

# Gradient Heterogeneity Complements Hessian Heterogeneity in Transformer Optimization

Anonymous authors  
Paper under double-blind review

## Abstract

Transformers are difficult to optimize with stochastic gradient descent (SGD) and largely rely on adaptive optimizers such as Adam. Despite extensive efforts, a theoretical explanation for Adam’s advantage over SGD in Transformer optimization is still incomplete. In this study, we analyze the optimization of Transformer models in the fine-tuning setting through the lens of *gradient heterogeneity*, defined as the variation in gradient norms across parameter blocks. We provide a theoretical analysis showing that gradient heterogeneity, together with Hessian heterogeneity, degrades the convergence of gradient-based methods such as SGD, while sign-based methods are substantially less sensitive to this effect. Adam’s coordinate-wise normalization makes its update directions depend mainly on gradient signs, so Adam can be interpreted as a soft variant of SignSGD. Our analysis uses the fact that SGD and SignSGD follow steepest descent directions under different norms, and derives upper bounds on the iteration complexity with implications for learning rate scaling of SignSGD. We further investigate the origin of gradient heterogeneity in Transformer architectures and show that it is strongly influenced by the placement of layer normalization, with Post-LN architectures exhibiting particularly pronounced heterogeneity. Experimental results from fine-tuning Transformers in both NLP and vision domains validate our theoretical analysis.

## 1 Introduction

Transformers (Vaswani, 2017) have achieved significant success across a wide range of tasks, particularly in language models. In practice, training language models largely relies on adaptive optimization methods (Liu et al., 2024; Grattafiori et al., 2024) such as Adam (Kingma & Ba, 2015). In contrast, while stochastic gradient descent (SGD) has long been a standard optimizer in deep learning (Lecun et al., 1998; He et al., 2016), it often exhibits inferior optimization behavior in Transformer architectures (Schmidt et al., 2021; Choi, 2019; Zhang et al., 2020b; Kunstner et al., 2023; Zhang et al., 2024a; Kunstner et al., 2024).

However, the underlying reasons for the performance gap are not yet fully understood. In particular, Adam has been shown to outperform SGD even in full-batch settings, while SignSGD (Bernstein et al., 2018), which serves as an effective proxy for Adam (Xie & Li, 2024; Li et al., 2025), achieves comparable performance under the same conditions (Kunstner et al., 2023). These observations suggest that the difference between Adam and SGD cannot be explained solely by stochastic gradient noise, but rather reflects fundamental differences between SGD and adaptive optimization methods. Other explanations, such as Adam’s robustness to heavy-tailed label distributions (Kunstner et al., 2024), capture certain aspects of this gap but do not fully account for the behavior observed in fine-tuning regimes with a small amount of labeled data. More recently, Zhang et al. (2024a) associated the Adam–SGD gap with *Hessian heterogeneity* in Transformers, defined as differences in block-wise Hessian spectra, although the underlying mechanism remains unclear.

In this study, we take a step toward a better understanding of the difference between Adam and SGD through a theoretical analysis. Specifically, we compare their optimization behaviors by analyzing their *iteration complexity*, defined as the number of optimization steps required for the gradient norm to become sufficiently small. Our analysis reveals that *gradient heterogeneity* and Hessian heterogeneity (Zhang et al., 2024a) jointly

Table 1: Comparison with prior studies on Transformer optimization. ✓: Supported; -: Not supported.

| Paper                  | Transformer | Theoretical complexity | Heterogeneity          | Layer normalization |
|------------------------|-------------|------------------------|------------------------|---------------------|
| Zhang et al. (2020b)   | ✓           | ✓                      | -                      | -                   |
| Crawshaw et al. (2022) | ✓           | ✓                      | -                      | -                   |
| Kunstner et al. (2023) | ✓           | -                      | -                      | -                   |
| Pan & Li (2022)        | ✓           | -                      | -                      | -                   |
| Kunstner et al. (2024) | ✓           | -                      | -                      | -                   |
| Zhang et al. (2024a)   | ✓           | -                      | ✓ (Hessian)            | -                   |
| <b>Ours</b>            | ✓           | ✓                      | ✓ (Gradient & Hessian) | ✓                   |

play an important role in shaping these differences. Gradient heterogeneity is defined as the variation in gradient norms across parameter blocks and is amenable to empirical analysis.

