

8. Appendix

8.1. Additional Results

Empirical Evidence of BP Gradients Converging to Subspaces

Previous studies have shown that BP gradients during fine-tuning of LLMs rapidly converge to low-dimensional subspaces. Building on this, our empirical investigation using the OPT-1.3B architecture on the SST-2 dataset uncovers a persistent low-rank structure in gradient matrices throughout the optimization process. As shown in Fig. 3, singular value decomposition (SVD) of gradient matrices across layers consistently reveals pronounced spectral decay, with only a small subset of singular values dominating the gradient spectrum during LLM fine-tuning with SGD.

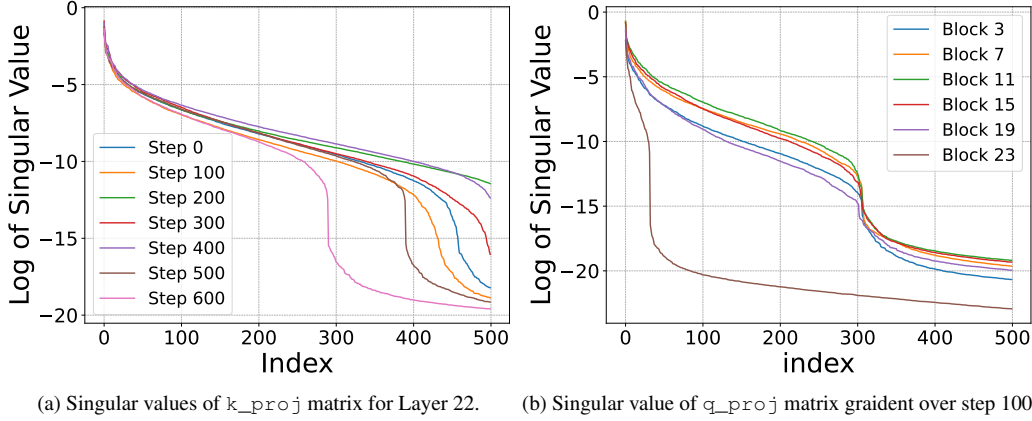


Figure 3. Low-rank structures in gradient matrices during LLM fine-tuning. (a) illustrates the temporal evolution of singular values for key projections within a single layer, showing consistent spectral decay across training steps. (b) presents the cross-layer comparison at step 100, revealing layer-invariant low-rank patterns in query projection gradients.

More Comparisons

In the main manuscript, we use the identical batch size for FO and ZO optimizers. Here, we adjust SGD with gradient accumulation to match the memory usage of ZO optimizers, and then compare their convergence speed and performance. The experimental settings are the same as those in Fig. 1, and the experimental results are shown in Fig. 4. With similar memory usage, SubZero attains a convergence rate nearly on par with SGD, surpasses MeZO, and achieves test accuracy comparable to that of SGD.

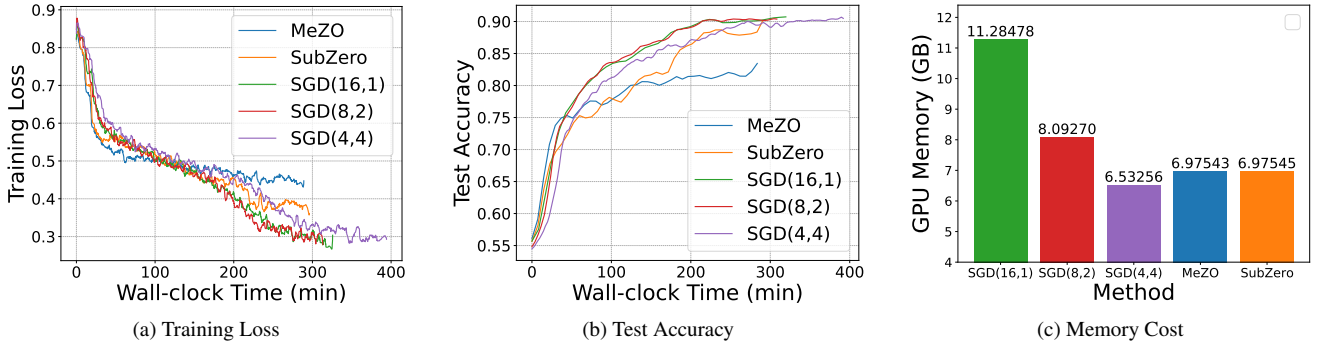


Figure 4. Visualization of training loss, test accuracy, and peak total GPU memory usage with OPT-1.3B on SST-2 in prompt tuning scheme. SGD(BS, GA) refers to SGD with a batch size of BS and GA times of gradient accumulation. All ZO methods utilize a batch size of 16, while SGD(BS, GA) applies gradient accumulation to ensure its memory usage aligns with that of the ZO optimizers. All methods are executed for 20K steps.

The vanilla LoRA is fine-tuned by Adam. We compare SubZero with SGD in the FT and LoRA schemes with vanilla LoRA using the pretrained OPT-1.3B model on SST-2. For Adam and SubZero with SGD, we apply the constant learning rate

schedule. The results are given in Table 9. We can see that SubZero with SGD in the FT scheme outperforms vanilla LoRA in both test accuracy and memory usage. SubZero with SGD in the LoRA scheme also achieves comparable test accuracy while maintaining minimal memory usage.

Table 9. Comparison with vanilla LoRA using the pretrained OPT-1.3B model on SST-2.

Method	Test Accuracy(%)	Total Memory (GB)
LoRA (Adam)	93.2	10.75
SubZero (FT)	93.4	6.88
SubZero (LoRA)	92.9	6.80

As described in Sec. 6.3, we compare the memory consumption and wall-clock time of ZO methods (MeZO and SubZero), SGD, and inference-only approaches (zero-shot and in-context learning (ICL)) using OPT-13B (see Table 10). Since inference-only methods do not involve fine-tuning, they have zero wall-clock time and their memory usage reflects only the inference load. For fine-tuning, all methods were run for 20K steps. The ZO methods, including SubZero, achieved over a 1.8 \times reduction in memory usage compared to SGD. Notably, SubZero’s memory footprint closely aligns with MeZO’s, while offering improved performance. We use per-layer weight updates for MeZO and SubZero (see Appendix 8.2), resulting in nearly identical memory usage for FT and LoRA schemes when one decimal place is reserved.

Although SubZero introduces computational and memory overhead due to QR decomposition when generating projection matrices, our empirical analysis reveals strictly bounded resource costs across all tested OPT model scales (see Table 11). Specifically, the additional time overhead remains below 8.5% (peaking at 8.36% for the 6.7B model), while the memory overhead stays under 1.8%, even under bfloat16 precision. This indicates that the computational cost of QR decomposition becomes asymptotically negligible as model complexity increases, thereby establishing SubZero’s practical scalability.

Table 10. Memory usage (GB) and wall-clock time (minutes) of fine-tuning OPT-13B, with SGD’s batch size being 8 for SQuAD and 16 for other tasks.

Task Method	SST-2		WIC		SQuAD	
	Mem.	Time	Mem.	Time	Mem.	Time
Zero-shot/ICL	24.2	0	24.8	0	27.2	0
SGD(FT)	48.9	190.3	48.9	257.3	122.7	623.7
MeZO(FT)	26.1	324.9	26.6	370.5	37.4	670.2
SubZero(FT)	26.5	337.3	27.1	385.3	37.8	690.5
MeZO(LoRA)	26.1	123.9	26.6	171.6	37.4	476.7
SubZero(LoRA)	26.1	130.3	26.6	179.7	37.4	486.5

Table 11. Peak memory usage (GB) and wall-clock time (seconds) for fine-tuning OPT models on the SST-2 dataset. All models fine-tuned on SST-2 for 20K steps. Precision strategy: FP32 for models ≤ 6.7 B, BF16 for ≥ 6.7 B.

OPT series	hidden size	MeZO Mem.	SubZero Mem.	Mem. Overhead (%)	MeZO Time	SubZero Time	Time Overhead (%)
1.3B	2048	6.80	6.88	+1.18%	11214.95	11683.39	+4.18%
2.7B	2560	12.40	12.60	+1.61%	21578.56	22335.05	+3.51%
6.7B	4096	13.98	14.20	+1.57%	9832.69	10654.59	+8.36%
13B	5120	26.08	26.53	+1.73%	18667.50	19245.10	+3.09%

Integration with Other ZO Optimizers

SubZero is orthogonal to many ZO methods and can be combined with them to further boost performance, such as the momentum mechanism in ZO-AdaMU [28], the sparsity pruning in S-MeZO [39], and the second-order information in HiZOO [68]. Here, we report the experimental results of integrating SubZero with ZO-AdaMU, as shown in Table 12 and

Fig. 5. This integration not only achieves faster convergence but also significantly reduces the memory overhead associated with ZO-AdaMU’s momentum mechanism.

Method	Memory (GB)	Accuracy (%)
MeZO	26.62	60.0
SubZero	27.06	60.8
ZO-AdaMU	50.77	60.7
SubZO-AdaMU	27.07	61.1

Table 12. Peak memory usage (GB) and test accuracy (%) for fine-tuning OPT-13B models on the WIC dataset.

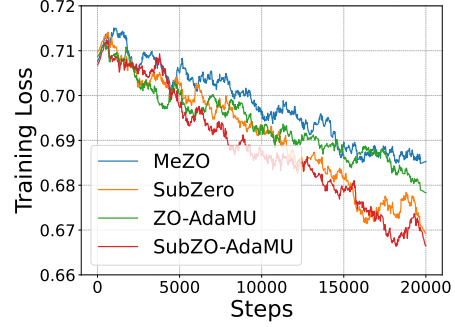


Figure 5. Training loss curves of Table 12.

