$, applying Theorem 8 to $\tilde{K} := K + 1$ yields

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(\mathbf{x}^{(t)})\|_2^2] \leq \frac{\tilde{K}\tau}{T} \cdot \frac{1}{\tilde{K}\tau} \sum_{t=0}^{\tilde{K}\tau-1} \mathbb{E}[\|\nabla f(\mathbf{x}^{(t)})\|_2^2] = \mathcal{O} \left(\frac{L\Delta}{\underline{\delta}^{5/2}T} + \sqrt{\frac{L\Delta\sigma^2}{\underline{\delta}^{7/2}T}} \right).$$

\square

C RESULTS FOR SPARSE SUBSPACE OPTIMIZATION

In this section, we illustrate how to transfer the main results of this paper to sparse subspace optimization algorithms. We first present the detailed algorithm formulation, then present the theoretical results corresponding to GaLore/GoLore. Although it only requires little effort to transfer results in GaLore/GoLore to sparse subspace optimization, we still include proofs for completeness.

1782 **Algorithm 5** GaSare / GoSare algorithms using stochastic / deterministic / large-batch gra-
1783 dients
1784

1785 **Input:** Initial point $\mathbf{x}^{(0)}$, data distribution \mathcal{D} , learning rate η , subspace changing frequency τ , rank
1786 $\{r_\ell\}_{\ell=1}^{N_L}$, optimizer hyperparameters $\beta_1, \beta_2, \epsilon$, large batch size \mathcal{B} .
1787 **Output:** $\{\mathbf{x}^{(t)}\}_{t=0}^T$.

1788 Initialize optimizer state $\{M_\ell^{(-1)}\}_{\ell=1}^{N_L}$ and $\{V_\ell^{(-1)}\}_{\ell=1}^{N_L}$ to zero;
1789 **for** $t = 0, 1, \dots, T - 1$ **do**
1790 **for** $\ell = 1, 2, \dots, N_L$ **do**
1791 **if** $t \equiv 0 \pmod{\tau}$ **then**
1792 $\mathbf{G}_\ell^{(t)} \leftarrow \nabla_\ell F(\mathbf{x}^{(t)}; \xi^{(t)});$ (stochastic)
1793 $\mathbf{G}_\ell^{(t)} \leftarrow \nabla_\ell f(\mathbf{x}^{(t)});$ (deterministic)
1794 $\mathbf{G}_\ell^{(t)} \leftarrow \frac{1}{\mathcal{B}} \sum_{b=1}^{\mathcal{B}} \nabla_\ell F(\mathbf{x}^{(t)}; \xi^{(t,b)});$ (large-batch)
1795 $\mathbf{S}_\ell^{(t)} \leftarrow \text{Top}_k(\mathbf{G}_\ell^{(t)});$ (GaSare)
1796 Sample $\mathbf{S}_\ell^{(t)} \sim \mathcal{U}(\text{Sp}_{m_\ell, n_\ell}^{k_\ell});$ (GoSare)
1797 **else**
1798 $\mathbf{G}_\ell^{(t)} \leftarrow \nabla_\ell F(\mathbf{x}^{(t)}; \xi^{(t)});$ (stochastic)
1799 $\mathbf{G}_\ell^{(t)} \leftarrow \nabla_\ell f(\mathbf{x}^{(t)});$ (deterministic)
1800 $\mathbf{G}_\ell^{(t)} \leftarrow \nabla_\ell F(\mathbf{x}^{(t)}; \xi^{(t)});$ (large-batch)
1801 $\mathbf{S}_\ell^{(t)} \leftarrow \mathbf{S}_\ell^{(t-1)};$
1802 **end if**
1803 $\mathbf{R}_\ell^{(t)} \leftarrow \mathbf{S}_\ell^{(t)} \odot \mathbf{G}_\ell^{(t)};$
1804 $\mathbf{M}_\ell^{(t)} \leftarrow (1 - \beta_1)\mathbf{S}_\ell^{(t)} \odot \mathbf{M}_\ell^{(t-1)} + \beta_1\mathbf{R}_\ell^{(t)};$
1805 $\mathbf{V}_\ell^{(t)} \leftarrow (1 - \beta_2)\mathbf{S}_\ell^{(t)} \odot \mathbf{V}_\ell^{(t-1)} + \beta_2\mathbf{R}_\ell^{(t)} \odot \mathbf{R}_\ell^{(t)};$
1806 **if** using Adam **then**
1807 $\mathbf{M}_\ell^{(t)} \leftarrow \mathbf{M}_\ell^{(t)} / (1 - \beta_1^t), \quad \mathbf{V}_\ell^{(t)} \leftarrow \mathbf{V}_\ell^{(t)} / (1 - \beta_2^t), \quad \mathbf{N}_\ell^{(t)} \leftarrow \mathbf{M}_\ell^{(t)} / (\sqrt{\mathbf{V}_\ell^{(t)}} + \epsilon);$
1808 **else if** using MSGD **then**
1809 $\mathbf{N}_\ell^{(t)} \leftarrow \mathbf{M}_\ell^{(t)};$
1810 **end if**
1811 $\mathbf{X}_\ell^{(t+1)} \leftarrow \mathbf{X}_\ell^{(t)} - \eta\mathbf{S}_\ell^{(t)} \odot \mathbf{N}_\ell^{(t)};$
1812 **end for**
1813 **end for**

1822 C.1 ALGORITHM DESIGN

1823
1824 While low-rank subspace optimization algorithms like GaLore/GoLore project full-parameter gra-
1825 dient $\mathbf{G} \in \mathbb{R}^{(m \times n)}$ into low-rank subspaces via projection like $\mathbf{P}^\top \mathbf{G}$, sparse subspace optimization
1826 algorithms use a sparse mask \mathbf{S} to get $\mathbf{S} \odot \mathbf{G}$. Specifically, consider the following set

$$1827 \text{Sp}_{m,n}^k = \{\mathbf{S} \in \{0, 1\}^{m \times n} \mid \|\mathbf{S}\|_F^2 = k\},$$

1828
1829 *i.e.*, a set of $m \times n$ matrices contains k ones and $(mn - k)$ zeros. Corresponding to the subspace
1830 selecting strategy in GaLore, we consider a Top- k strategy which places the k ones at indices cor-
1831 responding to \mathbf{G} 's elements with the k largest absolute values. We also consider a Rand- k strategy
1832 which samples the sparse mask matrix \mathbf{S} from the uniform distribution on $\text{Sp}_{m,n}^k$ corresponding to
1833 GoLore. For convenience, we name the algorithm using Top- k strategy as GaSare (**G**radient **S**parse
1834 **p**rojection), and the one using Rand- k strategy as GoSare (**G**radient **r**andom **S**parse
1835 **p**rojection). The concerned sparse subspace descent algorithms are described as in Alg. 5

1836 C.2 NOTATIONS AND USEFUL LEMMAS
1837

1838 We assume the model parameters consist of N_L weight matrices. We use $\mathbf{X}_\ell \in \mathbb{R}^{m_\ell \times n_\ell}$ to denote
1839 the ℓ -th weight matrix and $\mathbf{x} \in \mathbb{R}^d = (\text{vec}(\mathbf{X}_1)^\top, \dots, \text{vec}(\mathbf{X}_{N_L})^\top)^\top$ to denote the vector collect-
1840 ing all the parameters, $d = \sum_{\ell=1}^{N_L} m_\ell n_\ell$. We assume GaSare/GoSare applies sparse mask in $\text{Sp}_{m_\ell, n_\ell}^{k_\ell}$
1841 to the ℓ -th weight matrix and denote

$$1842 \delta_\ell = \frac{k_\ell}{m_\ell n_\ell}, \quad \underline{\delta} = \min_{1 \leq \ell \leq N_L} \delta_\ell, \quad \bar{\delta} = \max_{1 \leq \ell \leq N_L} \delta_\ell.$$

1843 We define $\tilde{\mathbf{M}}_\ell^{(t)} = \mathbf{S}_\ell^{(t)} \odot \mathbf{M}_\ell^{(t)}$ and $\tilde{\mathbf{m}} = (\text{vec}(\tilde{\mathbf{M}}_1)^\top, \dots, \text{vec}(\tilde{\mathbf{M}}_{N_L})^\top)^\top$. While using Alg. 5
1844 with MSGD, it holds that

$$1845 \tilde{\mathbf{M}}_\ell^{(t)} = \begin{cases} \beta_1 \mathbf{S}_\ell^{(0)} \odot \mathbf{G}_\ell^{(0)}, & t = 0; \\ \mathbf{S}_\ell^{(t)} \odot \left((1 - \beta_1) \tilde{\mathbf{M}}_\ell^{(t-1)} + \beta_1 \mathbf{G}_\ell^{(t)} \right), & t = k\tau, k \in \mathbb{N}^*; \\ (1 - \beta_1) \tilde{\mathbf{M}}_\ell^{(t-1)} + \beta_1 \mathbf{S}_\ell^{(t)} \odot \mathbf{G}_\ell^{(t)}, & t = k\tau + r, k \in \mathbb{N}, 1 \leq r < \tau; \end{cases}$$

1846 and that

$$1847 \mathbf{X}_\ell^{(t+1)} = \mathbf{X}_\ell^{(t)} - \eta \tilde{\mathbf{M}}_\ell^{(t)}.$$

1848 We use $\mathbf{E}_{m,n}$ to denote the all-one $m \times n$ matrix, i.e.,

$$1849 \mathbf{E}_{m,n} = \begin{pmatrix} 1 & 1 & \cdots & 1 \\ 1 & 1 & \cdots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \cdots & 1 \end{pmatrix} \in \mathbb{R}^{m \times n}.$$

1850 **Lemma 12** (Error of GaSare’s projection). *Let \mathbf{S} be the Top- k mask of $\mathbf{G} \in \mathbb{R}^{m \times n}$, it holds that*

$$1851 \|\mathbf{S} \odot \mathbf{G} - \mathbf{G}\|_F^2 \leq \left(1 - \frac{k}{mn}\right) \|\mathbf{G}\|_F^2.$$

1852 *Proof.* Let g_1, g_2, \dots, g_{mn} be elements of \mathbf{G} such that $|g_1| \geq |g_2| \geq \dots \geq |g_{mn}|$. It holds that

$$1853 \begin{aligned} 1854 \|\mathbf{S} \odot \mathbf{G} - \mathbf{G}\|_F^2 &= \sum_{i=1}^k (g_k - g_i)^2 + \sum_{i=k+1}^{mn} (0 - g_i)^2 \\ 1855 &= \sum_{i=k+1}^{mn} g_i^2 \\ 1856 &\leq \left(1 - \frac{k}{mn}\right) \sum_{i=1}^{mn} g_i^2 \\ 1857 &= \left(1 - \frac{k}{mn}\right) \|\mathbf{G}\|_F^2, \end{aligned}$$

1858 where the inequality uses $\frac{1}{mn-k} \sum_{i=k+1}^{mn} g_i^2 \leq \frac{1}{k} \sum_{i=1}^k g_i^2$. □

1859 **Lemma 13** (Error of GoSare’s projection). *Let $\mathbf{S} \sim \mathcal{U}(\text{Sp}_{m,n}^k)$, it holds for all $\mathbf{G} \in \mathbb{R}^{m \times n}$ that*

$$1860 \mathbb{E}[\mathbf{S}] = \frac{k}{mn} \cdot \mathbf{E}_{m,n}, \tag{78}$$

1861 and

$$1862 \mathbb{E}[\|\mathbf{S} \odot \mathbf{G} - \mathbf{G}\|_F^2] = \left(1 - \frac{k}{mn}\right) \|\mathbf{G}\|_F^2. \tag{79}$$

1890 *Proof.* To prove (78), it suffices to note that for any element $S_{i,j}$ in \mathcal{S} , it holds that

$$1891 \mathbb{E}[S_{i,j}] = \mathbb{P}[S_{i,j} = 1] = \frac{(mn-1)!/[(mn-k)!(k-1)!]}{(mn)!/[(mn-k)!k!]} = \frac{k}{mn}.$$

1894 To prove (79), we have

$$1895 \mathbb{E}[\|\mathcal{S} \odot \mathbf{G} - \mathbf{G}\|_F^2] = \sum_{1 \leq i \leq m, 1 \leq j \leq n} \mathbb{P}[S_{i,j} = 0] \mathbf{G}_{i,j}^2 = \left(1 - \frac{k}{mn}\right) \|\mathbf{G}\|_F^2.$$

□

1900 C.3 NON-CONVERGENCE OF GASARE

1902 In this subsection, we present the non-convergence result of GaSare, similar to that of GaLore.

1903 **Theorem 9** (Non-convergence of GaSare). *There exists an objective function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ satisfying*
 1904 *Assumptions 1, 2, a stochastic gradient oracle (F, \mathcal{D}) satisfying Assumption 3, an initial point $\mathbf{x}^{(0)}$,*
 1905 *a constant $\epsilon_0 > 0$ such that for GaSare with any sparsity level $k_\ell < m_\ell n_\ell$, subspace changing*
 1906 *frequency τ and any subspace optimizer ρ with arbitrary hyperparameters and any $t > 0$, it holds*
 1907 *that*

$$1908 \|\nabla f(\mathbf{x}^{(t)})\|_2^2 \geq \epsilon_0.$$

1910 *Proof.* Consider target function $f(\mathbf{X}) = \frac{L}{2} \|(\mathbf{p}\mathbf{p}^\top) \odot \mathbf{X}\|_F^2$ where $L > 0$, $\mathbf{X} \in \mathbb{R}^{n \times n}$ with $n > 1$
 1911 and $\mathbf{p} = (1, 0, \dots, 0)^\top \in \mathbb{R}^n$. It holds that

$$1912 f(\mathbf{X}) = \frac{LX_{1,1}^2}{2} \geq 0,$$

1915 thus f satisfies Assumption 1. Since $\nabla f(\mathbf{X}) = L(\mathbf{p}\mathbf{p}^\top) \odot \mathbf{X}$, it holds that

$$1916 \|\nabla f(\mathbf{X}) - \nabla f(\mathbf{Y})\|_F = L\|(\mathbf{p}\mathbf{p}^\top) \odot (\mathbf{X} - \mathbf{Y})\|_F \leq L\|\mathbf{X} - \mathbf{Y}\|_F,$$

1917 thus f satisfies Assumption 2.