We begin by deriving upper bounds on the iteration complexity of gradient-based and sign-based optimization methods in both deterministic and stochastic settings. Our analysis uses the fact that SGD and SignSGD correspond to steepest descent directions under different norms. Our results suggest that gradient-based methods are more sensitive to gradient and Hessian heterogeneity than sign-based methods, and also provide implications for the learning rate of SignSGD. To further investigate the origin of heterogeneity, we analyze gradient heterogeneity in Transformers and examine how it relates to architectural design choices. In particular, we find that applying layer normalization after residual connections amplifies gradient heterogeneity.

Our contributions are summarized as follows. Table 1 compares prior studies with ours.

- We derive upper bounds on the iteration complexity for optimization algorithms in both deterministic and stochastic settings. Our analysis suggests that SGD is highly sensitive to heterogeneity across parameter blocks, whereas sign-based (Adam-like) methods are less affected (Theorems 4.6 and 4.8). The results yield implications for learning rate scaling of SignSGD.
- We investigate gradient heterogeneity in Transformers, identifying the position of layer normalization as a factor influencing it (Section 4.7).
- Overall, we emphasize that the sign-based nature of Adam helps address optimization challenges arising from heterogeneity across parameter blocks, which is a characteristic of Transformer architectures.

## 2 Related work

**Adam in deep learning.** Adam (Kingma & Ba, 2015) is a widely used optimization algorithm in deep learning with convergence properties (Zhang et al., 2022). However, the reasons for its superior performance are not yet fully understood. Jiang et al. (2024) empirically observed that Adam tends to converge to parameter regions with uniform diagonal elements in the Hessian, supported by theoretical analysis based on two-layer linear models. Rosenfeld & Risteski (2024) argued that the ability of Adam to handle outliers in features is a critical factor in its effectiveness. Additionally, Kunstner et al. (2024) attributed the performance of Adam in language models to its ability to manage heavy-tailed class imbalance. Orvieto & Gower (2025) showed that setting  $\beta_1 = \beta_2$  preserves Adam’s performance and enables a mean-field variational interpretation. In this study, we provide a theoretical explanation of Adam’s advantage by focusing on heterogeneity across parameter blocks.

**Sign-based optimization and variants.** SignSGD, also known as sign descent, is an optimization method that is computationally efficient and memory-efficient, making it suitable for distributed training (Bernstein et al., 2018). Adam can be interpreted as a variance-adapted variant of SignSGD (Balles & Hennig, 2018). For example, Xie & Li (2024) analyzed the convergence properties of Adam from this perspective. Consistent with this interpretation, Zhao et al. (2025) found that sign-based optimizers restore the stability and performance of Adam and proposed using adaptive learning rates for each layer. Several variants of sign-based optimization have been proposed, such as block-wise adaptive learning rates (Zhang et al., 2024b) and error-feedback

schemes that mitigate bias and improve convergence (Karimireddy et al., 2019). Through program search, a sign-based optimization algorithm called Lion (evolved sign momentum) was discovered (Chen et al., 2024b), and its effectiveness was shown by Chen et al. (2024a). Our analysis theoretically clarifies why sign-based methods are less sensitive to gradient and Hessian heterogeneity than gradient-based methods.

**Optimization challenges in Transformers.** A key aspect of Transformer optimization is the notable superiority of Adam over SGD. Zhang et al. (2020b) attributed this to the heavy-tailed gradient noise, but Kunstner et al. (2023) later challenged this, arguing that the superior performance of Adam can be attributed to sign-based characteristics rather than gradient noise, supported by full-batch experiments. Li et al. (2025) demonstrated the similarity between Adam and SignSGD in the optimization and generalization of two-layer transformers. Pan & Li (2022) showed that, in Transformers, Adam updates exhibit lower directional sharpness than SGD. Ahn et al. (2024) demonstrated that linear Transformers exhibit optimization behaviors similar to standard Transformers. Zhang et al. (2024a) revealed that the Hessian spectrum of the loss function in Transformers is heterogeneous and suggested that this is one cause of the Adam-SGD performance gap. This heterogeneity was later confirmed by Ormaniec et al. (2025), who explicitly derived the Hessian of Transformers. Our work complements these studies by offering a theoretical analysis that highlights the heterogeneity across parameter blocks in Transformer optimization.