ZO-AdaMU accelerates convergence by introducing stochastic momentum, a mechanism that is orthogonal to SubZero. In all experiments, we use a shared learning rate and perturbation scale. SubZero-specific hyperparameters, such as rank and update interval, are set as detailed in Table 3, while all momentum-related settings remain consistent with the original ZO-AdaMU paper [28]. Notably, vanilla ZO-AdaMU requires storing the full momentum history, leading to substantial memory overhead. In contrast, when combined with SubZero, only low-dimensional momentum $\mathbf{M}_i^t \in \mathbb{R}^{r \times r}$ needs to be maintained, i.e.,

$$\mathbf{M}_i^t = \beta_i^t \mathbf{M}_i^{t-1} + (1 - \beta_i^t) \mathbf{Z}_i^t, \quad (17)$$

$$\mathbf{W}_i^t = \mathbf{W}_i^{t-1} - \eta_i^t \mathbf{U}_i \mathbf{M}_i^t \mathbf{V}_i^T. \quad (18)$$

When the subspaces are updated at time steps t such that $t \bmod F \equiv 0$ (as specified in Algorithm 3), the old momentum needs to be projected from the old subspaces \mathbf{U}^{t-1} and \mathbf{V}^{t-1} onto the new subspaces \mathbf{U}^t and \mathbf{V}^t . This projection ensures the optimizer’s continuity and maintains the effectiveness of the momentum mechanism. The projection is formulated as the following optimization problem:

$$\min_M \|\mathbf{U}^{t-1} \mathbf{M}^{t-1} (\mathbf{V}^{t-1})^T - \mathbf{U}^t \mathbf{M} (\mathbf{V}^t)^T\|. \quad (19)$$

More Ablation Studies

We investigate the effects of random seed, batch size, and combination with Adam for SubZero. We also provide more results on the reshaping strategy.

We first fine-tune the OPT-1.3B model on the SST-2 dataset in prompt tuning scheme with three random seeds. We present the results in Table 13, and hyperparameters are presented in Table 18. For various random seeds, the variance of MeZO is quite large, whereas the variance of SubZero is small, and its average performance is superior.

Then we examine the impact of batch size for ZO optimizers using the RoBERTa-large model on SST-2 in full-parameter tuning scheme. The results are shown in Table 14. The training epochs are 100K in Table 6, while they are 20K in Table 14. The remaining hyperparameters are consistent with Table 6, as detailed in Appendix 8.4. For ZO optimizers, a large batch size always gets better performance. Across various batch sizes, SubZero demonstrates better fine-tuning performance compared to MeZO.

Table 13. The impact of random seed with the pretrained OPT-1.3B model on SST-2 in prompt tuning scheme.

Seed	42	0	1234	AVG.
MeZO	85.9	83.3	80.7	83.3
SubZero	89.1	89.4	89.2	89.2

Next, we assess the impact of the Adam optimizer. We fine-tune the OPT-1.3B model on the SST-2 dataset, and the experimental results are displayed in Table 15. For ZO optimizers with Adam, we perform a grid search on the hyperparameters

Table 14. The impact of batch size with the pretrained RoBERTa-large model in full-parameter tuning scheme.

Batch Size	Method	SST-2	SST-5	SNLI	MNLI	AVG.
16	MeZO	91.7	44.7	77.3	53.0	66.7
	SubZero	91.9	45.9	77.5	52.8	67.0
32	MeZO	92.9	45.4	78.3	53.2	67.5
	SubZero	93.0	45.5	79.6	54.0	68.0

and find that keeping the learning rate and perturbation scale consistent with those with SGD resulted in good convergence, as detailed in Table 18. We utilize the linear and the constant learning rate schedules for SGD and Adam, respectively. For all ZO optimizers with SGD, we apply the constant learning rate schedule. For all ZO optimizers with Adam, we test the constant and the cosine annealing schedules. We note that SubZero surpasses MeZO when employing the Adam optimizer with both constant and cosine annealing schedules. Also, Adam does not provide an advantage over SGD for ZO optimization, which aligns with the conclusions of previous studies [22, 65].

Table 15. Comparison of test accuracy (%) for the pretrained OPT-1.3B model fine-tuned on SST-2 with SGD and Adam.

Method	FT	LoRA	Prefix	Prompt	AVG.
SGD	93.2	93.0	93.1	90.7	92.5
Adam	92.6	93.2	92.9	93.3	93.0
MeZO_SGD	92.3	92.8	91.6	85.9	90.7
SubZero_SGD	93.4	92.9	92.2	89.1	91.9
MeZO_Adam(constant)	92.3	93.3	90.7	84.6	90.2
SubZero_Adam(constant)	93.2	92.4	90.9	89.3	91.5
MeZO_Adam(cosine)	91.9	93.1	86.1	78.7	87.5
SubZero_Adam(cosine)	91.7	92.0	86.6	83.4	88.4

Finally, we provide more ablations on the reshaping strategy with OPT-1.3B on Winogrande in the PEFT schemes. The Winogrande dataset [48] is a benchmark for commonsense reasoning and available at <https://winogrande.allenai.org/>. The results are shown in Table 16. We can see that the reshaping strategy clearly enhances performance, aligning with the conclusion presented in Table 8.

Table 16. Reshaping strategy for non-square matrices with the pretrained OPT-1.3B model fine-tuned on Winogrande in PEFT schemes.

Method	LoRA	Prefix	Prompt	AVG.
SGD	58.3	56.9	58.4	57.9
SubZero(w/o)	56.6	56.6	56.5	56.6
SubZero(w/)	57.8	57.3	57.6	57.6

8.2. Implementation Details

We use one A800 GPU with the PyTorch 2.1.0+CUDA 11.8 framework for ZO methods and, if needed, two A800 GPUs for SGD.

The gradient estimation in SubZero is applicable to parameter matrices, while LLMs mainly consist of dense layers. For other trainable parameters, such as biases and layer normalization parameters, we recommend using the gradient estimation in MeZO [40], as these layers contain fewer parameters.

We introduce two useful strategies to implement our SubZero efficiently in memory.

In-place Operation. As indicated in Eqn. (7), directly computing the loss difference ρ requires twice the memory of inference, as it must store both the parameter matrix set \mathcal{W} and the perturbation matrix set $\tilde{\mathcal{Z}}$. To mitigate this, we draw inspiration from MeZO and utilize in-place operations. By employing the random seed trick, we store a random seed to compute ρ (see lines 9-12 in Algorithm 3 and Algorithm 2) and regenerate the low-dimensional perturbation matrices $\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_l$ (see line 15 in Algorithm 3). Consequently, the memory cost for fine-tuning with SubZero is nearly equivalent to that of inference (see Table 2 and Table 10).

Per-layer Weight Update. FO optimizers update all model parameters after BP by storing the entire gradients in memory. In contrast, ZO optimizers like SubZero calculate gradient estimates by first determining the loss value difference from two forward passes, then calculating the gradient estimate for each layer using this difference along with the layer’s perturbation. To reduce memory usage during training, we can implement the parameter update with `optimizer.step()` after calculating the gradient estimate for each layer.

SubZero significantly reduces GPU memory consumption with the two implementation strategies. It should note that we use the per-layer weight update strategy for MeZO in all experiments.

To simplify hyperparameter tuning, we employ a norm alignment trick, allowing SubZero to directly utilize hyperparameter settings, such as the learning rate, from MeZO [40]. For a random perturbation matrix $\mathbf{Z} \in \mathbb{R}^{m \times n}$, and its low-rank approximation is $\hat{\mathbf{Z}} = \mathbf{U}\mathbf{Z}'\mathbf{V}^\top$, where $\mathbf{U} \in \mathbb{R}^{m \times r}$, $\mathbf{V} \in \mathbb{R}^{n \times r}$, and $\mathbf{Z}' \in \mathbb{R}^{r \times r}$. If \mathbf{Z} and \mathbf{Z}' are Gaussian random matrices, and \mathbf{U} and \mathbf{V} are column-orthogonal matrices, then we have:

$$\mathbb{E}[\|\mathbf{Z}\|_F] = \sqrt{\frac{m \times n}{r^2}} \mathbb{E}[\|\hat{\mathbf{Z}}\|_F]. \quad (20)$$

Define $\mu = \sqrt{\frac{m \times n}{r^2}}$. Let MeZO’s learning rate be η and perturbation scale be ε . There are two equivalent approaches to obtain the perturbation for SubZero. The first approach involves multiplying the random low-dimensional perturbation matrix by μ , with SubZero adopting MeZO’s hyperparameters directly: $\eta' = \eta$ and $\varepsilon' = \varepsilon$. The second approach keeps the random low-dimensional perturbation matrix fixed and sets SubZero’s learning rate and perturbation scale as follows:

$$\eta' = \eta\mu^2, \varepsilon' = \varepsilon\mu.$$

We argue that norm alignment is crucial for SubZero, as changing the rank r affects the norm of the gradient estimate, complicating the fine-tuning of the associated learning rate.

S-MeZO [39], a new ZO method, aims to improve MeZO’s performance and convergence speed. However, its source code and detailed layer-wise hyperparameter configurations have not been released. Yang et al. [60] reproduce S-MeZO using a fixed sparsity ratio for each layer, selected based on the best overall result shown in Fig. 6 of their paper. So we perform S-MeZO with this non-official implementation code available at <https://github.com/yifanycc/AdaZeta>.