1919 Consider the following stochastic gradient oracle:

$$1920 F(\mathbf{X}; \xi) = f(\mathbf{X}) + \xi \tilde{\sigma} \cdot \text{tr}(\mathbf{Q}\mathbf{X}), \quad \text{and} \quad \mathbb{P}_{\xi \sim \mathcal{D}}[\xi = 1] = \mathbb{P}_{\xi \sim \mathcal{D}}[\xi = -1] = 0.5,$$

1922 where $\tilde{\sigma} = \sigma / \sqrt{n^2(n^2 - 1)/2}$ and

$$1923 \mathbf{Q} = \begin{pmatrix} 0 & \sqrt{n} & \cdots & \sqrt{n^2 - n} \\ \sqrt{1} & \sqrt{n+1} & \cdots & \sqrt{n^2 - n + 1} \\ \vdots & \vdots & \ddots & \vdots \\ \sqrt{n-1} & \sqrt{2n-1} & \cdots & \sqrt{n^2 - 1} \end{pmatrix} \in \mathbb{R}^{n \times n}.$$

1928 Note that $\nabla F(\mathbf{X}; \xi) = \nabla f(\mathbf{X}) + \xi \tilde{\sigma} \mathbf{Q}$, it holds for any $\mathbf{X} \in \mathbb{R}^{n \times n}$ that

$$1929 \mathbb{E}_{\xi \sim \mathcal{D}}[\nabla F(\mathbf{X}; \xi)] = \nabla f(\mathbf{X})$$

$$1930 \mathbb{E}_{\xi \sim \mathcal{D}}[\|\nabla F(\mathbf{X}; \xi) - \nabla f(\mathbf{X})\|_F^2] = \tilde{\sigma}^2 \|\mathbf{Q}\|_F^2 = \frac{\sigma^2}{n^2(n^2 - 1)/2} \cdot \sum_{i=1}^{n^2-1} i = \sigma^2,$$

1934 thus oracle (F, \mathcal{D}) satisfies Assumption 3.

1936 Consider the initial point $\mathbf{X}^{(0)}$ with $X_{1,1}^{(0)} = \lambda$, where $0 < \lambda < \tilde{\sigma}/L$ is a scalar. We show that
 1937 GaSare with the above objective function f , stochastic gradient oracle (F, \mathcal{D}) , initial point $\mathbf{X}^{(0)}$,
 1938 arbitrary sparsity level $0 < k < n^2$, arbitrary subspace changing frequency τ and arbitrary subspace
 1939 optimizer ρ , can only output points $\mathbf{X}^{(t)}$ with $\|\nabla f(\mathbf{X}^{(t)})\|_F^2 \geq \epsilon_0$ for $\epsilon_0 = L^2 \lambda^2 > 0$.

1941 When $\tau \mid t$, GaSare recomputes the sparse mask matrix at iteration t . If $X_{1,1}^{(t)} = \lambda$, the stochastic
 1942 gradient is given by

$$1943 \mathbf{G}^{(t)} = L(\mathbf{p}\mathbf{p}^\top) \odot \mathbf{X} + \xi^{(t)} \tilde{\sigma} \mathbf{Q}.$$

since $L\lambda < \bar{\sigma}$, the Top- k mask $\mathbf{S} \in \mathbb{R}^{n \times n}$ satisfies

$$\text{vec}(\mathbf{S}) = \underbrace{(0, 0, \dots, 0)}_{(n^2-k) \times} \underbrace{(1, 1, \dots, 1)}_{k \times} \in \mathbb{R}^{n^2},$$

Using this mask matrix, the subspace updates in the following τ iterations is as

$$\mathbf{X}^{(t+\Delta_t)} = \mathbf{X}^{(t)} + \mathbf{S}^{(t)} \odot \left(\sum_{s=0}^{\Delta_t-1} \rho^{(t+s)} (\mathbf{S}^{(t)} \odot \mathbf{G}^{(t)}) \right) \Rightarrow X_{1,1}^{(t+\Delta_t)} = X_{1,1}^{(t)} = \lambda,$$

for $\Delta_t = 1, 2, \dots, \tau$. Since $X_{1,1}^{(0)} = \lambda$, it holds for all $t > 0$ that $\mathbf{X}_{1,1}^{(t)} = \lambda$ and thus

$$\|\nabla f(\mathbf{X}^{(t)})\|_F^2 = L^2 \lambda^2 = \epsilon_0.$$

□

C.4 CONVERGENCE OF DETERMINISTIC GASARE

In this subsection, we prove the convergence properties of GaSare with deterministic gradients. The results and proofs are similar to those of deterministic GaLore in Appendix B.3.

Lemma 14 (Momentum contraction). *In deterministic GaSare using MSGD (Alg. 5), if $0 < \beta_1 \leq 1$, term $\tilde{\mathbf{M}}_\ell^{(t)}$ has the following contraction properties:*

- When $t = 0$, it holds that

$$\begin{aligned} \|\tilde{\mathbf{M}}_\ell^{(0)} - \nabla_\ell f(\mathbf{X}^{(0)})\|_F^2 &\leq (\tau - 1)(1 - \delta_\ell \beta_1) \sum_{r=0}^{\tau-2} \|\nabla_\ell f(\mathbf{x}^{(r+1)}) - \nabla_\ell f(\mathbf{x}^{(r)})\|_F^2 \\ &\quad + \frac{2(1 - \delta_\ell \beta_1)}{\tau} \sum_{r=0}^{\tau-1} \|\nabla_\ell f(\mathbf{x}^{(r)})\|_F^2; \end{aligned} \quad (80)$$

- When $t = k\tau$, $k \in \mathbb{N}^*$, it holds that

$$\begin{aligned} &\|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2 - \left(1 - \left(1 - \frac{\delta_\ell}{4}\right) \beta_1\right) \|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2 \\ &\leq \frac{2(1 - \delta_\ell)}{\tau} \sum_{r=0}^{\tau-1} \|\nabla_\ell f(\mathbf{x}^{(k\tau+r)})\|_F^2 + \frac{5(1 - \beta_1)}{\delta_\ell \beta_1} \|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2 \\ &\quad + (\tau - 1)(1 - \delta_\ell) \sum_{r=0}^{\tau-2} \|\nabla_\ell f(\mathbf{x}^{(k\tau+r+1)}) - \nabla_\ell f(\mathbf{x}^{(k\tau+r)})\|_F^2; \end{aligned} \quad (81)$$

- When $t = k\tau + r$, $k \in \mathbb{N}$, $1 \leq r < \tau$, it holds that

$$\begin{aligned} &\mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] - \left(1 - \left(1 - \frac{\delta_\ell}{4}\right) \beta_1\right) \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2] \\ &\leq \left(1 - \frac{\delta_\ell}{2}\right) \beta_1 \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + \frac{5(1 - \beta_1)}{\delta_\ell \beta_1} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2] \\ &\quad + \frac{10r\beta_1}{\delta_\ell} \sum_{i=1}^r \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(k\tau+i)}) - \nabla_\ell f(\mathbf{x}^{(k\tau+i-1)})\|_F^2]. \end{aligned} \quad (82)$$

Proof. For convenience we use \mathbf{E} to denote $\mathbf{E}_{m_\ell, n_\ell}$. When $t = 0$, we have

$$\begin{aligned} \|\tilde{\mathbf{M}}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2 &= \|\beta_1 (\mathbf{S}_\ell^{(0)} - \mathbf{E}) \odot \nabla_\ell f(\mathbf{x}^{(0)}) - (1 - \beta_1) \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2 \\ &\leq \beta_1 (1 - \delta_\ell) \|\nabla_\ell f(\mathbf{x}^{(0)})\|_F^2 + (1 - \beta_1) \|\nabla_\ell f(\mathbf{x}^{(0)})\|_F^2 \\ &= (1 - \delta_\ell \beta_1) \|\nabla_\ell f(\mathbf{x}^{(0)})\|_F^2, \end{aligned} \quad (83)$$

where the inequality uses Lemma 12 and Jensen's inequality. Applying Lemma 2 to (83) yields (80).

When $t = k\tau$, $k \in \mathbb{N}^*$, we have

$$\begin{aligned}
& \|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2 \\
&= \|\mathbf{S}_\ell^{(t)} \odot [(1 - \beta_1)\tilde{\mathbf{M}}_\ell^{(t-1)} + \beta_1 \mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})] - (\mathbf{E} - \mathbf{S}_\ell^{(t)}) \odot \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2 \\
&= \|\mathbf{S}_\ell^{(t)} \odot [(1 - \beta_1)(\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)}))]\|_F^2 + \|(\mathbf{E} - \mathbf{S}_\ell^{(t)}) \odot \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2 \\
&\leq \|(1 - \beta_1)(\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)}))\|_F^2 + (1 - \delta_\ell) \|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2, \tag{84}
\end{aligned}$$

where the inequality uses Lemma 12. By Young's inequality, we have

$$\begin{aligned}
& \|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2 \\
&= \|(\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t-1)})) - (\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(t-1)}))\|_F^2 \\
&\leq \left(1 + \frac{\delta_\ell \beta_1}{4}\right) \|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2 + \left(1 + \frac{4}{\delta_\ell \beta_1}\right) \|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2. \tag{85}
\end{aligned}$$

Applying Lemma 2 and (85) to (84) yields (81).

When $t = k\tau + r$, $k \in \mathbb{N}$, $1 \leq r < \tau$, we have

$$\begin{aligned}
& \|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2 \\
&= \|(1 - \beta_1)(\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)})) + \beta_1(\mathbf{S}_\ell^{(t)} - \mathbf{E}) \odot \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2 \\
&\leq (1 - \beta_1) \|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2 + \beta_1 \|(\mathbf{E} - \mathbf{S}_\ell^{(k\tau)}) \odot \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2, \tag{86}
\end{aligned}$$

where the inequality uses Jensen's inequality and $\mathbf{S}_\ell^{(t)} = \mathbf{S}_\ell^{(t-1)} = \dots = \mathbf{S}_\ell^{(k\tau)}$. The first term can be similarly upper bounded as (85). For the second term, we have

$$\begin{aligned}
& \|(\mathbf{E} - \mathbf{S}_\ell^{(k\tau)}) \odot \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2 \\
&\leq \left(1 + \frac{\delta_\ell}{4}\right) \|(\mathbf{E} - \mathbf{S}_\ell^{(k\tau)}) \odot \nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2 \\
&\quad + \left(1 + \frac{4}{\delta_\ell}\right) \|(\mathbf{E} - \mathbf{S}_\ell^{(k\tau)}) \odot (\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(k\tau)}))\|_F^2 \\
&\leq \left(1 + \frac{\delta_\ell}{4}\right) (1 - \delta_\ell) \|\nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2 + \frac{5}{\delta_\ell} \|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2, \tag{87}
\end{aligned}$$

where the first inequality uses Young's inequality and the second inequality uses Lemma 12. By Young's inequality, we have

$$\|\nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2 \leq \left(1 + \frac{\delta_\ell}{4}\right) \|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2 + \left(1 + \frac{4}{\delta_\ell}\right) \|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2. \tag{88}$$

Note that $t = k\tau + r$, we further have

$$\begin{aligned}
\|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2 &= \left\| \sum_{i=1}^r \nabla_\ell f(\mathbf{x}^{(k\tau+i)}) - \nabla_\ell f(\mathbf{x}^{(k\tau+i-1)}) \right\|_F^2 \\
&\leq r \sum_{i=1}^r \|\nabla_\ell f(\mathbf{x}^{(k\tau+i)}) - \nabla_\ell f(\mathbf{x}^{(k\tau+i-1)})\|_F^2, \tag{89}
\end{aligned}$$

where the inequality uses Cauchy's inequality. Applying (88)(89) to (87) yields

$$\begin{aligned}
& \|(\mathbf{E} - \mathbf{S}_\ell^{(k\tau)}) \odot \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2 \\
&\leq \left(1 - \frac{\delta_\ell}{2}\right) \|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2 + \frac{10r}{\delta_\ell} \sum_{i=1}^r \|\nabla_\ell f(\mathbf{x}^{(k\tau+i)}) - \nabla_\ell f(\mathbf{x}^{(k\tau+i-1)})\|_F^2. \tag{90}
\end{aligned}$$

Applying (85)(90) to (86) yields (82). \square

Based on Lemma 14, we can prove the convergence properties of deterministic GaSare similarly as the proofs of Lemma 7, Theorem 6 and Corollary 1. Below we directly present the final convergence results.