### 3 Preliminaries

This section introduces the notation and outlines the optimization methods relevant to our study.

#### 3.1 Notation and setup

**Vectors and matrices.** The  $k$ -th element of a vector  $\mathbf{a}$  is denoted by  $\mathbf{a}_k$ , and for a matrix  $\mathbf{A}$ , we use  $\mathbf{A}_{k,:}$ ,  $\mathbf{A}_{:,l}$ , and  $A_{k,l}$  to denote the  $k$ -th row,  $l$ -th column, and element at  $(k, l)$ , respectively. When a vector or matrix is split into blocks,  $[\cdot]_b$  denotes the  $b$ -th block. The  $\ell_q$  norm is denoted by  $\|\cdot\|_q$  for vectors and represents the operator norm for matrices. The all-ones vector and identity matrix of size  $a$  are denoted by  $\mathbf{1}_a$  and  $\mathbf{I}_a$ , respectively. The operator  $\text{blockdiag}(\cdot)$  constructs block diagonal matrices. Derivatives are computed using the numerator layout.

**Model and training.** We consider a classification task with  $C$  classes and sample space  $\mathcal{X}$ . The model  $\mathbf{f}(\cdot; \boldsymbol{\theta}) : \mathcal{X} \rightarrow \mathbb{R}^C$  is parameterized by  $\boldsymbol{\theta} \in \mathbb{R}^P$ , which is divided into  $B$  blocks, denoted as  $[\boldsymbol{\theta}]_b \in \mathbb{R}^{P_b}$ , with  $\sum_{b=1}^B P_b = P$ . The training dataset  $\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$  consists of  $N$  samples  $\mathbf{x}^{(i)} \in \mathcal{X}$  and the corresponding labels  $y^{(i)} \in \{1, \dots, C\}$ . The training objective is to minimize the training loss  $L(\boldsymbol{\theta}) := \frac{1}{N} \sum_{i=1}^N \ell(\mathbf{f}(\mathbf{x}^{(i)}; \boldsymbol{\theta}), y^{(i)})$ . Here,  $\ell : \mathbb{R}^C \times \{1, \dots, C\} \rightarrow \mathbb{R}$  denotes the loss function. The element-wise sign function is denoted by  $\text{sign}(\cdot)$ . The mini-batch loss is denoted by  $\widehat{L}(\boldsymbol{\theta})$ , and the learning rate at step  $t$  is represented by  $\eta_t$ .

#### 3.2 Optimization algorithms

**Adam.** Adam (Kingma & Ba, 2015) is widely used in deep learning. It uses the first and second moment estimates of the gradient  $\nabla \widehat{L}(\boldsymbol{\theta}_t)$ , denoted as  $\mathbf{m}_t$  and  $\mathbf{v}_t$ , computed using an exponential moving average to reduce mini-batch noise. The update is performed coordinate-wise as:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta_t \frac{\widehat{\mathbf{m}}_t}{\sqrt{\widehat{\mathbf{v}}_t + \epsilon}},$$

where  $\widehat{\cdot}$  denotes bias correction and  $\epsilon$  is a small constant for numerical stability.

**Adaptive learning rate and SignSGD.** A key feature of Adam is its *adaptive learning rate*, which is computed in a coordinate-wise manner. When the hyperparameter  $\epsilon$ , which is typically set close to zero, is ignored and the ratio  $|\widehat{\mathbf{m}}_t|/\sqrt{\widehat{\mathbf{v}}_t}$  is close to 1, Adam behaves similarly to SignSGD (Balles & Hennig, 2018; Bernstein et al., 2018). SignSGD updates the parameters with momentum  $\mathbf{m}_t$  as:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta_t \text{sign}(\mathbf{m}_t).$$

This method has the property that the updates are invariant to the scale of the gradient. In this sense, Adam can be seen as a soft version of SignSGD. Additionally, the optimizer RMSProp (Tieleman & Hinton, 2017), which inspired Adam, was originally motivated by the idea of using the sign of the gradient in a mini-batch setting. RMSProp is similar to Adam but without the momentum term.