8.3. Datasets

Following [40], we use SuperGLUE [57] for OPT experiments, including BoolQ [10], CB [13], COPA [47], MultiRC [29], ReCoRD [63], RTE [4, 5, 11, 21], WiC [44], and WSC [34]. We also utilize SST-2 [50] and two question answering (QA) datasets, SQuAD [45] and DROP [18]. For each task, we randomly sampled 1000 examples for training, 500 for validation, and 1000 for testing.

For LLama2-7B and Mistral-7B, we use CB [13] in the full-parameter tuning and three PEFT schemes. For OPT-1.3B, we utilize SST-2 [50] in the full-parameter tuning and three PEFT schemes.

For RoBERTa-large, we consider classification datasets: SST-2 [50], SST-5 [50], MNLI [59], and SNLI [6]. Following [40], the test set has 1000 examples for fast iteration, while we have 512 examples per class for both training and validation.

8.4. Hyperparameters

Using a larger batch size can consistently reduce the variance in ZO optimization, thus enhancing fine-tuning performance [20, 40, 60]. However, this increase in batch size also raises the time for forward passes and significantly elevates memory usage. We focus on developing ZO methods that minimize variance and improve performance with small batch sizes, with a default

Table 17. The hyperparameter search grids for OPT-13B. For each task, we run 20K steps for ZO methods (MeZO, S-MeZO, and SubZero) and SGD. We record the best model checkpoint based on the validation loss every 500 training steps.

Experiment	Hyperparameter	Value
MeZO(FT)	batch size	16
	learning rate	{ 1e-7, 2e-7, 5e-7, 1e-6 }
	ϵ	1e-3
MeZO(LoRA)	batch size	16
	learning rate	{ 1.5e-5, 3e-5, 5e-5 }
	ϵ	1e-3
S-MeZO(FT)	batch size	16
	learning rate	{ 1e-6, 5e-6 }
	ϵ	1e-3
S-MeZO(LoRA)	sparse rate	0.75
	batch size	16
	learning rate	{ 5e-5, 1e-4, 1e-3 }
SubZero(FT)	ϵ	1e-3
	rank	{ 32, 64, 128, 256 }
	subspace change frequency	{ 500, 1000, 2000 }
SubZero(LoRA)	batch size	16
	learning rate	{ 1.5e-5, 3e-5, 5e-5 }
	ϵ	1e-3
SGD(FT)	rank	{ 4, 8, 16 }
	subspace change frequency	{ 500, 1000, 2000 }
	batch size	16
SGD(FT)	learning rate	{ 1e-4, 1e-3, 5e-3 }

setting of 16. In some SGD experiments, like on MultiRC and SQuAD, the batch size is reduced to 8 due to limited GPU resources.

Consistent with previous studies [39, 40, 60, 65], we employ SGD without momentum by default to maintain memory efficiency. SGD utilizes the linear learning rate schedule, while all ZO methods with SGD apply a constant learning rate schedule, with weight decay set to 0.

For the projection matrix generation experiments summarized in Table 1, we perform full-parameter fine-tuning on the SST-2 dataset using three models: the OPT-1.3B model with a rank of 24 and a subspace update frequency of 1000, the LLaMA2-7B model with a rank of 24 and a subspace update frequency of 1000 and the OPT-13B model with a rank of 128 and a subspace update frequency of 1000. All other hyperparameters, including learning rate and perturbation scale, remain consistent across experiments. While the decomposition of the weight and Gaussian random matrices is relatively straightforward, the activation matrix relies on the AdaBK algorithm [61]. Specifically, we compute the second-order statistics of activations (see line 7 of Algorithm 1 in the paper of AdaBK) and employ an adaptive damping strategy based on the maximum singular value to improve the condition number, followed by QR decomposition. For dimensionality alignment, the second-order statistics of output and input features are decomposed to derive \mathbf{U} and \mathbf{V} , respectively. Finally, for the decomposition of the historical zeroth-order gradient, its matrix from the previous iteration is readily obtained and subsequently QR-decomposed, facilitated by our random seed technique and the saved historical projection matrix.

For RoBERTa, we run Adam for 1K steps and ZO methods for 100K steps. In the rest experiments, we run Adam for 5 epochs and SGD and ZO methods for 20K steps.

We follow previous work to set the hyperparameters in PEFT schemes [40, 65]. For LoRA, the rank is set to 8 and α is set to 16. For prefix tuning, the length of prefix tokens is set to 5, and we initialize these tunable representations by randomly sampling tokens from the vocabulary and then passing them through the LLM to get their keys and values at different attention

layers. For prompt tuning, the length of prompt virtual tokens is set to 10, and the prompt tokens are initialized with actual token values from the model’s embedding.

We present the hyperparameter search grids in Tables 17 and 18 to assist with result reproduction. For OPT-1.3B, we utilize the same hyperparameter settings as in Table 18. For RoBERTa-large, we use a learning rate of $\{1e-6, 5e-6\}$ and $\varepsilon=1e-3$ for MeZO and SubZero, with a batch size of 64. The rank for SubZero is set to $\{8, 16, 24\}$, and subspace change frequency is adjusted to $\{1000, 2000\}$.

For the ablation study, we evaluate the effectiveness of the orthogonal projection matrix using the OPT-13B model in full-parameter tuning scheme on the RTE and WSC datasets, and the results are presented in Table 5. The hyperparameter settings are consistent with those in Table 3, and further details are available in Table 17. The subspace dimensionality remains fixed across all experiments. It is noteworthy that both orthogonal and non-orthogonal projection matrices can utilize the same learning rate and perturbation scale. This is because the overall perturbation matrix is scaled by a factor of $\frac{1}{r}$, following a similar norm alignment strategy as detailed in Eqn. (20). We also perform ablation studies on the rank and subspace update frequency for SubZero, with results shown in Table 7. Full-parameter tuning scheme is conducted on the RTE dataset using the OPT-13B model, with specific experimental settings outlined in Table 17. All experiments employ the same learning rate and perturbation scale, enabled by the norm alignment technique described in Eqn. (20).

8.5. Prompt Templates

For autoregressive LLMs, we have three task types: classification, multiple-choice, and question answering. We adopt the prompt templates for various tasks in [40], which are summarized in Table 19. For masked LLMs, we also adopt the prompt templates in [40] and present them in Table 20.

8.6. Proofs

In practice, SubZero employs smaller and layer-specific low-rank perturbation matrices instead of a large model-scale projection matrix. However, it is more convenient to prove SubZero’s properties using a model-scale projection. Fortunately, the following lemma shows that the low-rank perturbation matrix for each layer can be represented as a layer-scale projection matrix, which is column orthogonal.

Lemma 1. *Let $\tilde{Z} = UZV^T$, where $U \in \mathbb{R}^{m \times r}$, $Z \in \mathbb{R}^{r \times r}$, $V \in \mathbb{R}^{n \times r}$, and $U^T U = V^T V = I_r$. Then we have $\text{vec}(\tilde{Z}) = P \text{vec}(Z)$ and $P^T P = I_{r^2}$, where $P = V \otimes U$.*

Proof. Since $\text{vec}(UZV^T) = (V \otimes U) \text{vec}(Z)$, we only need to show $(V \otimes U)^T (V \otimes U) = I_{r^2}$. In fact

$$(V \otimes U)^T (V \otimes U) = (V^T \otimes U^T)(V \otimes U) = (V^T V) \otimes (U^T U) = I_r \otimes I_r = I_{r^2}.$$

The proof is completed. \square

We can also demonstrate that the low-rank perturbation matrices across all layers can be represented as a model-scale projection matrix. We first give the following lemma.

Lemma 2. *Let a block diagonal matrix $P = \text{bdiag}(P_1, P_2, \dots, P_l)$ and $\tilde{z}_i = P_i z_i$, where $P_i^T P_i = I_{r^2}$ and $i = 1, 2, \dots, l$. Then we have $\tilde{z} = Pz$, where $\tilde{z} = [\tilde{z}_1^T, \dots, \tilde{z}_l^T]^T$, $z = [z_1^T, \dots, z_l^T]^T$ and $P^T P = I_{lr^2}$.*

Proof. It is easy to check that $\tilde{z} = Pz$. Besides, we have

$$P^T P = \text{bdiag}(P_1^T, \dots, P_l^T) \text{bdiag}(P_1, \dots, P_l) = \text{bdiag}(P_1^T P_1, \dots, P_l^T P_l) = I_{lr^2}.$$

The proof is completed. \square

We may define $P = \text{bdiag}(V_1 \otimes U_1, V_2 \otimes U_2, \dots, V_l \otimes U_l)$ that satisfies $P^T P = I$, $z = [\text{vec}(Z_1)^T, \text{vec}(Z_2)^T, \dots, \text{vec}(Z_l)^T]^T$, and $\tilde{z} = [\text{vec}(\tilde{Z}_1)^T, \text{vec}(\tilde{Z}_2)^T, \dots, \text{vec}(\tilde{Z}_l)^T]^T$. Then according to Lemma 2, the perturbation vector of SubZero is $\tilde{z} = Pz$, which is similar as existing random subspace methods in Eqn. (4), but with SubZero’s projection matrix being block diagonal and column orthogonal.

To prove Theorem 1 and Theorem 2, we first introduce some definitions and lemmas about Gaussian distribution.

Table 18. The hyperparameter search grids for LLama2-7B and Mistral-7B. For each task, we run 20K steps for ZO methods (MeZO and SubZero) and SGD. We record the best model checkpoint based on the validation loss every 500 training steps.