Theorem 10 (Convergence of deterministic GaSare). *Under Assumptions 1-2, if hyperparameters*

$$0 < \beta_1 \leq 1, \quad \tau \geq \frac{64}{3\beta_1\delta}, \quad 0 < \eta \leq \min \left\{ \frac{1}{4L}, \sqrt{\frac{3\delta\beta_1^2}{80L^2}}, \sqrt{\frac{3\delta}{80\tau^2L^2}}, \sqrt{\frac{3\beta_1}{16\tau L^2}} \right\},$$

GaSare using deterministic gradients and MSGD (Alg. 5) converges as

$$\frac{1}{K\tau} \sum_{t=0}^{K\tau-1} \|\nabla f(\mathbf{x}^{(t)})\|_2^2 \leq \frac{16\Delta}{\delta\eta K\tau}$$

for any $K \geq 1$, where $\Delta = f(\mathbf{x}^{(0)}) - \inf_{\mathbf{x}} f(\mathbf{x})$. If $T \geq 64/(3\delta)$ and we further choose

$$\begin{aligned} \beta_1 &= 1 \\ \tau &= \left\lceil \frac{64}{3\delta\beta_1} \right\rceil \\ \eta &= \left(4L + \sqrt{\frac{80L^2}{3\delta\beta_1^2}} + \sqrt{\frac{80\tau^2L^2}{3\delta}} + \sqrt{\frac{16\tau L^2}{3\beta_1}} \right)^{-1}, \end{aligned}$$

GaSare using deterministic gradients and MSGD (Alg. 5) converges as

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}^{(t)})\|_2^2 = \mathcal{O} \left(\frac{L\Delta}{\delta^{5/2}T} \right).$$

Consequently, the computation complexity to reach an ε -accurate solution \mathbf{x} such that $\|\nabla f(\mathbf{x})\|_2^2 \leq \varepsilon$ is $\mathcal{O} \left(\frac{L\Delta}{\delta^{5/2}\varepsilon} + \frac{1}{\delta} \right)$.

C.5 CONVERGENCE OF LARGE-BATCH GASARE

In this subsection, we present the convergence properties of GaSare with large-batch stochastic gradients. The results and proofs are similar to those of large-batch GaLore in Appendix B.4.

Lemma 15 (Momentum contraction). *Under Assumption 3, in large-batch GaSare using MSGD (Alg. 5), if $0 < \beta_1 \leq 1$, term $\tilde{\mathbf{M}}_\ell^{(t)}$ has the following contraction properties:*

- When $t = 0$, it holds that

$$\begin{aligned} \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(0)} - \nabla_\ell f(\mathbf{X}^{(0)})\|_F^2] &\leq 2(\tau - 1)(1 - \delta_\ell\beta_1) \sum_{r=0}^{\tau-2} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(r+1)}) - \nabla_\ell f(\mathbf{x}^{(r)})\|_F^2] \\ &\quad + \frac{4(1 - \delta_\ell\beta_1)}{\tau} \sum_{r=0}^{\tau-1} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(r)})\|_F^2] + \frac{4\beta_1\sigma_\ell^2}{\mathcal{B}}; \end{aligned} \quad (91)$$

- When $t = k\tau$, $k \in \mathbb{N}^*$, it holds that

$$\begin{aligned} &\mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] - \left(1 - \left(1 - \frac{\delta_\ell}{4} \right) \beta_1 \right) \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2] \\ &\leq \frac{4(1 - \delta_\ell)}{\tau} \sum_{r=0}^{\tau-1} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(k\tau+r)})\|_F^2] + \frac{5(1 - \beta_1)}{\delta_\ell\beta_1} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2] \\ &\quad + 2(\tau - 1)(1 - \delta_\ell) \sum_{r=0}^{\tau-2} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(k\tau+r+1)}) - \nabla_\ell f(\mathbf{x}^{(k\tau+r)})\|_F^2] + \frac{5\sigma_\ell^2}{\mathcal{B}}; \end{aligned} \quad (92)$$

- When $t = k\tau + r$, $k \in \mathbb{N}$, $1 \leq r < \tau$, it holds that

$$\begin{aligned}
& \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] - \left(1 - \left(1 - \frac{\delta_\ell}{4}\right)\beta_1\right) \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2] \\
& \leq \left(1 - \frac{\delta_\ell}{2}\right)\beta_1 \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + \frac{5(1-\beta_1)}{\delta_\ell\beta_1} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2] \\
& \quad + \frac{15r\beta_1}{\delta_\ell} \sum_{i=1}^r \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(k\tau+i)}) - \nabla_\ell f(\mathbf{x}^{(k\tau+i-1)})\|_F^2] + \left(\frac{11\beta_1}{\delta_\ell\mathcal{B}} + \beta_1^2\right)\sigma_\ell^2. \quad (93)
\end{aligned}$$

Proof. For convenience we use \mathbf{E} to denote $\mathbf{E}_{m_\ell, n_\ell}$. When $t = 0$, we have

$$\begin{aligned}
& \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2] \\
& = \mathbb{E}[\|\beta_1 \mathbf{S}_\ell^{(0)} \odot \mathbf{G}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2] \\
& = \mathbb{E}[\|\beta_1(\mathbf{S}_\ell^{(0)} - \mathbf{E}) \odot \mathbf{G}_\ell^{(0)} + \beta_1(\mathbf{G}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)})) - (1 - \beta_1)\nabla_\ell f(\mathbf{x}^{(0)})\|_F^2] \\
& \leq \beta_1 \mathbb{E}[\|(\mathbf{S}_\ell^{(0)} - \mathbf{E}) \odot \mathbf{G}_\ell^{(0)} + \mathbf{G}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2] + (1 - \beta_1)\|\nabla_\ell f(\mathbf{x}^{(0)})\|_F^2, \quad (94)
\end{aligned}$$

where the inequality uses Jensen's inequality. For the first term we have

$$\begin{aligned}
& \mathbb{E}[\|(\mathbf{S}_\ell^{(0)} - \mathbf{E}) \odot \mathbf{G}_\ell^{(0)} + \mathbf{G}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2] \\
& \leq 2\mathbb{E}[\|(\mathbf{E} - \mathbf{S}_\ell^{(0)}) \odot \mathbf{G}_\ell^{(0)}\|_F^2] + 2\mathbb{E}[\|\mathbf{G}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2] \\
& \leq 2(1 - \delta_\ell)\mathbb{E}[\|\mathbf{G}_\ell\|_F^2] + 2\mathbb{E}[\|\mathbf{G}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2] \\
& \leq 2(1 - \delta_\ell)\|\nabla_\ell f(\mathbf{x}^{(0)})\|_F^2 + \frac{(4 - 2\delta_\ell)\sigma_\ell^2}{\mathcal{B}}, \quad (95)
\end{aligned}$$

where the first inequality uses Cauchy's inequality, the second inequality uses Lemma 12, the third inequality uses $\mathbb{E}[\|\mathbf{G}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2] \leq \sigma_\ell^2/\mathcal{B}$ (Assumption 3). Applying (95) and Lemma 2 to (94) yields (91).

When $t = k\tau$, $k \in \mathbb{N}^*$, we have

$$\begin{aligned}
& \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
& = \mathbb{E}[\|\mathbf{S}_\ell^{(t)} \odot [(1 - \beta_1)\tilde{\mathbf{M}}_\ell^{(t-1)} + \beta_1\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})] - (\mathbf{E} - \mathbf{S}_\ell^{(t)}) \odot \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
& = \mathbb{E}[\|\mathbf{S}_\ell^{(t)} \odot [(1 - \beta_1)\tilde{\mathbf{M}}_\ell^{(t-1)} + \beta_1\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})]\|_F^2] + \mathbb{E}[\|(\mathbf{E} - \mathbf{S}_\ell^{(t)}) \odot \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2]. \quad (96)
\end{aligned}$$

We further have

$$\begin{aligned}
& \mathbb{E}[\|\mathbf{S}_\ell^{(t)} \odot [(1 - \beta_1)\tilde{\mathbf{M}}_\ell^{(t-1)} + \beta_1\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})]\|_F^2] \\
& \leq \mathbb{E}[\|(1 - \beta_1)\tilde{\mathbf{M}}_\ell^{(t-1)} + \beta_1\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
& = \mathbb{E}[\|(1 - \beta_1)(\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)})) + \beta_1(\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)}))\|_F^2] \\
& \leq \mathbb{E}[\|(1 - \beta_1)(\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)}))\|_F^2] + \beta_1^2 \mathbb{E}[\|\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2], \quad (97)
\end{aligned}$$

where the last inequality uses the unbiasedness of $\mathbf{G}_\ell^{(t)}$ (Assumption 3). By Young's inequality, we have

$$\begin{aligned}
& \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
& = \mathbb{E}[\|(\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t-1)})) - (\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(t-1)}))\|_F^2] \\
& \leq \left(1 + \frac{\delta_\ell\beta_1}{4}\right) \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2] + \left(1 + \frac{4}{\delta_\ell\beta_1}\right) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2]. \quad (98)
\end{aligned}$$

2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213

Applying (98) to (97) yields

$$\begin{aligned} & \mathbb{E}[\|\mathbf{S}_\ell^{(t)} \odot [(1 - \beta_1)\tilde{\mathbf{M}}_\ell^{(t-1)} + \beta_1\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})]\|_F^2] \\ & \leq \left(1 - \left(1 - \frac{\delta_\ell}{4}\right)\beta_1\right) \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2] + \frac{\beta_1^2\sigma^2}{\mathcal{B}} \\ & \quad + \frac{5(1 - \beta_1)}{\delta_\ell\beta_1} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2]. \end{aligned} \quad (99)$$

For the second term in (96), we have

$$\begin{aligned} & \mathbb{E}[\|(\mathbf{E} - \mathbf{S}_\ell^{(t)}) \odot \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\ & \leq 2\mathbb{E}[\|(\mathbf{E} - \mathbf{S}_\ell^{(t)}) \odot \mathbf{G}_\ell^{(t)}\|_F^2] + 2\mathbb{E}[\|(\mathbf{E} - \mathbf{S}_\ell^{(t)}) \odot (\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)}))\|_F^2] \\ & \leq 2(1 - \delta_\ell)\mathbb{E}[\|\mathbf{G}_\ell^{(t)}\|_F^2] + 2\mathbb{E}[\|\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\ & \leq 2(1 - \delta_\ell)\mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + \frac{4\sigma_\ell^2}{\mathcal{B}}, \end{aligned} \quad (100)$$

where the first inequality uses Cauchy's inequality, the second inequality uses Lemma 12, the third inequality uses Assumption 3. Applying (99)(100) to (96) and using Lemma 2 yields (92).

When $t = k\tau + r$, $k \in \mathbb{N}$, $1 \leq r < \tau$, we have

$$\begin{aligned} & \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\ & = \mathbb{E}[\|(1 - \beta_1)(\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)})) + \beta_1(\mathbf{S}_\ell^{(t)} \odot \mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)}))\|_F^2] \\ & = \mathbb{E}[\|(1 - \beta_1)(\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)})) + \beta_1(\mathbf{S}_\ell^{(t)} - \mathbf{E}) \odot \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\ & \quad + \beta_1^2\mathbb{E}[\|\mathbf{S}_\ell^{(t)} \odot (\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)}))\|_F^2] \\ & \leq (1 - \beta_1)\mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + \beta_1\mathbb{E}[\|(\mathbf{E} - \mathbf{S}_\ell^{(t)}) \odot \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\ & \quad + \beta_1^2\mathbb{E}[\|\mathbf{S}_\ell^{(t)} \odot (\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)}))\|_F^2], \end{aligned} \quad (101)$$

where the second equality uses the unbiasedness of $\mathbf{G}_\ell^{(t)}$ and the independence implied by $\mathbf{S}_\ell^{(t)} = \mathbf{S}_\ell^{(t-1)}$, the inequality uses Jensen's inequality. The first term is similarly bounded as (98). For the second term, we have

$$\begin{aligned} & \mathbb{E}[\|(\mathbf{E} - \mathbf{S}_\ell^{(k\tau)}) \odot \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\ & \leq \left(1 + \frac{\delta_\ell}{4}\right) \mathbb{E}[\|(\mathbf{E} - \mathbf{S}_\ell^{(k\tau)}) \odot \mathbf{G}_\ell^{(k\tau)}\|_F^2] \\ & \quad + \left(1 + \frac{4}{\delta_\ell}\right) \mathbb{E}[\|(\mathbf{E} - \mathbf{S}_\ell^{(k\tau)}) \odot (\nabla_\ell f(\mathbf{x}^{(t)}) - \mathbf{G}_\ell^{(k\tau)})\|_F^2] \\ & \leq \left(1 - \frac{3\delta_\ell}{4}\right) \mathbb{E}[\|\mathbf{G}_\ell^{(k\tau)}\|_F^2] + 2\left(1 + \frac{4}{\delta_\ell}\right) \mathbb{E}[\|\mathbf{G}_\ell^{(k\tau)} - \nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2] \\ & \quad + 2\left(1 + \frac{4}{\delta_\ell}\right) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2], \end{aligned} \quad (102)$$

where the first inequality uses Young's inequality, the second inequality uses Lemma 12 and Cauchy's inequality. We further have

$$\begin{aligned} & \left(1 - \frac{3\delta_\ell}{4}\right) \mathbb{E}[\|\mathbf{G}_\ell^{(k\tau)}\|_F^2] + 2\left(1 + \frac{4}{\delta_\ell}\right) \mathbb{E}[\|\mathbf{G}_\ell^{(k\tau)} - \nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2] \\ & \leq \left(1 - \frac{3\delta_\ell}{4}\right) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2] + \frac{11}{\delta_\ell} \mathbb{E}[\|\mathbf{G}_\ell^{(k\tau)} - \nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2] \\ & \leq \left(1 - \frac{3\delta_\ell}{4}\right) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2] + \frac{11\sigma_\ell^2}{\delta_\ell\mathcal{B}} \\ & \leq \left(1 - \frac{\delta_\ell}{2}\right) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + \left(1 + \frac{4}{\delta_\ell}\right) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2] + \frac{11\sigma_\ell^2}{\delta_\ell\mathcal{B}}, \end{aligned} \quad (103)$$

where the first inequality uses unbiasedness of $\mathbf{G}_\ell^{(k\tau)}$, the second inequality uses Assumption 3, the third inequality uses Young's inequality.