**SGD and gradient clipping.** SGD can also be modified to achieve scale invariance. A simple way to introduce scale invariance is to normalize the learning rate by the gradient norm, a technique known as normalized gradient descent. This method has been shown to be equivalent to gradient clipping up to a constant factor in the learning rate (Zhang et al., 2020a). Gradient clipping is commonly used to stabilize training, particularly in cases where large gradient magnitudes cause instability, and is often applied alongside other optimizers. However, a key difference between Adam and SGD is that SGD does not adapt the learning rate in a coordinate-wise manner.

**Steepest descent direction.** SGD and SignSGD can be interpreted as updating in the direction of *steepest descent* (Boyd & Vandenberghe, 2004; Xie & Li, 2024; Bernstein & Newhouse, 2024):

$$\Delta_t \in \arg \min_{\|\Delta\| \leq 1} \nabla \widehat{L}(\boldsymbol{\theta}_t)^\top \Delta.$$

The steepest descent direction associated with the norms  $\|\cdot\|_2$  and  $\|\cdot\|_\infty$  corresponds to the updates of SGD and SignSGD, respectively.

The steepest descent direction satisfies

$$\nabla \widehat{L}(\boldsymbol{\theta}_t)^\top \Delta = -\|\nabla \widehat{L}(\boldsymbol{\theta}_t)\|_*,$$

where  $\|\cdot\|_*$  denotes the dual norm of  $\|\cdot\|$ . Thus, evaluating the gradient using the dual norm is natural for analyzing SGD and SignSGD, as it quantifies the steepest decrease in a given descent direction under a unit-norm constraint.

## 4 Main results

In this section, we theoretically analyze optimization methods. We first introduce the setting, assumptions (Section 4.1), and complexity measures (Section 4.2), and then examine gradient–Hessian correlations (Section 4.3). Next, we derive upper bounds on the iteration complexity in deterministic (Section 4.4) and stochastic (Section 4.5) settings, together with implications for the learning rate of SignSGD. Finally, we investigate gradient heterogeneity in Transformers (Section 4.7). Our analysis suggests that heterogeneity across parameter blocks, a characteristic of Transformers, contributes to the Adam–SGD performance gap.

### 4.1 Setting and assumptions

**Gradient-based and sign-based sequences.** Kunstner et al. (2023) showed that in full-batch settings without gradient noise, SignSGD performs similarly to Adam and outperforms SGD. This suggests that the performance gap between Adam and SGD arises from differences between SignSGD and SGD. Other studies have also used SignSGD as a proxy for Adam in their analyses (Balles & Hennig, 2018; Li et al., 2025; Kunstner et al., 2024).

On the basis of these insights, we analyze the difference between parameter sequences  $\{\boldsymbol{\theta}_t^{\text{Grad}}\}_{t=0}^\infty$  and  $\{\boldsymbol{\theta}_t^{\text{Sign}}\}_{t=0}^\infty$ , referred to as the gradient-based and sign-based sequences, respectively. These sequences correspond to updates performed by gradient-based and sign-based optimization. In deterministic settings, these updates are defined as follows:

$$\begin{aligned} \boldsymbol{\theta}_{t+1}^{\text{Grad}} &= \boldsymbol{\theta}_t^{\text{Grad}} - \eta_t \nabla L(\boldsymbol{\theta}_t^{\text{Grad}}), \\ \boldsymbol{\theta}_{t+1}^{\text{Sign}} &= \boldsymbol{\theta}_t^{\text{Sign}} - \eta_t \text{sign}(\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})). \end{aligned}$$

In the stochastic setting, the loss  $L$  is replaced with the mini-batch loss  $\widehat{L}$ .

We consider fine-tuning settings, in which the parameter  $\theta$  can typically be assumed to remain within a region  $\mathcal{R}_{\text{FT}}$  throughout training. This assumption restricts  $\theta$  to the localized region  $\mathcal{R}_{\text{FT}}$ , allowing further assumptions to be applied within this region.

**Assumption 4.1** (Fine-tuning). The parameter  $\theta$  remains within the region  $\mathcal{R}_{\text{FT}}$  throughout training and there exists  $\theta_* \in \mathcal{R}_{\text{FT}}$  such that  $L_* := L(\theta_*) = \min_{\theta \in \mathcal{R}_{\text{FT}}} L(\theta)$ .