Experiment	Hyperparameter	Value
MeZO(FT)	batch size	16
	learning rate	{ 1e-7, 5e-7, 1e-6}
	ε	1e-3
MeZO(LoRA)	batch size	16
	learning rate	{ 1e-6, 5e-6, 1e-5, 3e-5}
	ε	1e-3
MeZO(Prefix)	batch size	16
	learning rate	{ 1e-3, 5e-3, 1e-2}
	ε	1e-1
MeZO(Prompt)	batch size	16
	learning rate	{ 1e-3, 5e-3, 1e-2}
	ε	1e-2
SubZero(FT)	batch size	16
	learning rate	{ 1e-7, 5e-7, 1e-6}
	ε	1e-3
	rank	{24, 48}
SubZero(LoRA)	subspace change frequency	1000
	batch size	16
	learning rate	{ 1e-6, 5e-6, 1e-5, 3e-5}
	ε	1e-3
SubZero(Prefix)	rank	{4, 8}
	subspace change frequency	1000
	batch size	16
	learning rate	{ 1e-3, 5e-3, 1e-2}
SubZero(Prompt)	ε	1e-1
	rank	{4, 8}
	subspace change frequency	1000
	batch size	16
SGD(FT)	learning rate	{ 1e-3, 5e-3, 1e-2}
	ε	1e-2
	rank	{16, 24}
	subspace change frequency	1000
SGD(FT)	batch size	16
	learning rate	{ 1e-5, 1e-4, 1e-3, 5e-3}

Definition 1. We say \mathbf{z} is a standard n -dimensional Gaussian vector (denote by $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$), if its probability density function $p(\mathbf{z}) = \frac{1}{\kappa} e^{-\frac{1}{2}\|\mathbf{z}\|^2}$, where $\kappa > 0$ satisfies $\int_{\mathbb{R}^n} \frac{1}{\kappa} e^{-\frac{1}{2}\|\mathbf{z}\|^2} d\mathbf{z} = 1$.

Definition 2. Let $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$. We say x is a chi-square random variable with degrees of freedom n (denote by $x \sim \chi^2(n)$), if $x = \|\mathbf{z}\|^2$.

Lemma 3. Let $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$. For any orthogonal $(n \times n)$ -matrix \mathbf{Q} and continuous function f , we have $\mathbb{E}_{\mathbf{z}}[f(\mathbf{z})] = \mathbb{E}_{\mathbf{z}}[f(\mathbf{Q}\mathbf{z})]$.

Lemma 4. If $x \sim \chi^2(n)$, then we have

$$\mathbb{E}_x[x] = n, \quad \text{Var}_x[x] = 2n.$$

Lemma 5. [41] Let $f \in C_{L_2}^{2,2}(\mathbb{R}^n)$. Then for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, we have

$$|f(\mathbf{y}) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle - \frac{1}{2} \langle \nabla^2 f(\mathbf{x})(\mathbf{y} - \mathbf{x}), \mathbf{y} - \mathbf{x} \rangle| \leq \frac{L_2}{6} \|\mathbf{y} - \mathbf{x}\|^3.$$

Table 19. The prompt templates used in the OPT-1.3B, OPT-13B, LLama2-7B, and Mistral-7B experiments.

Task	Type	Prompt
SST-2	cls.	<text> It was terrible/great
RTE	cls.	<premise> Does this mean that "<hypothesis>" is true? Yes or No? Yes or No
CB	cls.	Does this mean that "<hypothesis>" is true? Yes or No? Yes/No/Maybe
BoolQ	cls.	<passage> <question>? Yes/No
WSC	cls.	<text> In the previous sentence, does the pronoun "<span2>" refer to <span1>? Yes or No? Yes/No
WIC	cls.	Does the word "<word>" have the same meaning in these two sentences? Yes, No? <sentence1> <sentence2> Yes/No
MultiRC	cls.	<paragraph> Question: <question> I found this answer "<answer>". Is that correct? Yes or No? Yes/No
COPA	mch.	<premise> so/because <candidate>
ReCoRD	mch.	<passage> <query>.replace("@placeholder", <candidate>)
SQuAD	QA	Title: <title> Context: <context> Question: <question> Answer:
DROP	QA	Passage: <context> Question: <question> Answer:

Table 20. The prompt templates used in RoBERTa-large experiments. C is the number of classification categories.

Task	C	Type	Prompt
SST-2	2	sentiment cls.	<sentence1> It was great/terrible
SST-5	5	sentiment cls.	<sentence1> It was great/good/okay/bad/terrible
MNLI	3	NLI	<sentence1> ? Yes/Maybe/No , <sentence2>
SNLI	3	NLI	<sentence1> ? Yes/Maybe/No , <sentence2>

Lemma 6. [41] Let $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$. For $0 \leq t \leq 2$, we have

$$\mathbb{E}_{\mathbf{z}}[\|\mathbf{z}\|^t] \leq n^{t/2}.$$

For $t \geq 2$, we have

$$n^{t/2} \leq \mathbb{E}_{\mathbf{z}}[\|\mathbf{z}\|^t] \leq (n+t)^{t/2}.$$

Lemma 7. Let $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$. For all $\mathbf{y} \in \mathbb{R}^n$, we have

$$\mathbb{E}_{\mathbf{z}}[\|\langle \mathbf{y}, \mathbf{z} \rangle \mathbf{z}\|^2] = (n+2)\|\mathbf{y}\|^2.$$

Proof. Note that for any orthogonal $(n \times n)$ -matrix \mathbf{Q} , we have

$$\|\langle \mathbf{y}, \mathbf{Q}\mathbf{z} \rangle \mathbf{Q}\mathbf{z}\|^2 = \|\langle \mathbf{Q}^\top \mathbf{y}, \mathbf{z} \rangle \mathbf{z}\|^2, \quad \|\mathbf{Q}^\top \mathbf{y}\| = \|\mathbf{y}\|.$$

In accordance with Lemma 3, we can set $\mathbf{y} = [1, 0, \dots, 0]^\top$, and only need to prove $\mathbb{E}_{\mathbf{z}}[\|\langle \mathbf{y}, \mathbf{z} \rangle \mathbf{z}\|^2] = n + 2$. Equipped with Lemma 4, we get

$$\mathbb{E}_{\mathbf{z}}[\|\langle \mathbf{y}, \mathbf{z} \rangle \mathbf{z}\|^2] = \mathbb{E}_{\mathbf{z}} \left[\sum_{i=1}^n z_1^2 z_i^2 \right] = \sum_{i=1}^n \mathbb{E}_{\mathbf{z}}[z_1^2 z_i^2] = \mathbb{E}_{z_1}[z_1^4] + \mathbb{E}_{z_1}[z_1^2] \sum_{i=2}^n \mathbb{E}_{\mathbf{z}}[z_i^2] = n + 2.$$

The proof is completed. \square

Theorem 1. For the gradient estimation in Eqn. (8), the following two properties hold.

a) By using gradient estimation in (8), our estimated gradient $\hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z})$ is equivalent to

$$\hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z}) = \frac{f(\mathbf{x} + \varepsilon \mathbf{P} \mathbf{z}) - f(\mathbf{x} - \varepsilon \mathbf{P} \mathbf{z})}{2\varepsilon} \mathbf{P} \mathbf{z}, \quad (13)$$

where $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_q)$, $\varepsilon > 0$, $\mathbf{P} \in \mathbb{R}^{d \times q}$ satisfies $\mathbf{P}^\top \mathbf{P} = \mathbf{I}_q$ with $d = \sum_{i=1}^l m_i n_i$ and $q = lr^2$.

b) Let $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_q)$, and $f \in C_{L_2}^{2,2}(\mathbb{R}^d)$. Then we have

$$\Phi(\mathbf{x}) = \|\mathbb{E}_{\mathbf{z}}[\hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z})] - \mathbf{P} \mathbf{P}^\top \nabla f(\mathbf{x})\|_2 \leq \frac{\varepsilon^2}{6} L_2 (q + 4)^2.$$

Proof. a) Evidently, the conclusion is established based on Lemma 1 and Lemma 2.

b) Let $a_{\mathbf{z}}(\tau) = f(\mathbf{x} + \tau \mathbf{z}) - f(\mathbf{x}) - \tau \langle \nabla f(\mathbf{x}), \mathbf{z} \rangle - \frac{\tau^2}{2} \langle \nabla^2 f(\mathbf{x}) \mathbf{z}, \mathbf{z} \rangle$. Lemma 5 implies that

$$|a_{\mathbf{z}}(\pm \varepsilon)| \leq \frac{\varepsilon^3}{6} L_2 \|\mathbf{z}\|^3.$$

Note that

$$\begin{aligned} & \mathbb{E}_{\mathbf{z}}[\hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z})] - \mathbf{P} \mathbf{P}^\top \nabla f(\mathbf{x}) \\ &= \frac{\mathbf{P}}{2\kappa\varepsilon} \int_{\mathbb{R}^q} [f(\mathbf{x} + \varepsilon \mathbf{P} \mathbf{z}) - f(\mathbf{x} - \varepsilon \mathbf{P} \mathbf{z}) - 2\varepsilon \langle \nabla f(\mathbf{x}), \mathbf{P} \mathbf{z} \rangle] \mathbf{z} e^{-\frac{1}{2}\|\mathbf{z}\|^2} d\mathbf{z}. \end{aligned}$$