Applying (103) to (102) and applying Cauchy's inequality yields

$$\begin{aligned} & \mathbb{E}[\|(\mathbf{E} - \mathbf{S}_\ell^{(k\tau)}) \odot \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\ & \leq \left(1 - \frac{\delta_\ell}{2}\right) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + \frac{11\sigma_\ell^2}{\delta_\ell \mathcal{B}} + \frac{15r}{\delta_\ell} \sum_{i=1}^r \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(k\tau+i)}) - \nabla_\ell f(\mathbf{x}^{(k\tau+i-1)})\|_F^2]. \end{aligned} \quad (104)$$

For the third term, we have

$$\mathbb{E}[\|\mathbf{S}_\ell^{(k\tau)} \odot (\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)}))\|_F^2] \leq \mathbb{E}[\|\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \leq \sigma_\ell^2, \quad (105)$$

where the second inequality uses Assumption 3.

Applying (98)(104)(105) to (101) yields (93). \square

Based on Lemma 15, we can prove the convergence properties of large-batch GaSare similarly as the proofs of Lemma 9, Theorem 7 and Corollary 2. Below we directly present the final convergence results.

Theorem 11 (Convergence of large-batch GaSare). *Under Assumptions 1-3, if hyperparameters*

$$0 < \beta_1 \leq 1, \quad \tau \geq \frac{128}{3\beta_1 \underline{\delta}}, \quad 0 < \eta \leq \min \left\{ \frac{1}{4L}, \sqrt{\frac{3\underline{\delta}\beta_1^2}{80L^2}}, \sqrt{\frac{\underline{\delta}}{40\tau^2 L^2}}, \sqrt{\frac{3\beta_1}{32\tau L^2}} \right\},$$

GaSare using large-batch stochastic gradients and MSGD (Alg. 5) converges as

$$\frac{1}{K\tau} \sum_{t=0}^{K\tau-1} \mathbb{E}[\|\nabla f(\mathbf{x}^{(t)})\|_2^2] \leq \frac{16\Delta}{\underline{\delta}\eta K\tau} + \left(\frac{160}{3\beta_1 \underline{\delta}\tau \mathcal{B}} + \frac{352}{3\underline{\delta}^2 \mathcal{B}} + \frac{32\beta_1}{3\underline{\delta}} \right) \sigma^2$$

for any $K \geq 1$, where $\Delta = f(\mathbf{x}^{(0)}) - \inf_{\mathbf{x}} f(\mathbf{x})$. If $T \geq 2 + 256/(3\underline{\delta}) + (256\sigma)^2/(9\sqrt{\underline{\delta}}L\Delta)$ and we further choose

$$\begin{aligned} \beta_1 &= \left(1 + \sqrt{\frac{\underline{\delta}^{3/2} \sigma^2 T}{L\Delta}} \right)^{-1}, \\ \tau &= \left\lceil \frac{128}{3\underline{\delta}\beta_1} \right\rceil, \\ \eta &= \left(4L + \sqrt{\frac{80L^2}{3\underline{\delta}\beta_1^2}} + \sqrt{\frac{40\tau^2 L^2}{\underline{\delta}}} + \sqrt{\frac{32\tau L^2}{3\beta_1}} \right)^{-1}, \\ \mathcal{B} &= \left\lceil \frac{1}{\underline{\delta}\beta_1} \right\rceil, \end{aligned}$$

GaSare using large-batch stochastic gradients and MSGD (Alg. 5) converges as

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(\mathbf{x}^{(t)})\|_2^2] = \mathcal{O} \left(\frac{L\Delta}{\underline{\delta}^{5/2} T} + \sqrt{\frac{L\Delta\sigma^2}{\underline{\delta}^{7/2} T}} \right).$$

Consequently, the computation complexity to reach an ε -accurate solution \mathbf{x} such that $\|\nabla f(\mathbf{x})\|_2^2 \leq \varepsilon$ is $\mathcal{O} \left(\frac{L\Delta\sigma^2}{\underline{\delta}^{7/2}\varepsilon^2} + \frac{L\Delta}{\underline{\delta}^{5/2}\varepsilon} + \frac{\sigma^2}{\underline{\delta}^{1/2}L\Delta} + \frac{1}{\underline{\delta}} \right)$.

C.6 CONVERGENCE OF GoSARE

In this subsection, we present the convergence properties of GoSare with small-batch stochastic gradients. The results and proofs are similar to those of GoLore in Appendix B.5.

Lemma 16 (Momentum contraction). *Under Assumption 3, in GoSare using MSGD (Alg. 5), if $0 < \beta_1 \leq 1$, term $\tilde{\mathbf{M}}_\ell^{(t)}$ has the following contraction properties:*

- When $t = 0$, it holds that

$$\begin{aligned} \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(0)} - \nabla_\ell f(\mathbf{X}^{(0)})\|_F^2] &\leq (\tau - 1)(1 - \delta_\ell \beta_1) \sum_{r=0}^{\tau-2} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(r+1)}) - \nabla_\ell f(\mathbf{x}^{(r)})\|_F^2] \\ &\quad + \frac{2(1 - \delta_\ell \beta_1)}{\tau} \sum_{r=0}^{\tau-1} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(r)})\|_F^2] + \delta_\ell \beta_1^2 \sigma_\ell^2; \end{aligned} \quad (106)$$

- When $t = k\tau$, $k \in \mathbb{N}^*$, it holds that

$$\begin{aligned} &\mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] - \delta_\ell \left(1 - \left(1 - \frac{\delta_\ell}{4}\right) \beta_1\right) \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2] \\ &\leq \frac{2(1 - \delta_\ell)}{\tau} \sum_{r=0}^{\tau-1} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(k\tau+r)})\|_F^2] + \frac{5(1 - \beta_1)}{\beta_1} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2] \\ &\quad + (\tau - 1)(1 - \delta_\ell) \sum_{r=0}^{\tau-2} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(k\tau+r+1)}) - \nabla_\ell f(\mathbf{x}^{(k\tau+r)})\|_F^2] + \delta_\ell \beta_1^2 \sigma_\ell^2; \end{aligned} \quad (107)$$

- When $t = k\tau + r$, $k \in \mathbb{N}$, $1 \leq r < \tau$, it holds that

$$\begin{aligned} &\mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] - \left(1 - \left(1 - \frac{\delta_\ell}{4}\right) \beta_1\right) \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2] \\ &\leq \left(1 - \frac{\delta_\ell}{2}\right) \beta_1 \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + \frac{5(1 - \beta_1)}{\delta_\ell \beta_1} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2] \\ &\quad + \frac{10r\beta_1}{\delta_\ell} \sum_{i=1}^r \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(k\tau+i)}) - \nabla_\ell f(\mathbf{x}^{(k\tau+i-1)})\|_F^2] + \beta_1^2 \sigma_\ell^2. \end{aligned} \quad (108)$$

Proof. For convenience we use \mathbf{E} to denote $\mathbf{E}_{m_\ell, n_\ell}$. When $t = 0$, we have

$$\begin{aligned} &\mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2] \\ &= \mathbb{E}[\|\beta_1 \mathbf{S}_\ell^{(0)} \odot \mathbf{G}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2] \\ &= \mathbb{E}[\|(\beta_1 \mathbf{S}_\ell^{(0)} - \mathbf{E}) \odot \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2] + \beta_1^2 \mathbb{E}[\|\mathbf{S}_\ell^{(0)} \odot (\mathbf{G}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)}))\|_F^2], \end{aligned} \quad (109)$$

where the second equality uses unbiasedness of $\mathbf{G}_\ell^{(0)}$. By Lemma 5 we have

$$\begin{aligned} &\mathbb{E}[\|(\beta_1 \mathbf{S}_\ell^{(0)} - \mathbf{E}) \odot \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2] \\ &= \sum_{1 \leq i \leq m_\ell, 1 \leq j \leq n_\ell} \mathbb{E}[(\beta_1 [S_\ell^{(0)}]_{i,j} - 1)^2] [\nabla_\ell f(\mathbf{x}^{(0)})]_{i,j}^2 \\ &= \sum_{1 \leq i \leq m_\ell, 1 \leq j \leq n_\ell} (1 - 2\beta_1 \delta_\ell + \beta_1^2 \delta_\ell) [\nabla_\ell f(\mathbf{x}^{(0)})]_{i,j}^2 \\ &\leq (1 - \delta_\ell \beta_1) \|\nabla_\ell f(\mathbf{x}^{(0)})\|_F^2. \end{aligned} \quad (110)$$

Similarly, by Lemma 5 we have

$$\begin{aligned} &\mathbb{E}[\|\mathbf{S}_\ell^{(0)} \odot (\mathbf{G}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)}))\|_F^2] \\ &= \sum_{1 \leq i \leq m_\ell, 1 \leq j \leq n_\ell} \mathbb{E}[[S_\ell^{(0)}]_{i,j}^2] [\mathbf{G}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)})]_{i,j}^2 \\ &= \delta_\ell \mathbb{E}[\|\mathbf{G}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2] \\ &\leq \delta_\ell \sigma_\ell^2, \end{aligned} \quad (111)$$

2322 where the inequality uses Assumption 3. Applying (110)(111) and Lemma 2 to (109) yields (106).

2323 When $t = k\tau$, $k \in \mathbb{N}^*$, we have

$$\begin{aligned}
2324 & \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
2325 & = \mathbb{E}[\|\mathbf{S}_\ell^{(t)} \odot [(1 - \beta_1)\tilde{\mathbf{M}}_\ell^{(t-1)} + \beta_1 \mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})] - (\mathbf{E} - \mathbf{S}_\ell^{(t)}) \odot \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
2326 & = \delta_\ell \mathbb{E}[\|(1 - \beta_1)\tilde{\mathbf{M}}_\ell^{(t-1)} + \beta_1 \mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + (1 - \delta_\ell) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2], \quad (112)
\end{aligned}$$

2330 where the second equality uses Lemma 13. For the first term, we have

$$\begin{aligned}
2331 & \mathbb{E}[\|(1 - \beta_1)\tilde{\mathbf{M}}_\ell^{(t-1)} + \beta_1 \mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
2332 & = \mathbb{E}[\|(1 - \beta_1)(\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)})) + \beta_1(\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)}))\|_F^2] \\
2333 & \leq \mathbb{E}[\|(1 - \beta_1)(\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)}))\|_F^2] + \beta_1^2 \mathbb{E}[\|\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
2334 & \leq (1 - \beta_1) \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + \beta_1^2 \sigma_\ell^2, \quad (113)
\end{aligned}$$

2337 where both inequalities use Assumption 3. By Young's inequality, we have

$$\begin{aligned}
2338 & \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
2339 & = \mathbb{E}[\|(\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t-1)})) - (\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(t-1)}))\|_F^2] \\
2340 & \leq \left(1 + \frac{\delta_\ell \beta_1}{4}\right) \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2] + \left(1 + \frac{4}{\delta_\ell \beta_1}\right) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2]. \quad (114)
\end{aligned}$$