We assume Lipschitz continuity for the Hessian matrix of the loss function, a standard assumption in optimization analysis (Nesterov, 2013).

**Assumption 4.2** (Lipschitz continuity (Nesterov, 2013)). Within the region  $\mathcal{R}_{\text{FT}}$ , the loss function  $L$  is twice differentiable, and its Hessian matrix is  $\rho_H$ -Lipschitz continuous

$$\|\nabla^2 L(\theta) - \nabla^2 L(\theta')\|_2 \leq \rho_H \|\theta - \theta'\|_2.$$

Additionally, empirical studies have shown that Hessian matrices of deep learning models often exhibit a near-block-diagonal structure (Maes et al., 2024; Kunstner et al., 2024; Collobert, 2004; Zhang et al., 2024a; Zhao et al., 2025). The block-diagonal approximation is also used in optimization methods (Martens & Grosse, 2015; Zhang et al., 2017). Thus, we assume that the Hessian matrix of the loss function is close to block-diagonal.

**Assumption 4.3** (Near block-diagonal Hessian). Within the region  $\mathcal{R}_{\text{FT}}$ , the Hessian matrix can be approximated by a block-diagonal matrix with an approximation error  $\delta_D$ :

$$\|\nabla^2 L(\theta) - \nabla^2 L_D(\theta)\|_2 \leq \delta_D, \tag{1}$$

for all  $\theta \in \mathcal{R}_{\text{FT}}$ , where

$$\nabla^2 L_D(\theta) := \text{blockdiag}(\{[\nabla^2 L(\theta)]_b\}_{b=1}^B)$$

represents the block-diagonal approximation.

Note that in Eq. (1), the left-hand side is bounded above by the sum of squared elements in the non-diagonal blocks, following the relationship between  $\|\cdot\|_2$  and the Frobenius norm.

## 4.2 Gradient heterogeneity and complexity measure

**Gradient heterogeneity.** We define *gradient heterogeneity* as the disparity in gradient norms across different parameter blocks,  $\{\|\nabla L(\theta)\|_2\}_{b=1}^B$ .

This concept complements *Hessian heterogeneity*, introduced by Zhang et al. (2024a) (referred to as “block heterogeneity” in their paper), which is defined in terms of differences in the Hessian spectrum and is generally more difficult to analyze empirically than gradient heterogeneity. We characterize gradient heterogeneity quantitatively through visualizations (Figures 3 and S.2) and Gini coefficients (Table S.5), offering concrete measures.

**Weighted Hessian norms.** To analyze the complexity of optimization, we define the following two measures.

**Definition 4.4** (Weighted Hessian norms). The gradient-weighted Hessian norm  $\Lambda_G$  and parameter-weighted Hessian norm  $\Lambda_P$  are defined as:

$$\Lambda_G := \sup_{\theta \in \mathcal{R}_{\text{FT}}^+} \sum_{b=1}^B \frac{\|[\nabla L(\theta)]_b\|_2^2}{\|\nabla L(\theta)\|_2^2} \|[\nabla^2 L(\theta)]_b\|_2,$$

$$\Lambda_P := \sup_{\theta \in \mathcal{R}_{\text{FT}}} \sum_{b=1}^B \frac{P_b}{P} \|[\nabla^2 L(\theta)]_b\|_2.$$

Here, we define  $\mathcal{R}_{\text{FT}}^+ := \{\theta \in \mathcal{R}_{\text{FT}} : \|\nabla L(\theta)\|_2 > 0\}$ .

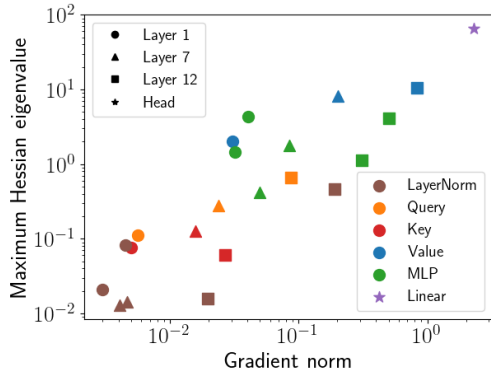


Figure 1: **Correlation between gradient norm and maximum Hessian eigenvalue.** Each point denotes the mean value computed over a parameter block (pre-trained RoBERTa on RTE).