Therefore, in accordance with Lemma 6, we have

$$\begin{aligned} & \|\mathbb{E}_{\mathbf{z}}[\hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z})] - \mathbf{P} \mathbf{P}^\top \nabla f(\mathbf{x})\| \\ & \leq \frac{1}{2\kappa\varepsilon} \int_{\mathbb{R}^q} |f(\mathbf{x} + \varepsilon \mathbf{P} \mathbf{z}) - f(\mathbf{x} - \varepsilon \mathbf{P} \mathbf{z}) - 2\varepsilon \langle \nabla f(\mathbf{x}), \mathbf{P} \mathbf{z} \rangle| \|\mathbf{z}\| e^{-\frac{1}{2}\|\mathbf{z}\|^2} d\mathbf{z} \\ &= \frac{1}{2\kappa\varepsilon} \int_{\mathbb{R}^q} |a_{\mathbf{P} \mathbf{z}}(\varepsilon) - a_{\mathbf{P} \mathbf{z}}(-\varepsilon)| \|\mathbf{z}\| e^{-\frac{1}{2}\|\mathbf{z}\|^2} d\mathbf{z} \\ & \leq \frac{\varepsilon^2 L_2}{6\kappa} \int_{\mathbb{R}^q} \|\mathbf{z}\|^4 e^{-\frac{1}{2}\|\mathbf{z}\|^2} d\mathbf{z} \leq \frac{\varepsilon^2}{6} L_2 (q + 4)^2. \end{aligned}$$

The proof is completed. \square

Theorem 2. Let $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{H} \mathbf{x}$ and $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_q)$, where $\mathbf{H} \in \mathbb{R}^{d \times d}$ is positive definite. We have

$$\mathbb{E}_{\mathbf{z}}[\hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z})] = \mathbf{P} \mathbf{P}^\top \nabla f(\mathbf{x}), \quad (14)$$

$$\mathbb{E}_{\mathbf{z}}[\|\hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z})\|^2] = (q + 2) \|\mathbf{P}^\top \nabla f(\mathbf{x})\|^2, \quad (15)$$

$$\mathbb{E}_{\mathbf{z}} \left[\frac{\langle \nabla f(\mathbf{x}), \hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z}) \rangle^2}{\|\mathbf{P}^\top \nabla f(\mathbf{x})\|^2 \|\hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z})\|^2} \right] = \frac{1}{q}. \quad (16)$$

Proof. It is easy to check that $\hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z}) = \mathbf{P} \langle \mathbf{P}^\top \nabla f(\mathbf{x}), \mathbf{z} \rangle \mathbf{z}$. Thus we have $\mathbb{E}_{\mathbf{z}}[\hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z})] = \mathbf{P} \mathbf{P}^\top \nabla f(\mathbf{x})$. Combined with Lemma 7, we get $\mathbb{E}_{\mathbf{z}}[\|\hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z})\|^2] = (q + 2) \|\mathbf{P}^\top \nabla f(\mathbf{x})\|^2$. Note that for any orthogonal $(q \times q)$ -matrix \mathbf{Q} , we have

$$\begin{aligned}
\mathbb{E}_{\mathbf{z}} \left[\frac{\langle \nabla f(\mathbf{x}), \hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z}) \rangle^2}{\|\mathbf{P}^\top \nabla f(\mathbf{x})\|^2 \|\hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z})\|^2} \right] &= \mathbb{E}_{\mathbf{z}} \left[\frac{\langle \mathbf{P}^\top \nabla f(\mathbf{x}), \mathbf{z} \rangle^2}{\|\mathbf{P}^\top \nabla f(\mathbf{x})\|^2 \|\mathbf{z}\|^2} \right] \\
&= \mathbb{E}_{\mathbf{z}} \left[\frac{\langle \mathbf{P}^\top \nabla f(\mathbf{x}), \mathbf{Q}\mathbf{z} \rangle^2}{\|\mathbf{P}^\top \nabla f(\mathbf{x})\|^2 \|\mathbf{Q}\mathbf{z}\|^2} \right] \\
&= \mathbb{E}_{\mathbf{z}} \left[\frac{\langle \mathbf{Q}^\top \mathbf{P}^\top \nabla f(\mathbf{x}), \mathbf{z} \rangle^2}{\|\mathbf{Q}^\top \mathbf{P}^\top \nabla f(\mathbf{x})\|^2 \|\mathbf{z}\|^2} \right].
\end{aligned}$$

In accordance with Lemma 3, we can set $\mathbf{P}^\top \nabla f(\mathbf{x}) = [1, 0, \dots, 0]^\top$. Thus we have

$$\mathbb{E}_{\mathbf{z}} \left[\frac{\langle \nabla f(\mathbf{x}), \hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z}) \rangle^2}{\|\mathbf{P}^\top \nabla f(\mathbf{x})\|^2 \|\hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z})\|^2} \right] = \mathbb{E}_{\mathbf{z}} \left[\frac{z_1^2}{\|\mathbf{z}\|^2} \right] = \frac{1}{q}.$$

The proof is completed. \square

To illustrate the convergence of Subzero with SGD, our analysis is divided into two main segments. We first investigate the convergence behavior of SubZero solution process while keeping the projection matrix \mathbf{P} constant. Next, we evaluate the effects of the lazy updates to \mathbf{P} . Based on these evaluations, we establish the global convergence of Subzero. Without loss of generality, we concentrate on the scenario where the number of layers is 1.

First, when the subspace \mathbf{P} is fixed, the original problem of SubZero can be reformulated as an optimization problem within the subspace. Define $h(\mathbf{y}) = f(\mathbf{x} + \mathbf{P}\mathbf{y})$, $h_\varepsilon(\mathbf{y}) = \mathbb{E}_{\mathbf{z}}[h(\mathbf{y} + \varepsilon\mathbf{z})]$, and $g_\varepsilon(\mathbf{y}) = \frac{h(\mathbf{y} + \varepsilon\mathbf{z}) - f(\mathbf{y})}{\varepsilon}\mathbf{z}$. According to Lemma 8, if f is first L_1 -smooth, then h is also first L_1 -smooth.

Lemma 8. Let $h(\mathbf{y}) = f(\mathbf{x} + \mathbf{P}\mathbf{y})$, where $f \in C_{L_1}^{1,1}(\mathbb{R}^d)$, and $\mathbf{P}^\top \mathbf{P} = \mathbf{I}$, then we have $h \in C_{L_1}^{1,1}(\mathbb{R}^q)$.

Proof. The following proves that if f is first L_1 -smooth, then h is also first L_1 -smooth. For any $\mathbf{y}_1 \in \mathbb{R}^q$ and $\mathbf{y}_2 \in \mathbb{R}^q$, we have

$$\begin{aligned}
\|\nabla h(\mathbf{y}_1) - \nabla h(\mathbf{y}_2)\| &= \|\mathbf{P}^\top \nabla(f(\mathbf{x} + \mathbf{P}\mathbf{y}_1)) - \mathbf{P}^\top \nabla(f(\mathbf{x} + \mathbf{P}\mathbf{y}_2))\| \\
&\leq \|\mathbf{P}^\top\| \|\nabla(f(\mathbf{x} + \mathbf{P}\mathbf{y}_1)) - \nabla(f(\mathbf{x} + \mathbf{P}\mathbf{y}_2))\| \\
&\leq L_1 \|\mathbf{P}(\mathbf{y}_1 - \mathbf{y}_2)\| \\
&= L_1 \|\mathbf{y}_1 - \mathbf{y}_2\|.
\end{aligned}$$

The proof is completed. \square

Now, we can analyze the convergence of SubZero when fixing the subspace.

Lemma 9. [41] Let $f \in C_{L_1}^{1,1}(\mathbb{R})$. Then, for any $\mathbf{x} \in \mathbb{R}$, we have

$$\mathbb{E}_{\mathbf{z}}[\|g_\varepsilon(\mathbf{x})\|^2] = \mathbb{E}_{\mathbf{z}} \left[\left\| \frac{f(\mathbf{x} + \varepsilon\mathbf{z}) - f(\mathbf{x})}{\varepsilon} \right\|^2 \right] \leq 4(n+4)\|\nabla f_\varepsilon(\mathbf{x})\|^2 + 3\varepsilon^2 L_1^2(f)(n+4)^3, \quad (21)$$

and

$$\|\nabla f(\mathbf{x})\|^2 \leq 2\|\nabla f_\varepsilon(\mathbf{x})\|^2 + \frac{\varepsilon^2}{2} L_1^2(f)(n+6)^3, \quad (22)$$

where $f_\varepsilon(\mathbf{x}) = \mathbb{E}_{\mathbf{z}}[f(\mathbf{x} + \varepsilon\mathbf{z})]$.