2346 Applying (113)(114) and Lemma 2 to (112) yields (107).

2347 When $t = k\tau + r$, $k \in \mathbb{N}$, $1 \leq r < \tau$, we have

$$\begin{aligned}
2348 & \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
2349 & = \mathbb{E}[\|(1 - \beta_1)(\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)})) + \beta_1(\mathbf{S}_\ell^{(t)} \odot \mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)}))\|_F^2] \\
2350 & = \mathbb{E}[\|(1 - \beta_1)(\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)})) + \beta_1(\mathbf{S}_\ell^{(t)} - \mathbf{E}) \odot \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
2351 & \quad + \beta_1^2 \mathbb{E}[\|\mathbf{S}_\ell^{(t)} \odot (\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)}))\|_F^2] \\
2352 & \leq (1 - \beta_1) \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + \beta_1 \mathbb{E}[\|(\mathbf{E} - \mathbf{S}_\ell^{(t)}) \odot \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
2353 & \quad + \beta_1^2 \mathbb{E}[\|\mathbf{S}_\ell^{(t)} \odot (\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)}))\|_F^2], \quad (115)
\end{aligned}$$

2358 where the second equality uses the unbiasedness of $\mathbf{G}_\ell^{(t)}$ and the independence implied by $\mathbf{S}_\ell^{(t)} =$
2359 $\mathbf{S}_\ell^{(t-1)}$, the inequality uses Jensen's inequality. The first term is similarly bounded as (114). For the
2360 second term, we have

$$\begin{aligned}
2361 & \mathbb{E}[\|(\mathbf{E} - \mathbf{S}_\ell^{(k\tau)}) \odot \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
2362 & \leq \left(1 + \frac{\delta_\ell}{4}\right) \mathbb{E}[\|(\mathbf{E} - \mathbf{S}_\ell^{(k\tau)}) \odot \nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2] \\
2363 & \quad + \left(1 + \frac{4}{\delta_\ell}\right) \mathbb{E}[\|(\mathbf{E} - \mathbf{S}_\ell^{(k\tau)}) \odot (\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(k\tau)}))\|_F^2] \\
2364 & \leq \left(1 - \frac{3\delta_\ell}{4}\right) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2] + \left(1 + \frac{4}{\delta_\ell}\right) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2], \quad (116)
\end{aligned}$$

2370 where the first inequality uses Young's inequality, the second inequality uses Lemma 13. By Young's
2371 inequality, we have

$$\begin{aligned}
2372 & \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2] \leq \left(1 + \frac{\delta_\ell}{4}\right) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + \left(1 + \frac{4}{\delta_\ell}\right) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2]. \quad (117) \\
2373 & \\
2374 & \\
2375 &
\end{aligned}$$

Applying (117) to (116) and applying Cauchy’s inequality yields

$$\begin{aligned} & \mathbb{E}[\|(\mathbf{E} - \mathbf{S}_\ell^{(k\tau)}) \odot \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\ & \leq \left(1 - \frac{\delta_\ell}{2}\right) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + \frac{10r}{\delta_\ell} \sum_{i=1}^r \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(k\tau+i)}) - \nabla_\ell f(\mathbf{x}^{(k\tau+i-1)})\|_F^2]. \end{aligned} \quad (118)$$

For the third term, we have

$$\mathbb{E}[\|\mathbf{S}_\ell^{(k\tau)} \odot (\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)}))\|_F^2] \leq \mathbb{E}[\|\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \leq \sigma_\ell^2, \quad (119)$$

where the second inequality uses Assumption 3.

Applying (114)(118)(119) to (115) yields (108). \square

Based on Lemma 16, we can prove the convergence properties of GoSare similarly as the proofs of Lemma 11, Theorem 8 and Corollary 3. Below we directly present the final convergence results.

Theorem 12 (Convergence of GoSare). *Under Assumptions 1-3, if hyperparameters*

$$0 < \beta_1 \leq 1, \quad \tau \geq \frac{64}{3\beta_1\delta}, \quad 0 < \eta \leq \min \left\{ \frac{1}{4L}, \sqrt{\frac{3\delta\beta_1^2}{80L^2}}, \sqrt{\frac{3\delta}{80\tau^2L^2}}, \sqrt{\frac{3\beta_1}{16\tau L^2}} \right\},$$

GoSare using small-batch stochastic gradients and MSGD (Alg. 5) converges as

$$\frac{1}{K\tau} \sum_{t=0}^{K\tau-1} \mathbb{E}[\|\nabla f(\mathbf{x}^{(t)})\|_2^2] \leq \frac{16\Delta}{\delta\eta K\tau} + \frac{32\beta_1\sigma^2}{3\delta}$$

for any $K \geq 1$, where $\Delta = f(\mathbf{x}^{(0)}) - \inf_{\mathbf{x}} f(\mathbf{x})$. If $T \geq 2 + 128/(3\delta) + (128\sigma)^2/(9\sqrt{\delta}L\Delta)$ and we further choose

$$\begin{aligned} \beta_1 &= \left(1 + \sqrt{\frac{\delta^{3/2}\sigma^2 T}{L\Delta}}\right)^{-1}, \\ \tau &= \left\lceil \frac{64}{3\delta\beta_1} \right\rceil, \\ \eta &= \left(4L + \sqrt{\frac{80L^2}{3\delta\beta_1^2}} + \sqrt{\frac{80\tau^2L^2}{3\delta}} + \sqrt{\frac{16\tau L^2}{3\beta_1}}\right)^{-1}, \end{aligned}$$

GoSare using small-batch stochastic gradients and MSGD (Alg. 5) converges as

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(\mathbf{x}^{(t)})\|_2^2] = \mathcal{O} \left(\frac{L\Delta}{\delta^{5/2}T} + \sqrt{\frac{L\Delta\sigma^2}{\delta^{7/2}T}} \right).$$

Consequently, the computation complexity to reach an ε -accurate solution \mathbf{x} such that $\|\nabla f(\mathbf{x})\|_2 \leq \varepsilon$ is $\mathcal{O} \left(\frac{L\Delta\sigma^2}{\delta^{7/2}\varepsilon^2} + \frac{L\Delta}{\delta^{5/2}\varepsilon} + \frac{\sigma^2}{\delta^{1/2}L\Delta} + \frac{1}{\delta} \right)$.

D THE ReLoRA-LIKE IMPLEMENTATION

An equivalent, ReLoRA-like implementation of Alg. 1 is as illustrated in Alg. 6, where we only present the case with small-batch stochastic gradients for convenience. In fact, applying ReLoRA with a fixed \mathbf{A} or \mathbf{B} is not our contribution, as it has already been used in several previous works (Hao et al., 2024; Loeschke et al., 2024). While leading to the same results, this ReLoRA-like implementation (Alg. 6) can potentially save computation as it computes the subspace gradient directly without computing the full-parameter one. Consider the case where $m \leq n$ and we use MSGD and a batch size of b . The computation complexity of GaLore’s original implementation is $2bmn$ for forward propagation, $4bmn$ for backward propagation, $4rmn$ for projection, $3rn$ for momentum update and $2mn$ for weight update. The computational complexity of our ReLoRA-like implementation is $2bmn + 2brm + 2brn$ for forward propagation, $2bmn + 2brm + 2brn$ for backward propagation, $3rn$ for momentum updates and $2rn$ for weight updates. As illustrated in Table 1, our implementation can potentially reduce computation with little memory overhead.

2430
2431
2432
2433
2434
2435
2436
2437
2438
2439
2440
2441
2442
2443
2444
2445
2446
2447
2448
2449
2450
2451
2452
2453
2454
2455
2456
2457
2458
2459
2460
2461
2462
2463
2464
2465
2466
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
2483

Algorithm 6 ReLoRA-like implementation of **GaLore** / **GoLore** algorithm using stochastic gradients **with** / **without** momentum projection

Input: Initial point $\mathbf{x}^{(0)}$, data distribution \mathcal{D} , learning rate η , subspace changing frequency τ , rank $\{r_\ell\}_{\ell=1}^{N_L}$, optimizer hyperparameters $\beta_1, \beta_2, \epsilon$, large batch size \mathcal{B} .

Output: $\{\mathbf{x}^{(t)}\}_{t=0}^T$.

Initialize LoRA adaptation $\mathbf{X}_\ell = \mathbf{W}_\ell + \mathbf{B}_\ell \mathbf{A}_\ell$ for $\ell = 1, 2, \dots, N_L$, where $\mathbf{W}_\ell^{(0)} = \mathbf{X}_\ell^{(0)}$, $\mathbf{A}_\ell^{(0)} = 0$ and $\mathbf{B}_\ell^{(0)} = 0$;

Initialize optimizer state $\{\mathbf{M}_\ell^{(-1)}\}_{\ell=1}^{N_L}$ and $\{\mathbf{V}_\ell^{(-1)}\}_{\ell=1}^{N_L}$ to zero;

for $t = 0, 1, \dots, T - 1$ **do**

for $\ell = 1, 2, \dots, N_L$ **do**

if $t \equiv 0 \pmod{\tau}$ **then**

$\mathbf{G}_\ell^{(t)} \leftarrow \nabla_\ell F(\mathbf{x}^{(t)}; \xi^{(t)});$

$\mathbf{U}, \Sigma, \mathbf{V} \leftarrow \text{SVD}(\mathbf{G}_\ell^{(t)}), \mathbf{P}_\ell^{(t)} \leftarrow \mathbf{U}[:, :r_\ell], \mathbf{Q}_\ell^{(t)} \leftarrow \mathbf{V}[:, :r_\ell];$ (GaLore)

 Sample $\mathbf{P}_\ell^{(t)} \sim \mathcal{U}(\text{St}_{m_\ell, r_\ell}), \mathbf{Q}_\ell^{(t)} \sim \mathcal{U}(\text{St}_{n_\ell, r_\ell});$ (GoLore)

$\mathbf{R}_\ell^{(t)} \leftarrow \begin{cases} (\mathbf{P}_\ell^{(t)})^\top \mathbf{G}_\ell^{(t)}, & \text{if } m_\ell \leq n_\ell; \\ \mathbf{G}_\ell^{(t)} \mathbf{Q}_\ell^{(t)}, & \text{if } m_\ell > n_\ell; \end{cases}$

else

$\mathbf{R}_\ell^{(t)} \leftarrow \begin{cases} \nabla_{\mathbf{A}_\ell} F(\mathbf{x}^{(t)}; \xi^{(t)}), & \text{if } m_\ell \leq n_\ell; \\ \nabla_{\mathbf{B}_\ell} F(\mathbf{x}^{(t)}; \xi^{(t)}), & \text{if } m_\ell > n_\ell; \end{cases}$

end if

$\mathbf{M}_\ell^{(t)} \leftarrow \begin{cases} (1 - \beta_1)(\mathbf{P}_\ell^{(t)})^\top \mathbf{B}_\ell^{(t)} \mathbf{M}_\ell^{(t-1)} + \beta_1 \mathbf{R}_\ell^{(t)}, & \text{if } m_\ell \leq n_\ell; \\ (1 - \beta_1) \mathbf{M}_\ell^{(t-1)} \mathbf{A}_\ell^{(t)} \mathbf{Q}_\ell^{(t)} + \beta_1 \mathbf{R}_\ell^{(t)}, & \text{if } m_\ell > n_\ell; \end{cases}$ (with MP)

$\mathbf{M}_\ell^{(t)} \leftarrow (1 - \beta_1) \mathbf{M}_\ell^{(t-1)} + \beta_1 \mathbf{R}_\ell^{(t)};$ (without MP)

$\mathbf{V}_\ell^{(t)} \leftarrow (1 - \beta_2) \mathbf{V}_\ell^{(t-1)} + \beta_2 \mathbf{R}_\ell^{(t)} \odot \mathbf{R}_\ell^{(t)};$

if using Adam **then**

$\mathbf{M}_\ell^{(t)} \leftarrow \mathbf{M}_\ell^{(t)} / (1 - \beta_1^t), \mathbf{V}_\ell^{(t)} \leftarrow \mathbf{V}_\ell^{(t)} / (1 - \beta_2^t), \mathbf{N}_\ell^{(t)} \leftarrow \mathbf{M}_\ell^{(t)} / (\sqrt{\mathbf{V}_\ell^{(t)}} + \epsilon);$

else if using MSGD **then**

$\mathbf{N}_\ell^{(t)} \leftarrow \mathbf{M}_\ell^{(t)};$

end if

if $t \equiv 0 \pmod{\tau}$ **then**

$\mathbf{W}_\ell^{(t+1)} \leftarrow \mathbf{W}_\ell^{(t)} + \mathbf{B}_\ell^{(t)} \mathbf{A}_\ell^{(t)};$

$\mathbf{A}_\ell^{(t+1)} \leftarrow \begin{cases} -\eta \mathbf{N}_\ell^{(t)}, & \text{if } m_\ell \leq n_\ell; \\ (\mathbf{Q}_\ell^{(t)})^\top, & \text{if } m_\ell > n_\ell; \end{cases}$

$\mathbf{B}_\ell^{(t+1)} \leftarrow \begin{cases} \mathbf{P}_\ell^{(t)}, & \text{if } m_\ell \leq n_\ell; \\ -\eta \mathbf{N}_\ell^{(t)}, & \text{if } m_\ell > n_\ell; \end{cases}$

else

$\mathbf{W}_\ell^{(t+1)} \leftarrow \mathbf{W}_\ell^{(t)};$

$\mathbf{A}_\ell^{(t+1)} \leftarrow \begin{cases} \mathbf{A}_\ell^{(t)} - \eta \mathbf{N}_\ell^{(t)}, & \text{if } m_\ell \leq n_\ell; \\ \mathbf{A}_\ell^{(t)}, & \text{if } m_\ell > n_\ell; \end{cases}$

$\mathbf{B}_\ell^{(t+1)} \leftarrow \begin{cases} \mathbf{B}_\ell^{(t)}, & \text{if } m_\ell \leq n_\ell; \\ \mathbf{B}_\ell^{(t)} - \eta \mathbf{N}_\ell^{(t)}, & \text{if } m_\ell > n_\ell; \end{cases}$

end if

end for

end for

E EXPERIMENTAL SPECIFICATIONS

In this section, we elaborate on the missing details concerned with the experiments we present in Sec. 7.

GaLore’s non-convergence. We compared Galore, large-batch GaLore, GoLore and full-parameter training on the constructed quadratic problem defined in (1). We used a batch size of 128 for large-batch GaLore and a batch size of 1 for the others.

Pre-training tasks on C4 dataset. We pre-trained LLaMA-60M on C4 dataset for 10,000 iterations on 4 NVIDIA A100 40G GPUs. We use batch size 128, learning rate 1.0e-3, rank 128, scaling factor $\alpha = 1$, subspace changing frequency $\tau = 200$, and a max sequence length of 256. Results under 8-bit training are shown in Fig. 6. Fig. 7 presents the results of different algorithms after trained on more tokens.

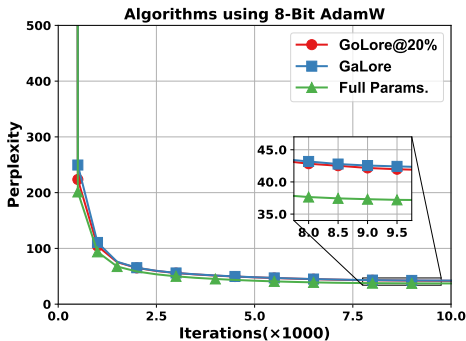


Figure 6: Pre-training curves of algorithms using 8-bit AdamW.

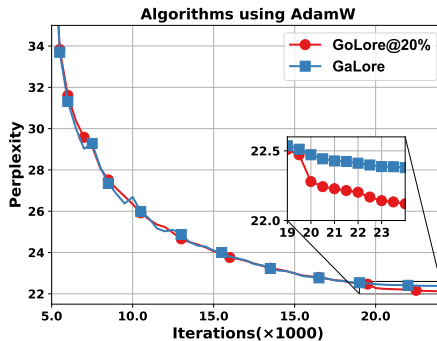


Figure 7: Pre-training curves of algorithms using AdamW.

Fine-tuning tasks on WinoGrande dataset. We fine-tune pre-trained LLaMA2-7B model on the WinoGrande dataset for 30 epochs on 4 NVIDIA A100 80G GPUs. We use batch size 1, rank 1024, subspaces changing frequency $\tau = 500$ and a max sequence length of 2048. The learning rate and scaling factor are set as 1.0e-4 and $\alpha = 4$ for GaLore/GoLore, thus corresponds to a learning rate of 4.0e-4 in full-parameter fine-tuning. The test accuracy is presented in Table 3, where GoLore performs similarly to GaLore due to overfitting.

Fine-tuning tasks on BoolQ dataset. We fine-tune pre-trained LLaMA2-7B model on the BoolQ (Clark et al., 2019) dataset on 4 NVIDIA A100 80G GPUs. We use batch size 1, rank 1024, subspaces changing frequency $\tau = 500$ and a max sequence length of 2048. We use MSGD as the subspace optimizer, where the learning rate and scaling factor are set as 1.0e-4 and $\alpha = 4$ for GaLore/GoLore, corresponding to a learning rate of 4.0e-4 in full-parameter fine-tuning. Table 3 presents the test accuracy of different algorithms, where GoLore outperforms GaLore. We further fine-tune pre-trained OPT-13B (Zhang et al., 2022) for 1 epoch using the same experimental setup, whose results are shown in Table 4.

Table 3: Evaluating GaLore/GoLore for fine-tuning on WinoGrande and BoolQ using pre-trained LLaMA2-7B.

Algorithm	BoolQ (1 epoch)	BoolQ (3 epochs)	WinoGrande (80 epochs)
Full Params.	86.48	87.43	69.85
GaLore	84.89	86.79	68.51
GoLore@20%	85.81	86.88	68.51

Table 4: Results for fine-tuning pre-trained OPT-13B models on BoolQ. *OOM* stands for "out of memory".

Algorithm	Memory	Accuracy
Full Params.	OOM	-
GaLore	77.68 GB	79.79
GoLore@30%	77.68 GB	81.96

Fine-tuning tasks on GLUE benchmark. We fine-tune pre-trained RoBERTa-Base model on the GLUE benchmark for 30 epochs on a single GeForce RTX 4090. Training details including batch size, learning rate, rank, scaling factor α and max sequence length are illustrated in Table 5.

Table 5: Hyperparameters used in fine-tuning pre-trained RoBERTa-Base model on the GLUE benchmark.

Hyperparameter	CoLA	STS-B	MRPC	RTE	SST2	MNLI	QNLI	QQP
batch size	32	16	16	16	16	16	16	16
Learning Rate	2.5e-5	2.0e-5	3.5e-5	7.0e-6	1.0e-5	1.0e-5	1.0e-5	1.0e-5
Rank	4	4	4	4	4	4	4	4
GaLore’s α	4	4	4	4	4	4	4	4
FLORA’s α	4	4	4	4	4	4	4	4
GoLore’s α	4	4	4	4	4	4	4	4
Frequency τ	500	500	500	500	500	500	500	500
Max Seq. Len.	512	512	512	512	512	512	512	512

F CONNECTIONS WITH OTHER ALGORITHMS

Connection with zero-th order methods. Zero-th order methods (Malladi et al., 2023; Zhang et al., 2023; Chen et al., 2024) are another line of works on memory-efficient training. While these algorithms randomly select a direction to estimate the directional derivatives by finite difference, GoLore computes subspace gradients via back propagation. The directions used in zero-th order methods change every iteration, while GoLore applies a more lazily strategy changing its subspace every τ iterations.

Connection with gradient sketching methods. Gradient sketching methods like Hanzely et al. (2018) and Wang et al. (2024) uses gradient sketches in algorithm iterates. These methods recover gradient estimates from projected gradients and retains full-size gradients and optimizer states. In comparison, GoLore directly updates with projected gradients and retains compressed gradients and optimizer states, which is more memory-efficient.

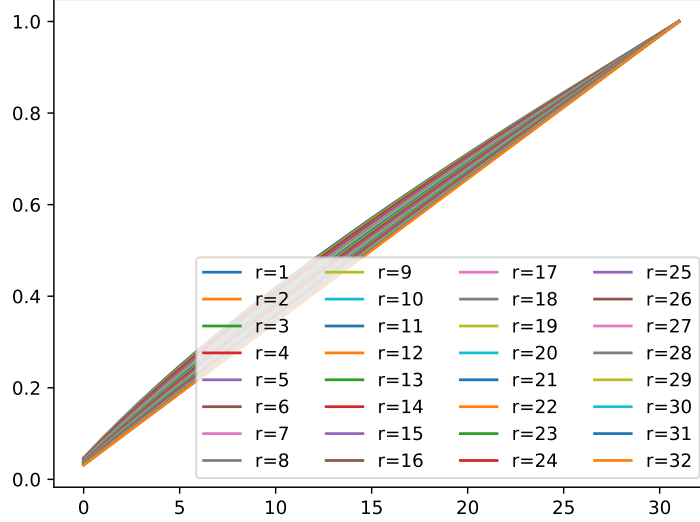
G CONVERGENCE OF GALORE UNDER ISOTROPIC NOISE ASSUMPTIONS

Based on the anisotropic gradient noise we use to construct the counter-example in the proof of GaLore’s non-convergence under standard assumptions, an interesting open question is whether GaLore is guaranteed to converge if the noise are further assumed isotropic. In this section, we consider the following additional assumption:

Assumption 4 (Isotropic noise). *The distribution of stochastic noise for each gradient matrix is invariant under orthogonal transformations, i.e., it holds for any layer $\ell = 1, \dots, N_L$, parameter $\mathbf{x} \in \mathbb{R}^d$ and orthogonal matrix $\mathbf{O}_1 \in \mathbb{R}^{m_\ell \times m_\ell}$, $\mathbf{O}_2 \in \mathbb{R}^{n_\ell \times n_\ell}$ that*

$$\nabla_\ell F(\mathbf{x}; \xi) - \nabla_\ell f(\mathbf{x}) \stackrel{\text{dist}}{=} \mathbf{O}_1 [\nabla_\ell F(\mathbf{x}; \xi) - \nabla_\ell f(\mathbf{x})] \mathbf{O}_2,$$

where $A \stackrel{\text{dist}}{=} B$ represents A and B shares the same distribution.

Figure 8: Observations with a small noise scale $\sigma = 0.1$.

Remark. The property in Assumption 4 can be satisfied by multivariate Gaussian distribution, *e.g.*, $\text{vec}(\nabla_{\ell}F(\mathbf{x}; \xi) - \nabla_{\ell}f(\mathbf{x})) \sim \mathcal{N}(0, \frac{\sigma_{\ell}^2}{m_{\ell}n_{\ell}} \cdot \mathbf{I}_{m_{\ell} \times n_{\ell}})$.

Besides Assumption 4, we consider an additional assumption, which is crucial in analyzing the projection error.

Assumption 5 (Leading property). Let $\mathcal{D}_{\ell}(\mathbf{x})$ denotes the distribution of gradient noise $\nabla_{\ell}F(\mathbf{x}; \xi) - \nabla_{\ell}f(\mathbf{x})$. We assume $\mathcal{D}_{\ell}(\mathbf{x})$ satisfies the following "leading property": if $\mathbf{A} \sim \mathcal{D}_{\ell}(\mathbf{x})$, $\mathbf{B} \in \mathbb{R}^{m_{\ell} \times n_{\ell}}$ satisfies $B_{11} \geq B_{22} \geq \dots \geq B_{\min\{m_{\ell}, n_{\ell}\}, \min\{m_{\ell}, n_{\ell}\}} \geq 0$ and $B_{ij} = 0$ for $i \neq j$, the SVD decomposition $\mathbf{U}\Sigma\mathbf{V}^{\top}$ satisfies

$$\begin{cases} \frac{1}{r} \sum_{i=1}^k \sum_{j=1}^r \mathbb{E}[U_{ij}^2] \geq \frac{k}{n}, & \forall 1 \leq k, r \leq m, & \text{if } m_{\ell} \leq n_{\ell}; \\ \frac{1}{r} \sum_{i=1}^k \sum_{j=1}^r \mathbb{E}[V_{ij}^2] \geq \frac{k}{n}, & \forall 1 \leq k, r \leq n, & \text{if } m_{\ell} > n_{\ell}. \end{cases}$$

Though not fully established in theory, we can empirically validate that multivariate Gaussian distribution may satisfy Assumption 5.

Specifically, we consider the following experiment setup. Let $\text{vec}(\mathbf{A}) \sim \mathcal{N}(0, \sigma^2 \cdot \mathbf{I}_{32 \times 32})$ for some noise scale $\sigma > 0$ and select a fixed matrix \mathbf{B} with $B_{11} \geq B_{22} \geq \dots \geq B_{32,32} \geq 0$. In order to validate the properties in expectation, we sample matrix \mathbf{A} for 200,000 times and uses the empirical expectations $\hat{\mathbb{E}}[U_{ij}]$'s to estimate the true expectations $\mathbb{E}[U_{ij}]$'s. Figures 8, 9, 10 represent results under different noise scales $\sigma = 10, 1, 0.1$, respectively, where " $r = r_0$ " in each figure plots the line connecting points $(k, \frac{1}{r_0} \sum_{i=1}^k \sum_{j=1}^{r_0} \hat{\mathbb{E}}[U_{ij}^2])$ for $k = 1, 2, \dots, 32$. As presented, all lines " $r = r_0$ " with $r_0 < 32$ are above the line " $r = 32$ ", which is guaranteed to pass through the points $(k, \frac{k}{32})$, $k = 1, 2, \dots, 32$, in theory. Consequently, we have good reason to believe that multivariate Gaussian distribution can empirically satisfy Assumption 5.

With Assumptions 4 and 5, we can establish new error bounds for GaLore's SVD projection.