Lemma 10. Let $\mathbf{y}^* = \arg \min_{\mathbf{x} \in \mathbb{R}^q} h(\mathbf{y})$, where $h \in C_{L_1}^{1,1}(\mathbb{R}^q)$ and h is non-convex. Suppose $\mathcal{E}_k = (\mathbf{z}_0, \mathbf{z}_1, \dots, \mathbf{z}_{k-1}, \mathbf{z}_k)$, where $\mathbf{z}_k \sim \mathcal{N}(0, \mathbf{I}_q)$ and $\eta = \frac{1}{4(q+4)L_1}$. $\{\mathbf{y}_k\}_{k>0}$ is the sequence generated by Algorithm 3. Let $\phi_0 = h(\mathbf{y}_0)$, and for $k \geq 1$, $\phi_k = \mathbb{E}_{\mathcal{E}_{k-1}}[h(\mathbf{y}_k)]$. For the \mathbf{P} defined in (10), which is fixed, we have

$$\phi_{k+1} - \phi_k \leq -\frac{1}{4}\eta \mathbb{E}_{\mathcal{E}_k}[\|\nabla h(\mathbf{y}_k)\|^2] + \frac{\varepsilon^2(q+6)^3}{8} L_1^2 + \frac{3\varepsilon^2(q+4)}{32} L_1 \quad (23)$$

Proof. If a subspace $\mathbf{P} \in \mathbb{R}^{d \times q}$ is fixed, the optimization objective can be reformulated as

$$\min_{\mathbf{y} \in \mathbb{R}^q} h(\mathbf{y}) := f(\mathbf{x} + \mathbf{P}\mathbf{y}),$$

Let \mathbf{y}_0 be an initial point and $\{\eta_k\}_{k \geq 0}$ a sequence of positive real numbers. Consider the randomized gradient search algorithm $\mathcal{RG}_\varepsilon(\varepsilon > 0)$:

- 1) Generate \mathbf{z}_k and the corresponding $g_\varepsilon(\mathbf{y}_k)$, where $\mathbf{z}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_q)$.
- 2) Update $\mathbf{y}_{k+1} = \mathbf{y}_k - \eta_k g_\varepsilon(\mathbf{y}_k)$.

We aim to estimate the evolution of the function h_ε after one iteration of this algorithm.

Given that h is L_1 -Lipschitz continuous for the first derivative, and h_ε is L_ε -Lipschitz continuous for the first derivative (where $L_\varepsilon \leq L_1$) [41]. Thus, we have

$$h_\varepsilon(\mathbf{y}_{k+1}) \leq h_\varepsilon(\mathbf{y}_k) - \eta_k \langle \nabla h_\varepsilon(\mathbf{y}_k), g_\varepsilon(\mathbf{y}_k) \rangle + \frac{1}{2} \eta_k^2 L_\varepsilon \|g_\varepsilon(\mathbf{y}_k)\|^2.$$

Taking expectation with respect to \mathbf{z}_k , we obtain

$$\mathbb{E}_{\mathbf{z}_k}[h_\varepsilon(\mathbf{y}_{k+1})] \leq h_\varepsilon(\mathbf{y}_k) - \eta_k \|\nabla h_\varepsilon(\mathbf{y}_k)\|^2 + \frac{1}{2} \eta_k^2 L_\varepsilon \mathbb{E}_{\mathbf{z}_k}[\|g_\varepsilon(\mathbf{y}_k)\|^2].$$

Since $h \in C^{1,1}(\mathbb{R}^q)$, from Lemma 9, we have

$$\begin{aligned} \mathbb{E}_{\mathbf{z}_k}[h_\varepsilon(\mathbf{y}_{k+1})] &\leq h_\varepsilon(\mathbf{y}_k) - \eta_k \|\nabla h_\varepsilon(\mathbf{y}_k)\|^2 \\ &\quad + \frac{1}{2} \eta_k^2 L_1 (4(q+4) \|\nabla h_\varepsilon(\mathbf{y}_k)\|^2 + 3\varepsilon^2 L_1^2 (q+4)^3). \end{aligned}$$

Setting $\eta_k = \hat{\eta} = \frac{1}{4(q+4)L_1}$, we get

$$\mathbb{E}_{\mathbf{z}_k}[h_\varepsilon(\mathbf{y}_{k+1})] \leq h_\varepsilon(\mathbf{y}_k) - \frac{1}{2} \hat{\eta} \|\nabla h_\varepsilon(\mathbf{y}_k)\|^2 + \frac{3\varepsilon^2}{32} L_1 (q+4).$$

Taking the expectation with respect to \mathcal{E}_k , we get

$$\phi_{k+1} \leq \phi_k - \frac{1}{2} \hat{\eta} \mathbb{E}_{\mathcal{E}_k}[\|\nabla h_\varepsilon(\mathbf{y}_k)\|^2] + \frac{3\varepsilon^2(q+4)}{32} L_1,$$

From Lemma 9, we have $\mathbb{E}_{\mathcal{E}_k}[\|\nabla h(\mathbf{y}_k)\|^2] \leq 2\mathbb{E}_{\mathcal{E}_k}[\|\nabla h_\varepsilon(\mathbf{y}_k)\|^2] + \frac{\varepsilon^2(q+6)^3}{2} L_1^2$. Therefore,

$$\phi_{k+1} - \phi_k \leq -\frac{1}{4} \hat{\eta} \mathbb{E}_{\mathcal{E}_k}[\|\nabla h(\mathbf{y}_k)\|^2] + \frac{\varepsilon^2(q+6)^3}{8} L_1^2 + \frac{3\varepsilon^2(q+4)}{32} L_1. \quad (24)$$

The proof is completed. \square

Next, we need to measure the randomness of our random subspace. From Lemma 14, if the projection matrix is obtained by Algorithm 1, we have $\mathbb{E}[\mathbf{P}\mathbf{P}^T] = \frac{q}{d} \mathbf{I}$, where q represents the dimension of the subspace, d represents the dimension of the origin space, and $\mathbf{P} = \mathbf{V} \otimes \mathbf{U}$ (see Lemma 1).

Lemma 11. Let matrix $\mathbf{A} = (\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_r) \in \mathbb{R}^{n \times r}$ be composed of column vectors \mathbf{a}_k which are mutually independent and $\mathbf{a}_k \in \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$. Suppose Gram-Schmidt process $\mathbf{u}_k = \mathbf{a}_k - \sum_{s=1}^{k-1} \langle \mathbf{a}_k, \mathbf{e}_s \rangle \mathbf{e}_s$ and $\mathbf{e}_k = \frac{\mathbf{u}_k}{\|\mathbf{u}_k\|}$. $[\mathbf{a}_k]_i \leftrightarrow [\mathbf{a}_k]_j$ represents the exchange of the i -th element and the j -th element of \mathbf{a}_k , while all other elements remain unchanged. $[\mathbf{a}_k]_i = -1 \times [\mathbf{a}_k]_i$ signifies that only the i -th element of \mathbf{a}_k is multiplied by -1 , while all other elements remain unchanged. Suppose $f(\mathbf{A}, \mathbf{U}, \mathbf{E})$ be a function of the matrix \mathbf{A} , $\mathbf{U} = (\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r)$ and $\mathbf{E} = (\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_r)$, then

- (1) if $[\mathbf{a}_k]_i \leftrightarrow [\mathbf{a}_k]_j$ or $[\mathbf{a}_k]_i = -1 \times [\mathbf{a}_k]_i$, $\mathbb{E}[f]$ remain unchanged.
- (2) if $[\mathbf{a}_k]_i \leftrightarrow [\mathbf{a}_k]_j \Rightarrow [\mathbf{u}_k]_i \leftrightarrow [\mathbf{u}_k]_j$ and $[\mathbf{e}_k]_i \leftrightarrow [\mathbf{e}_k]_j$.

(3) if $[\mathbf{a}_k]_i = -1 \times [\mathbf{a}_k]_i \Rightarrow [\mathbf{u}_k]_i = -1 \times [\mathbf{u}_k]_i$, $[\mathbf{e}_k]_i = -1 \times [\mathbf{e}_k]_i$, $[\mathbf{u}_k]_j = 1 \times [\mathbf{u}_k]_j$, and $[\mathbf{e}_k]_j = 1 \times [\mathbf{e}_k]_j$, where $i \neq j$.

$$(4) \mathbb{E} \left[\frac{[\mathbf{u}_k]_i^2}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle} \right] = \frac{1}{n}.$$

$$(5) \mathbb{E} \left[\frac{[\mathbf{u}_k]_i [\mathbf{u}_k]_j}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle} \right] = 0, \text{ where } i \neq j.$$

Proof. According to real analysis, the matrix \mathbf{A} is full rank almost everywhere under a Gaussian distribution, and both \mathbf{u}_k and \mathbf{e}_k are non-zero almost everywhere.

(1) Since \mathbf{a}_k is independently and identically distributed, it obviously holds.

(2) For base case $k = 1$, it obviously holds. Assume the result holds for all $k = 1, 2, \dots, k-1$, where $k \geq 2$, then $[\mathbf{a}_k]_i \leftrightarrow [\mathbf{a}_k]_j \Rightarrow [\mathbf{u}_k]_i = [\mathbf{a}_k]_i - \sum_{s=1}^{k-1} \langle \mathbf{a}_k, \mathbf{e}_s \rangle [\mathbf{e}_s]_i$, $[\mathbf{u}_k]_j = [\mathbf{a}_k]_j - \sum_{s=1}^{k-1} \langle \mathbf{a}_k, \mathbf{e}_s \rangle [\mathbf{e}_s]_j$, $[\mathbf{e}_k]_i = \frac{[\mathbf{u}_k]_i}{\|\mathbf{u}_k\|}$, and $[\mathbf{e}_k]_j = \frac{[\mathbf{u}_k]_j}{\|\mathbf{u}_k\|}$.

Thus, by strong induction, we have $[\mathbf{u}_k]_i \leftrightarrow [\mathbf{u}_k]_j$ and $[\mathbf{e}_k]_i \leftrightarrow [\mathbf{e}_k]_j$.