Lemma 17 (Error of GaLore's projection under isotropic noise). Let $\mathbf{G} = \nabla_{\ell}f(\mathbf{x})$ and $\mathbf{E} = \nabla_{\ell}F(\mathbf{x}; \xi) - \nabla_{\ell}f(\mathbf{x})$, projection matrix $\mathbf{P} = \mathbf{U}[:, : r_{\ell}]$, $\mathbf{Q} = \mathbf{V}[:, : r_{\ell}]$ where $\mathbf{U}\Sigma\mathbf{V}^{\top} = \mathbf{G} + \mathbf{E}$ is the SVD of stochastic gradient $\nabla_{\ell}F(\mathbf{x}; \xi)$, it holds under Assumptions 4 and 5 for $m_{\ell} \leq n_{\ell}$ that

$$\mathbb{E}[\|\mathbf{P}\mathbf{P}^{\top}\mathbf{G} - \mathbf{G}\|_F^2] \leq \left(1 - \frac{r_{\ell}}{m_{\ell}}\right) \|\mathbf{G}\|_F^2,$$

2646
 2647
 2648
 2649
 2650
 2651
 2652
 2653
 2654
 2655
 2656
 2657
 2658
 2659
 2660
 2661
 2662
 2663
 2664
 2665
 2666
 2667
 2668
 2669
 2670
 2671
 2672
 2673
 2674
 2675
 2676
 2677
 2678
 2679
 2680
 2681
 2682
 2683
 2684
 2685
 2686
 2687
 2688
 2689
 2690
 2691
 2692
 2693
 2694
 2695
 2696
 2697
 2698
 2699

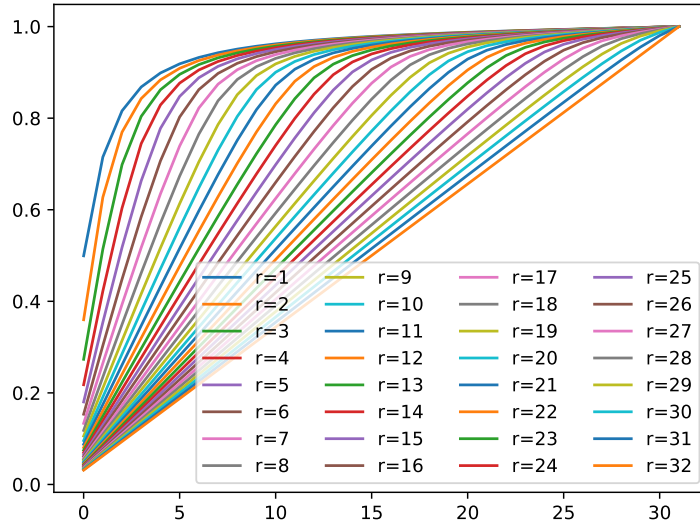


Figure 9: Observations with a medium noise scale $\sigma = 1$.

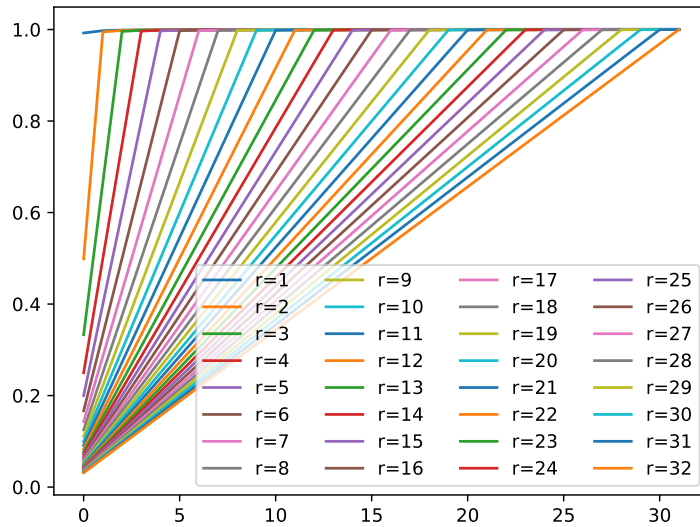


Figure 10: Observations with a large noise scale $\sigma = 10$.

2700 and for $m_\ell > n_\ell$ that
 2701

$$2702 \mathbb{E}[\|\mathbf{G}\mathbf{Q}\mathbf{Q}^\top - \mathbf{G}\|_F^2] \leq \left(1 - \frac{r_\ell}{n_\ell}\right) \|\mathbf{G}\|_F^2.$$

2703 *Proof.* We only consider the case where $m_\ell < n_\ell$, as the proof for the other case is similar. We first
 2704 conduct SVD of \mathbf{G} and get $\mathbf{G} = \mathbf{U}_0 \boldsymbol{\Sigma}_0 \mathbf{V}_0^\top$. It holds that

$$2705 \begin{aligned} 2706 \|\mathbf{P}\mathbf{P}^\top \mathbf{G}\|_F^2 &= \text{tr}(\mathbf{G}^\top \mathbf{P}\mathbf{P}^\top \mathbf{G}) \\ 2707 &= \text{tr}(\mathbf{V}_0 \boldsymbol{\Sigma}_0^\top \mathbf{U}_0^\top \mathbf{P}\mathbf{P}^\top \mathbf{U}_0 \boldsymbol{\Sigma}_0 \mathbf{V}_0^\top) \\ 2708 &= \text{tr}(\boldsymbol{\Sigma}_0 \boldsymbol{\Sigma}_0^\top \mathbf{U}_0^\top \mathbf{P}\mathbf{P}^\top \mathbf{U}_0). \end{aligned} \quad (120)$$

2709 Denote $\tilde{\mathbf{U}} = \mathbf{U}_0^\top \mathbf{U}$ and $\tilde{\mathbf{V}} = \mathbf{V}_0^\top \mathbf{V}$, it holds that $\tilde{\mathbf{U}} \boldsymbol{\Sigma}_0 \tilde{\mathbf{V}}^\top = (\mathbf{U}_0^\top \mathbf{U}) \boldsymbol{\Sigma}_0 (\mathbf{V}_0^\top \mathbf{V})^\top$ is SVD of
 2710 $\mathbf{U}_0^\top (\mathbf{G} + \mathbf{E}) \mathbf{V}_0 = \mathbf{U}_0^\top \mathbf{E} \mathbf{V}_0 + \boldsymbol{\Sigma}_0 \stackrel{\text{dist}}{=} \mathbf{E} + \boldsymbol{\Sigma}_0$. By Assumption 5 we have

$$2711 \frac{1}{r_\ell} \sum_{i=1}^k \sum_{j=1}^{r_\ell} \mathbb{E}[\tilde{U}_{ij}^2] \geq \frac{k}{m_\ell}, \quad k = 1, 2, \dots, m_\ell. \quad (121)$$

2712 Let $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{m_\ell} \geq 0$ represent the singular values of \mathbf{G} , taking expectations of (120) yields

$$2713 \begin{aligned} 2714 \mathbb{E}[\|\mathbf{P}\mathbf{P}^\top \mathbf{G}\|_F^2] &= \text{tr}(\boldsymbol{\Sigma}_0 \boldsymbol{\Sigma}_0^\top \mathbb{E}[\mathbf{U}_0^\top \mathbf{P}\mathbf{P}^\top \mathbf{U}_0]) \\ 2715 &= \sum_{i=1}^{m_\ell} \sigma_i^2 \sum_{j=1}^{r_\ell} \mathbb{E}[\tilde{U}_{ij}^2] \\ 2716 &\geq \sum_{i=1}^{m_\ell} \sigma_i^2 \cdot \frac{r_\ell}{m_\ell} = \frac{r_\ell}{m_\ell} \cdot \|\mathbf{G}\|_F^2, \end{aligned} \quad (122)$$

2717 where the inequality applies $\sigma_1^2 \geq \sigma_2^2 \geq \dots \geq \sigma_{m_\ell}^2$ and (121). Based on (122), we have

$$2718 \mathbb{E}[\|\mathbf{P}\mathbf{P}^\top \mathbf{G} - \mathbf{G}\|_F^2] = \|\mathbf{G}\|_F^2 - \mathbb{E}[\|\mathbf{P}\mathbf{P}^\top \mathbf{G}\|_F^2] \leq \left(1 - \frac{r_\ell}{m_\ell}\right) \|\mathbf{G}\|_F^2,$$

2719 which completes the proof. \square

2720 **Lemma 18** (Momentum contraction). *Under Assumption 3-5, in GaLore using MSGD with MP, if*
 2721 $0 < \beta_1 \leq 1$, term $\tilde{\mathbf{M}}_\ell^{(t)}$ has the following contraction properties:

- 2722 • When $t = 0$, it holds that

$$2723 \begin{aligned} 2724 \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(0)} - \nabla_\ell f(\mathbf{X}^{(0)})\|_F^2] &\leq (\tau - 1)(2 - \delta_\ell) \sum_{r=0}^{\tau-2} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(r+1)}) - \nabla_\ell f(\mathbf{x}^{(r)})\|_F^2] \\ 2725 &\quad + \frac{2(2 - \delta_\ell)}{\tau} \sum_{r=0}^{\tau-1} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(r)})\|_F^2] + \beta_1^2 \sigma_\ell^2; \end{aligned} \quad (123)$$

- 2726 • When $t = k\tau$, $k \in \mathbb{N}^*$, it holds that

$$2727 \begin{aligned} 2728 \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] &- \left(1 - \left(1 - \frac{\delta_\ell}{4}\right) \beta_1\right) \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2] \\ 2729 &\leq \frac{2(1 - \delta_\ell)}{\tau} \sum_{r=0}^{\tau-1} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(k\tau+r)})\|_F^2] + \frac{5(1 - \beta_1)}{\delta_\ell \beta_1} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2] \\ 2730 &\quad + (\tau - 1)(1 - \delta_\ell) \sum_{r=0}^{\tau-2} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(k\tau+r+1)}) - \nabla_\ell f(\mathbf{x}^{(k\tau+r)})\|_F^2] + \beta_1^2 \sigma_\ell^2; \end{aligned} \quad (124)$$

- When $t = k\tau + r$, $k \in \mathbb{N}$, $1 \leq r < \tau$, it holds that

$$\begin{aligned}
& \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] - \left(1 - \left(1 - \frac{\delta_\ell}{4}\right)\beta_1\right) \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2] \\
& \leq \left(1 - \frac{\delta_\ell}{2}\right)\beta_1 \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + \frac{5(1-\beta_1)}{\delta_\ell\beta_1} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2] \\
& \quad + \frac{10r\beta_1}{\delta_\ell} \sum_{i=1}^r \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(k\tau+i)}) - \nabla_\ell f(\mathbf{x}^{(k\tau+i-1)})\|_F^2] + \beta_1^2 \sigma_\ell^2. \tag{125}
\end{aligned}$$

Proof. Without loss of generality assume $m_\ell \leq n_\ell$ (the other case can be proved similarly). When $t = 0$, (123) is direct result of Lemma 8 by letting $\mathcal{B} = 1$. When $t = 0$, we have

$$\begin{aligned}
& \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2] \\
& = \mathbb{E}[\|\beta_1 \mathbf{P}_\ell^{(0)} (\mathbf{P}_\ell^{(0)})^\top \mathbf{G}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2] \\
& = \mathbb{E}[\|(\beta_1 \mathbf{P}_\ell^{(0)} (\mathbf{P}_\ell^{(0)})^\top - \mathbf{I}) \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2] + \beta_1^2 \mathbb{E}[\|\mathbf{P}_\ell^{(0)} (\mathbf{P}_\ell^{(0)})^\top (\mathbf{G}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)}))\|_F^2], \tag{126}
\end{aligned}$$

For the first term, we have

$$\begin{aligned}
& \mathbb{E}[\|(\beta_1 \mathbf{P}_\ell^{(0)} (\mathbf{P}_\ell^{(0)})^\top - \mathbf{I}) \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2] \\
& = (1 - \beta_1)^2 \mathbb{E}[\|\mathbf{P}_\ell^{(0)} (\mathbf{P}_\ell^{(0)})^\top \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2] + \mathbb{E}[\|(\mathbf{I} - \mathbf{P}_\ell^{(0)} (\mathbf{P}_\ell^{(0)})^\top) \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2] \\
& \leq ((1 - \beta_1)^2 + (1 - \delta_\ell)) \|\nabla_\ell f(\mathbf{x}^{(0)})\|_F^2 \leq (2 - \delta_\ell) \|\nabla_\ell f(\mathbf{x}^{(0)})\|_F^2, \tag{127}
\end{aligned}$$

where the first inequality uses Lemma 17. For the second term, we have

$$\mathbb{E}[\|\mathbf{P}_\ell^{(0)} (\mathbf{P}_\ell^{(0)})^\top (\mathbf{G}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)}))\|_F^2] \leq \mathbb{E}[\|\mathbf{G}_\ell^{(0)} - \nabla_\ell f(\mathbf{x}^{(0)})\|_F^2] \leq \sigma_\ell^2. \tag{128}$$

Applying (127)(128) to (126) and using Lemma 2 yields (123).

When $t = k\tau$, $k \in \mathbb{N}^*$, according to the proof of Lemma 8, we have

$$\begin{aligned}
& \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
& \leq \left(1 + \frac{\delta_\ell\beta_1}{4}\right) \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2] + \left(1 + \frac{4}{\delta_\ell\beta_1}\right) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(t-1)})\|_F^2], \tag{129}
\end{aligned}$$

and

$$\begin{aligned}
& \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
& = \mathbb{E}[\|\mathbf{P}_\ell^{(t)} (\mathbf{P}_\ell^{(t)})^\top [(1 - \beta_1)\tilde{\mathbf{M}}_\ell^{(t-1)} + \beta_1 \mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})]\|_F^2] \\
& \quad + \mathbb{E}[\|(\mathbf{I} - \mathbf{P}_\ell^{(t)} (\mathbf{P}_\ell^{(t)})^\top) \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
& \leq \mathbb{E}[\|(1 - \beta_1)(\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)}))\|_F^2] + \beta_1^2 \sigma_\ell^2 + (1 - \delta_\ell) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2], \tag{130}
\end{aligned}$$

where the last inequality applies Lemma 17. Applying (129) to (130) and using Lemma 2 yields (124).