(3) For base case $k = 1$, it obviously holds. Assume the result holds for all $k = 1, 2, \dots, k-1$, where $k \geq 2$, then

$$\begin{aligned} [\mathbf{a}_k]_i = -1 \times [\mathbf{a}_k]_i &\Rightarrow \begin{cases} [\mathbf{u}_k]_i = [\mathbf{a}_k]_i \times (-1) - \sum_{s=1}^{k-1} \langle \mathbf{a}_k, \mathbf{e}_s \rangle [\mathbf{e}_s]_i \times (-1) = [\mathbf{u}_k]_i \times (-1) \\ [\mathbf{u}_k]_j = [\mathbf{u}_k]_j \times 1, i \neq j \end{cases} \\ &\Rightarrow \begin{cases} [\mathbf{e}_k]_i \times (-1) = \frac{[\mathbf{u}_k]_i}{\|\mathbf{u}_k\|} \times (-1) \\ [\mathbf{e}_k]_j = [\mathbf{e}_k]_j \times 1, j \neq i \end{cases} \end{aligned}$$

By strong induction, we have $[\mathbf{u}_k]_i = -1 \times [\mathbf{u}_k]_i$, $[\mathbf{e}_k]_i = -1 \times [\mathbf{e}_k]_i$, $[\mathbf{u}_k]_j = 1 \times [\mathbf{u}_k]_j$, and $[\mathbf{e}_k]_j = 1 \times [\mathbf{e}_k]_j$, where $i \neq j$. \square

$$(4) \text{ Since } \left| \frac{[\mathbf{u}_k]_i^2}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle} \right| \leq 1, \mathbb{E} \left[\frac{[\mathbf{u}_k]_i^2}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle} \right] \text{ exists. } [\mathbf{a}_k]_i \leftrightarrow [\mathbf{a}_k]_j \Rightarrow \frac{[\mathbf{u}_k]_i^2}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle} \leftrightarrow \frac{[\mathbf{u}_k]_j^2}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle}.$$

$$\text{Thus, } \mathbb{E} \left[\frac{[\mathbf{u}_k]_i^2}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle} \right] \times n = \sum_{s=1}^n \mathbb{E} \left[\frac{[\mathbf{u}_k]_s^2}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle} \right] = \mathbb{E} \left[\frac{\langle \mathbf{u}_k, \mathbf{u}_k \rangle}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle} \right] = 1 \Rightarrow \mathbb{E} \left[\frac{[\mathbf{u}_k]_i^2}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle} \right] = \frac{1}{n}.$$

$$(5) \text{ Since } \left| \frac{[\mathbf{u}_k]_i [\mathbf{u}_k]_j}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle} \right| \leq \left| \frac{[\mathbf{u}_k]_i^2 + [\mathbf{u}_k]_j^2}{2 \langle \mathbf{u}_k, \mathbf{u}_k \rangle} \right| \leq 1, \mathbb{E} \left[\frac{[\mathbf{u}_k]_i [\mathbf{u}_k]_j}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle} \right] \text{ exists.}$$

$$[\mathbf{a}_k]_i = [\mathbf{a}_k]_i \times -1 \Rightarrow \mathbb{E} \left[\frac{[\mathbf{u}_k]_i [\mathbf{u}_k]_j}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle} \right] = \mathbb{E} \left[\frac{-[\mathbf{u}_k]_i [\mathbf{u}_k]_j}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle} \right] = 0, \text{ where } i \neq j.$$

Lemma 12. Let $\mathbf{A} \in \mathbb{R}^{n \times r}$ be a matrix with independent standard normal entries, i.e., each element of \mathbf{A} is an i.i.d. $\mathcal{N}(0, 1)$ random variable. Suppose \mathbf{A} undergoes QR decomposition via the Gram-Schmidt process to yield a column-orthogonal matrix $\mathbf{Q} \in \mathbb{R}^{n \times r}$ with orthonormal columns $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_r$ and an upper triangular matrix $\mathbf{R} \in \mathbb{R}^{r \times r}$. Then, for each $k = 1, 2, \dots, r$, the expected value of the outer product of the k -th orthonormal column vector \mathbf{e}_k of \mathbf{Q} is given by:

$$\mathbb{E}[\mathbf{e}_k \mathbf{e}_k^T] = \frac{1}{n} \mathbf{I},$$

where \mathbf{I} is the $n \times n$ identity matrix.

Proof. By the Gram-Schmidt process, we have $\mathbf{e}_k = \frac{\mathbf{u}_k}{\|\mathbf{u}_k\|}$, where $\mathbf{u}_k = \mathbf{a}_k - \sum_{s=1}^{k-1} \langle \mathbf{a}_k, \mathbf{e}_s \rangle \mathbf{e}_s$. Thus, $\mathbf{e}_k \mathbf{e}_k^T = \frac{\mathbf{u}_k \mathbf{u}_k^T}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle}$.

The (i, j) -th entry of $\mathbb{E}[\mathbf{e}_k \mathbf{e}_k^T]$ can be written as:

$$\mathbb{E}[(\mathbf{e}_k \mathbf{e}_k^T)_{ij}] = \mathbb{E} \left[\frac{[\mathbf{u}_k]_i [\mathbf{u}_k]_j}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle} \right].$$

For diagonal entries ($i = j$): When $i = j$, from Lemma 11(4), we have:

$$\mathbb{E}[(\mathbf{e}_k \mathbf{e}_k^T)_{ii}] = \mathbb{E}\left[\frac{[\mathbf{u}_k]_i^2}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle}\right] = \frac{1}{n}.$$

For off-diagonal entries ($i \neq j$): When $i \neq j$, from Lemma 11(5), we have:

$$\mathbb{E}[(\mathbf{e}_k \mathbf{e}_k^T)_{ij}] = \mathbb{E}\left[\frac{[\mathbf{u}_k]_i [\mathbf{u}_k]_j}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle}\right] = 0.$$

Combining these two cases, we conclude that $\mathbb{E}[\mathbf{e}_k \mathbf{e}_k^T]$ is a diagonal matrix with all diagonal entries equal to $\frac{1}{n}$. Thus,

$$\mathbb{E}[\mathbf{e}_k \mathbf{e}_k^T] = \frac{1}{n} \mathbf{I},$$

where \mathbf{I} is the $n \times n$ identity matrix. The proof is completed. \square

Lemma 13. Let $\mathbf{A} \in \mathbb{R}^{n \times r}$ be a matrix with independent standard normal entries, i.e., each element of \mathbf{A} is an i.i.d. $\mathcal{N}(0, 1)$ random variable. Suppose \mathbf{A} undergoes QR decomposition to yield an orthogonal matrix $\mathbf{Q} \in \mathbb{R}^{n \times r}$ with orthonormal columns and an upper triangular matrix $\mathbf{R} \in \mathbb{R}^{r \times r}$. Then, the expected value of the outer product of the matrix \mathbf{Q} with itself is given by:

$$\mathbb{E}[\mathbf{Q}\mathbf{Q}^T] = \frac{r}{n} \mathbf{I}$$

where \mathbf{I} is the $n \times n$ identity matrix.

Proof. The QR decomposition of \mathbf{A} is given by $\mathbf{A} = \mathbf{Q}\mathbf{R}$, where \mathbf{Q} is an orthogonal matrix with columns $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_r$ and \mathbf{R} is an upper triangular matrix. Since \mathbf{Q} is orthogonal, $\mathbf{Q}\mathbf{Q}^T = \mathbf{I}_r$, where \mathbf{I}_r is the $r \times r$ identity matrix. We aim to compute $\mathbb{E}[\mathbf{Q}\mathbf{Q}^T]$. By linearity of expectation and the fact that the columns of \mathbf{Q} are orthonormal, we have:

$$\mathbb{E}[\mathbf{Q}\mathbf{Q}^T] = \mathbb{E}\left[\sum_{k=1}^r \mathbf{e}_k \mathbf{e}_k^T\right] = \sum_{k=1}^r \mathbb{E}[\mathbf{e}_k \mathbf{e}_k^T].$$

From Lemma 12, we know that $\mathbb{E}[\mathbf{e}_k \mathbf{e}_k^T] = \frac{1}{n} \mathbf{I}$ for each k . Therefore:

$$\mathbb{E}[\mathbf{Q}\mathbf{Q}^T] = \sum_{k=1}^r \frac{1}{n} \mathbf{I} = \frac{r}{n} \mathbf{I}.$$

The proof is completed. \square

Lemma 14. Let $\mathbf{A}_1 \in \mathbb{R}^{m \times r}$ and $\mathbf{A}_2 \in \mathbb{R}^{n \times r}$ be matrices with independent standard normal entries, i.e., each element of \mathbf{A}_1 and \mathbf{A}_2 is an i.i.d. $\mathcal{N}(0, 1)$ random variable. Suppose \mathbf{A}_1 and \mathbf{A}_2 undergo QR decomposition to yield orthogonal matrices $\mathbf{Q}_1 \in \mathbb{R}^{m \times r}$ and $\mathbf{Q}_2 \in \mathbb{R}^{n \times r}$ with orthonormal columns, respectively. Define $\mathbf{P} = \mathbf{Q}_2 \otimes \mathbf{Q}_1$, where \otimes denotes the Kronecker product. Then, the expected value of the outer product of the matrix \mathbf{P} with itself is given by:

$$\mathbb{E}[\mathbf{P}\mathbf{P}^T] = \frac{r^2}{mn} \mathbf{I},$$

where \mathbf{I} is the $mn \times mn$ identity matrix.

Proof. The Kronecker product $\mathbf{P} = \mathbf{Q}_2 \otimes \mathbf{Q}_1$ results in a matrix $\mathbf{P} \in \mathbb{R}^{mn \times r^2}$. From Lemma 13, we have $\mathbf{Q}_1 \mathbf{Q}_1^T = \frac{r}{m} \mathbf{I}$ and $\mathbf{Q}_2 \mathbf{Q}_2^T = \frac{r}{n} \mathbf{I}$. We aim to compute $\mathbb{E}[\mathbf{P}\mathbf{P}^T]$. Using the properties of the Kronecker product, we have:

$$\mathbb{E}[\mathbf{P}\mathbf{P}^T] = \mathbb{E}[(\mathbf{Q}_2 \otimes \mathbf{Q}_1)(\mathbf{Q}_2^T \otimes \mathbf{Q}_1^T)] = \mathbb{E}[(\mathbf{Q}_2 \mathbf{Q}_2^T) \otimes (\mathbf{Q}_1 \mathbf{Q}_1^T)] = \frac{r^2}{mn} \mathbf{I} \otimes \mathbf{I} = \frac{r^2}{mn} \mathbf{I}$$

The proof is completed. \square

Now we can assess the impact of the lazy updates to P .

Theorem 3. Let $f \in C_{L_1}^{1,1}(\mathbb{R}^d)$ be a non-convex function bounded below by f^* . Suppose $\mathcal{E}_k = (\mathbf{z}_0, \mathbf{z}_1, \dots, \mathbf{z}_k)$ with $\mathbf{z}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_q)$, and let $\mathcal{P}_j = (\mathbf{P}_0, \mathbf{P}_1, \dots, \mathbf{P}_j)$, where \mathbf{P}_j is defined in (10) with a fixed update frequency F . Then, the sequence $\{\mathbf{x}_k\}_{k>0}$ generated by Algorithm 3 satisfies:

$$\frac{1}{T} \sum_{k=0}^{T-1} \mathbb{E}_{\mathcal{E}_k, \mathcal{P}_{\lfloor k/F \rfloor}} [\|\nabla f(\mathbf{x}_k)\|^2] \leq \epsilon$$

with $T = \mathcal{O}(\frac{d}{\epsilon})$ if the perturbation scale satisfies $\epsilon \leq \mathcal{O}\left(\frac{\epsilon^{1/2}}{q^{3/2}d^{1/2}L_1^{3/2}}\right)$. Here $T = KF$, where K denotes the total number of subspace updates.

Proof. Suppose $\mathcal{P}_j = (\mathbf{P}_0, \mathbf{P}_1, \dots, \mathbf{P}_j)$, where \mathbf{P}_j is the sequence generated by Eqn. (10) and $j \leq K$. In accordance with Lemma 8 and Lemma 10, if the subspace is fixed, we can transform the original problem $f \in C_{L_1}^{1,1}(\mathbb{R}^d)$ into $h \in C_{L_1}^{1,1}(\mathbb{R}^q)$ through transformation $h(\mathbf{y}) = f(\mathbf{x} + \mathbf{P}_j\mathbf{y})$. Consider the update rule:

$$\mathbf{y}_{j,0} = 0, h_j(\mathbf{y}) = f(\mathbf{x}_{jF} + \mathbf{P}_j\mathbf{y}), \forall j \in 0, 1, \dots, K-1 \quad (25)$$

$$\mathbf{y}_{j,k} = \mathbf{y}_{j,k-1} - \eta \widehat{\nabla} h_j(\mathbf{y}_{j,k-1}), \forall k \in 0, 1, \dots, F \quad (26)$$

$$\mathbf{x}_{jF+k} = \mathbf{x}_{jF} + \mathbf{P}_j\mathbf{y}_k, \quad (27)$$

In the j -th subspace, the projection matrix \mathbf{P}_j remains constant, hence we can accumulate the changes of ϕ within the current subspace. Using Lemma 10, we have

$$\phi_{(j+1)F} - \phi_{jF} \leq -\frac{1}{4}\hat{\eta} \sum_{i=0}^{K-1} \mathbb{E}_{\mathcal{E}_{jF+i}} [\|\nabla h_j(\mathbf{y}_{j,i})\|^2] + \frac{\epsilon^2(q+6)^3}{8}KL_1^2 + \frac{3\epsilon^2(q+4)}{32}KL_1 \quad (28)$$

$$\leq -\frac{1}{4}\hat{\eta}\mathbb{E}_{\mathcal{E}_{jF}} [\|\nabla h_j(\mathbf{y}_{j,0})\|^2] + \frac{\epsilon^2(q+6)^3}{8}KL_1^2 + \frac{3\epsilon^2(q+4)}{32}KL_1. \quad (29)$$

Additionally, we note that $\nabla h_j(\mathbf{y}_{j,0}) = (\mathbf{P}_j)^\top \nabla f(\mathbf{x}_{jF})$. Taking expectations over the overall historical projection matrix \mathcal{P}_j , and noting Lemma 14, $\mathbb{E}[\mathbf{P}_j(\mathbf{P}_j)^\top] = \frac{q}{d}\mathbf{I}$, with \mathbf{P}_j independent of \mathbf{x}_{jF} , we get

$$\mathbb{E}_{\mathcal{P}_{j+1}}[\phi_{(j+1)F}] - \mathbb{E}_{\mathcal{P}_j}[\phi_{jF}] \leq -\frac{1}{4}\hat{\eta}\mathbb{E}_{\mathcal{E}_{jF}, \mathcal{P}_j} [\|(\mathbf{P}_j)^\top \nabla f(\mathbf{x}_{jF})\|^2] + \frac{\epsilon^2(q+6)^3}{8}KL_1^2 + \frac{3\epsilon^2(q+4)}{32}KL_1 \quad (30)$$

$$= -\frac{q}{4d}\hat{\eta}\mathbb{E}_{\mathcal{E}_{jF}, \mathcal{P}_j} [\|\nabla f(\mathbf{x}_{jF})\|^2] + \frac{\epsilon^2(q+6)^3}{8}KL_1^2 + \frac{3\epsilon^2(q+4)}{32}KL_1. \quad (31)$$

Assuming $f(\mathbf{x}) \geq f^*$ holds for all $\mathbf{x} \in \mathbb{R}^d$, and letting $T = KF$, summing the inequality yields

$$\mathbb{E}_{\mathcal{P}_{K-1}}[\phi_T] \leq \mathbb{E}_{\mathcal{P}_0}[\phi_0] - \frac{q}{4d}\hat{\eta} \sum_{j=0}^{K-1} \mathbb{E}_{\mathcal{E}_{jF}, \mathcal{P}_j} [\|\nabla f(\mathbf{x}_{jF})\|^2] + T\frac{\epsilon^2(q+6)^3}{8}L_1^2 + T\frac{3\epsilon^2(q+4)}{32}L_1. \quad (32)$$

Since $\mathbb{E}_{\mathcal{P}_{K-1}}[\phi_T] \geq f^*$, we have:

$$f^* \leq \mathbb{E}_{\mathcal{P}_0}[\phi_0] - \frac{q}{4d}\hat{\eta} \sum_{j=0}^{K-1} \mathbb{E}_{\mathcal{E}_{jF}, \mathcal{P}_j} [\|\nabla f(\mathbf{x}_{jF})\|^2] + T\frac{\epsilon^2(q+6)^3}{8}L_1^2 + T\frac{3\epsilon^2(q+4)}{32}L_1. \quad (33)$$

Rearranging the inequality, we get

$$\frac{q}{4d}\hat{\eta} \sum_{j=0}^{K-1} \mathbb{E}_{\mathcal{E}_{jF}, \mathcal{P}_j} [\|\nabla f(\mathbf{x}_{jF})\|^2] \leq \mathbb{E}_{\mathcal{P}_0}[\phi_0] - f^* + T\frac{\epsilon^2(q+6)^3}{8}L_1^2 + T\frac{3\epsilon^2(q+4)}{32}L_1. \quad (34)$$

Substituting $\hat{\eta} = \frac{1}{4(q+4)L_1}$, we obtain:

$$\frac{q}{16d(q+4)L_1} \sum_{j=0}^{K-1} \mathbb{E}_{\mathcal{E}_{jF}, \mathcal{P}_j} [\|\nabla f(\mathbf{x}_{jF})\|^2] \leq \mathbb{E}_{\mathcal{P}_0}[\phi_0] - f^* + T \frac{\varepsilon^2(q+6)^3}{8} L_1^2 + T \frac{3\varepsilon^2(q+4)}{32} L_1. \quad (35)$$

Thus, we have

$$\frac{1}{T} \sum_{k=0}^{T-1} \mathbb{E}_{\mathcal{E}_k, \mathcal{P}_{\lfloor k/F \rfloor}} [\|\nabla f(\mathbf{x}_k)\|^2] \leq \frac{16(q+4)dL_1(\mathbb{E}_{\mathcal{P}_0}[\phi_0] - f^*)}{qT} + \frac{2\varepsilon^2(q+6)^3(q+4)d}{q} L_1^3 + \frac{3\varepsilon^2(q+4)^2d}{2q} L_1^2. \quad (36)$$

To ensure $\sum_{k=0}^{T-1} \mathbb{E}_{\mathcal{E}_k, \mathcal{P}_{\lfloor k/F \rfloor}} [\|\nabla f(\mathbf{x}_k)\|^2] \leq \epsilon$, we can choose

$$\varepsilon \leq \mathcal{O} \left(\frac{\epsilon^{1/2}}{q^{3/2}d^{1/2}L_1^{3/2}} \right).$$

Then, the upper bound for the expected number of steps is $\mathcal{O} \left(\frac{d}{\epsilon} \right)$. The proof is completed. \square