When $t = k\tau + r$, $k \in \mathbb{N}$, $1 \leq r < \tau$, we have the following results according to the proof of Lemma 8:

$$\begin{aligned}
& \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
& \leq (1 - \beta_1) \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + \beta_1 \mathbb{E}[\|(\mathbf{I} - \mathbf{P}_\ell^{(t)} (\mathbf{P}_\ell^{(t)})^\top) \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
& \quad + \beta_1^2 \mathbb{E}[\|\mathbf{P}_\ell^{(t)} (\mathbf{P}_\ell^{(t)})^\top (\mathbf{G}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)}))\|_F^2] \\
& \leq (1 - \beta_1) \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t-1)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + \beta_1 \mathbb{E}[\|(\mathbf{I} - \mathbf{P}_\ell^{(t)} (\mathbf{P}_\ell^{(t)})^\top) \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + \beta_1^2 \sigma_\ell^2, \tag{131}
\end{aligned}$$

For the second term, we have

$$\begin{aligned}
& \mathbb{E}[\|(I - \mathbf{P}_\ell^{(k\tau)})(\mathbf{P}_\ell^{(k\tau)})^\top \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
& \leq \left(1 + \frac{\delta_\ell}{4}\right) \mathbb{E}[\|(I - \mathbf{P}_\ell^{(k\tau)})(\mathbf{P}_\ell^{(k\tau)})^\top \nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2] \\
& \quad + \left(1 + \frac{4}{\delta_\ell}\right) \mathbb{E}[\|(I - \mathbf{P}_\ell^{(k\tau)})(\mathbf{P}_\ell^{(k\tau)})^\top (\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(k\tau)}))\|_F^2] \\
& \leq \left(1 - \frac{3\delta_\ell}{4}\right) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2] + \left(1 + \frac{4}{\delta_\ell}\right) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2] \\
& \leq \left(1 - \frac{\delta_\ell}{2}\right) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + 2 \left(1 + \frac{4}{\delta_\ell}\right) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)}) - \nabla_\ell f(\mathbf{x}^{(k\tau)})\|_F^2] \\
& \leq \left(1 - \frac{\delta_\ell}{2}\right) \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + \frac{10r}{\delta_\ell} \sum_{i=1}^r \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(k\tau+i)}) - \nabla_\ell f(\mathbf{x}^{(k\tau+i-1)})\|_F^2], \quad (132)
\end{aligned}$$

where the first inequality applies Young's inequality, the second inequality applies Lemma 17, the third inequality applies Young's inequality, the last inequality applies Cauchy's inequality. Applying (129)(132) to (131) yields (125). \square

Lemma 19 (Momentum error). *Under Assumption 2-5, if $0 < \beta_1 \leq 1$ in GaLore using MSGD and MP, it holds for any $K \geq 1$ that*

$$\begin{aligned}
& \sum_{t=0}^{K\tau-1} \mathbb{E}[\|\tilde{\mathbf{m}}^{(t)} - \nabla f(\mathbf{x}^{(t)})\|_2^2] \\
& \leq \left(\frac{5(1-\beta_1)}{(1-\bar{\delta}/4)\bar{\delta}\beta_1^2} + \frac{5\tau(\tau-1)}{(1-\bar{\delta}/4)\bar{\delta}} + \frac{2(\tau-1)}{(1-\bar{\delta}/4)\beta_1}\right) L^2 \sum_{t=0}^{K\tau-2} \mathbb{E}[\|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_2^2] \\
& \quad + \left(\frac{1-\bar{\delta}/2}{1-\bar{\delta}/4} + \frac{4}{(1-\bar{\delta}/4)\tau\beta_1}\right) \sum_{t=0}^{K\tau-1} \mathbb{E}[\|\nabla f(\mathbf{x}^{(t)})\|_2^2] + \frac{K\tau\beta_1\sigma^2}{1-\bar{\delta}/4}. \quad (133)
\end{aligned}$$

Proof. By Lemma 18 we have

$$\begin{aligned}
& \sum_{t=0}^{K\tau-1} \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] - \left(1 - \left(1 - \frac{\delta_\ell}{4}\right)\beta_1\right) \sum_{t=0}^{K\tau-2} \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
& \leq \left(\frac{5(1-\beta_1)}{\delta_\ell\beta_1} + \frac{5\tau(\tau-1)\beta_1}{\delta_\ell} + 2(\tau-1)\right) \sum_{t=0}^{K\tau-2} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t+1)}) - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
& \quad + \left(\frac{4}{\tau} + \left(1 - \frac{\delta_\ell}{4}\right)\beta_1\right) \sum_{t=0}^{K\tau-1} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + K\tau\beta_1^2\sigma_\ell^2,
\end{aligned}$$

which implies

$$\begin{aligned}
& \sum_{t=0}^{K\tau-1} \mathbb{E}[\|\tilde{\mathbf{M}}_\ell^{(t)} - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
& \leq \left(\frac{5(1-\beta_1)}{(1-\delta_\ell/4)\delta_\ell\beta_1^2} + \frac{5\tau(\tau-1)}{(1-\delta_\ell/4)\delta_\ell} + \frac{2(\tau-1)}{(1-\delta_\ell/4)\beta_1}\right) \sum_{t=0}^{K\tau-2} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t+1)}) - \nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] \\
& \quad + \left(\frac{1-\delta_\ell/2}{1-\delta_\ell/4} + \frac{4}{(1-\delta_\ell/4)\tau\beta_1}\right) \sum_{t=0}^{K\tau-1} \mathbb{E}[\|\nabla_\ell f(\mathbf{x}^{(t)})\|_F^2] + \frac{K\tau\beta_1\sigma_\ell^2}{1-\delta_\ell/4}. \quad (134)
\end{aligned}$$

Summing (134) for $\ell = 1, \dots, N_L$ and applying Assumption 2-3 yields (133). \square

Now we are ready to prove the convergence of GaLore with small-batch stochastic gradients under isotropic noise assumptions.

Theorem 13 (Convergence of Galore under isotropic noise assumptions). *Under Assumptions 1-5, if hyperparameters*

$$0 < \beta_1 \leq 1, \quad \tau \geq \frac{128}{3\beta_1\bar{\delta}}, \quad 0 < \eta \leq \min \left\{ \frac{1}{4L}, \sqrt{\frac{3\delta\beta_1^2}{80L^2}}, \sqrt{\frac{3\delta}{80\tau^2L^2}}, \sqrt{\frac{3\beta_1}{32\tau L^2}} \right\}, \quad (135)$$

GaLore using small-batch stochastic gradients and MSGD with MP converges as

$$\frac{1}{K\tau} \sum_{t=0}^{K\tau-1} \mathbb{E} \|\nabla f(\mathbf{x}^{(t)})\|_2^2 \leq \frac{16\Delta}{\bar{\delta}\eta K\tau} + \frac{32\beta_1\sigma^2}{3\bar{\delta}} \quad (136)$$

for any $K \geq 1$, where $\Delta = f(\mathbf{x}^{(0)}) - \inf_{\mathbf{x}} f(\mathbf{x})$.

Proof. By Lemma 4 we have

$$\begin{aligned} \sum_{t=0}^{K\tau-1} \mathbb{E} \|\nabla f(\mathbf{x}^{(t)})\|_2^2 &\leq \frac{2[f(\mathbf{x}^{(0)}) - \mathbb{E}[f(\mathbf{x}^{(K\tau)})]]}{\eta} + \sum_{t=0}^{K\tau-1} \mathbb{E} \|\tilde{\mathbf{m}}^{(t)} - \nabla f(\mathbf{x}^{(t)})\|_2^2 \\ &\quad - \left(\frac{1}{\eta^2} - \frac{L}{\eta} \right) \sum_{t=0}^{K\tau-1} \mathbb{E} \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_2^2. \end{aligned} \quad (137)$$

Applying Lemma 19 to (137) and using $\bar{\delta} \leq \delta < 1$ yields

$$\begin{aligned} &\left(\frac{\delta}{4} - \frac{16}{3\tau\beta_1} \right) \sum_{t=0}^{K\tau-1} \mathbb{E} \|\nabla f(\mathbf{x}^{(t)})\|_2^2 \\ &\leq \frac{2}{\eta} \mathbb{E}[f(\mathbf{x}^{(0)}) - f(\mathbf{x}^{(K\tau)})] + \frac{4K\tau\beta_1\sigma^2}{3} \\ &\quad - \left(\frac{1}{\eta^2} - \frac{L}{\eta} - \frac{20(1-\beta_1)L^2}{3\delta\beta_1^2} - \frac{20\tau(\tau-1)L^2}{3\bar{\delta}} - \frac{8(\tau-1)L^2}{3\beta_1} \right) \sum_{t=0}^{K\tau-1} \mathbb{E} \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_2^2. \end{aligned} \quad (138)$$

By (135) we have

$$\frac{\delta}{4} - \frac{16}{3\tau\beta_1} \geq \frac{\delta}{8}, \quad \text{and} \quad \frac{1}{4\eta^2} \geq \max \left\{ \frac{L}{\eta}, \frac{20(1-\beta_1)L^2}{3\delta\beta_1^2}, \frac{20\tau(\tau-1)L^2}{3\bar{\delta}}, \frac{8(\tau-1)L^2}{3\beta_1} \right\}. \quad (139)$$

Applying (139) to (138) yields (136). \square

Corollary 4 (Convergence complexity of GaLore under isotropic noise assumptions). *Under Assumptions 1-5, if $T \geq 2 + 256/(3\bar{\delta}) + (256\sigma^2)/(9\sqrt{\bar{\delta}}L\Delta)$ and we choose*

$$\begin{aligned} \beta_1 &= \left(1 + \sqrt{\frac{\delta^{3/2}\sigma^2 T}{L\Delta}} \right)^{-1}, \\ \tau &= \left\lceil \frac{128}{3\bar{\delta}\beta_1} \right\rceil, \\ \eta &= \left(4L + \sqrt{\frac{80L^2}{3\bar{\delta}\beta_1^2}} + \sqrt{\frac{80\tau^2L^2}{3\bar{\delta}}} + \sqrt{\frac{32\tau L^2}{3\beta_1}} \right)^{-1}, \end{aligned}$$

GaLore using small-batch stochastic gradients and MSGD with MP converges as

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(\mathbf{x}^{(t)})\|_2^2 = \mathcal{O} \left(\frac{L\Delta}{\bar{\delta}^{5/2}T} + \sqrt{\frac{L\Delta\sigma^2}{\bar{\delta}^{7/2}T}} \right), \quad (140)$$

where $\Delta = f(\mathbf{x}^{(0)}) - \inf_{\mathbf{x}} f(\mathbf{x})$. Consequently, the computation complexity to reach an ε -accurate solution \mathbf{x} such that $\|\nabla f(\mathbf{x})\|_2^2 \leq \varepsilon$ is $\mathcal{O} \left(\frac{L\Delta\sigma^2}{\bar{\delta}^{7/2}\varepsilon^2} + \frac{L\Delta}{\bar{\delta}^{5/2}\varepsilon} + \frac{\sigma^2}{\bar{\delta}^{1/2}L\Delta} + \frac{1}{\bar{\delta}} \right)$.

2916 *Proof.* $T \geq 2 + 256/(3\underline{\delta}) + (256\sigma)^2/(9\sqrt{\underline{\delta}}L\Delta)$ guarantees $T \geq \tau$. Let $T = K\tau + r$, where $K \in \mathbb{N}^*$
 2917 and $0 \leq r < \tau$. If $r = 0$, (140) is a direct result of Theorem 13. If $r > 0$, applying Theorem 13 to
 2918 $\tilde{K} := K + 1$ yields

$$2920 \quad \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(\mathbf{x}^{(t)})\|_2^2] \leq \frac{\tilde{K}\tau}{T} \cdot \frac{1}{\tilde{K}\tau} \sum_{t=0}^{\tilde{K}\tau-1} \mathbb{E}[\|\nabla f(\mathbf{x}^{(t)})\|_2^2] = \mathcal{O} \left(\frac{L\Delta}{\underline{\delta}^{5/2}T} + \sqrt{\frac{L\Delta\sigma^2}{\underline{\delta}^{7/2}T}} \right).$$

□

2923
2924
2925
2926
2927
2928
2929
2930
2931
2932
2933
2934
2935
2936
2937
2938
2939
2940
2941
2942
2943
2944
2945
2946
2947
2948
2949
2950
2951
2952
2953
2954
2955
2956
2957
2958
2959
2960
2961
2962
2963
2964
2965
2966
2967
2968
2969