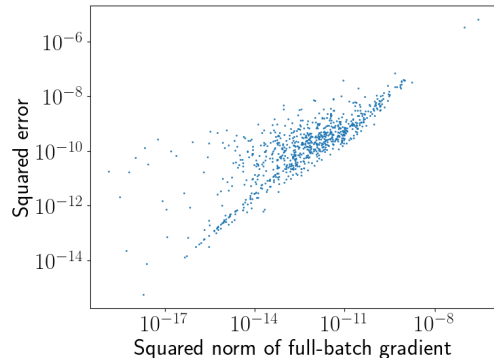


Figure 2: **Correlation between the full-batch gradient and gradient error.** Each point represents the absolute values of a coordinate (pre-trained RoBERTa on RTE).

We define  $\Lambda_G$  over  $\mathcal{R}_{\text{FT}}^+$  to avoid the degenerate stationary case; this does not affect our iteration-complexity analysis.

In these definitions,  $\Lambda_G$  weights the operator norm of each Hessian block by the squared gradient norm of the corresponding block, while  $\Lambda_P$  weights it by the parameter dimension. The definitions ensure that the weights of all Hessian blocks sum to 1, as shown by the equalities:  $\sum_{b=1}^B \frac{\|[\nabla L(\boldsymbol{\theta})]_b\|_2^2}{\|\nabla L(\boldsymbol{\theta})\|_2^2} = \sum_{b=1}^B \frac{P_b}{P} = 1$ .

### 4.3 Gradient-Hessian correlation

As shown in Figure 1, large Hessian operator norms  $\|[\nabla^2 L(\boldsymbol{\theta})]_b\|_2$  are often associated with large gradient magnitudes  $\|[\nabla L(\boldsymbol{\theta})]_b\|_2$ . In contrast, no such correlation is observed between the Hessian operator norm  $\|[\nabla^2 L(\boldsymbol{\theta})]_b\|_2$  and the parameter dimension  $P_b$ , as detailed in Appendix D.1. Under gradient-Hessian correlation, large gradient heterogeneity leads to an increase in  $\Lambda_G$ , whereas  $\Lambda_P$  remains relatively small.

**Approximate explanation.** If the loss function  $L$  is approximated in the region  $\mathcal{R}_{\text{FT}}$  by a second-order Taylor expansion around the optimum  $\boldsymbol{\theta}_* \in \mathcal{R}_{\text{FT}}$ , where  $\nabla L(\boldsymbol{\theta}_*)$  is close to  $\mathbf{0}$ , and the Hessian matrix is assumed to be block-diagonal, the following inequality approximately holds:

$$\|[\nabla L(\boldsymbol{\theta})]_b\|_2 \leq \|[\nabla^2 L(\boldsymbol{\theta}_*)]_b\|_2 \|\delta_{\boldsymbol{\theta}}\|_2,$$

where  $\delta_{\boldsymbol{\theta}} = \boldsymbol{\theta} - \boldsymbol{\theta}_*$ . This inequality suggests a positive correlation between the gradient norm and the Hessian norm.

**Support from prior studies.** This gradient-Hessian correlation was observed or assumed in previous studies. For instance, Zhang et al. (2024a); Jiang et al. (2024) demonstrated the relationship between  $|\nabla L(\boldsymbol{\theta})_i|$  and  $|\nabla^2 L(\boldsymbol{\theta})_{i,i}|$ . Additionally, the  $(L_0, L_1)$ -smoothness assumption (Zhang et al., 2020a) and its coordinate-wise generalization (Crawshaw et al., 2022) reflect this correlation.

### 4.4 Complexity bound

To analyze optimization algorithms, we define a complexity measure inspired by Carmon et al. (2020); Zhang et al. (2020a); Crawshaw et al. (2022). This measure reflects the number of parameter updates needed to achieve a sufficiently small gradient norm, with higher complexity indicating slower convergence.

**Definition 4.5** (Iteration complexity). We define the iteration complexity of a parameter sequence  $\{\boldsymbol{\theta}_t\}_{t=0}^{\infty}$  for  $\boldsymbol{\theta}_t \in \mathbb{R}^P$  with the loss function  $L$  and the norm  $\|\cdot\|_q$ :

$$\mathcal{T}_{\varepsilon}(\{\boldsymbol{\theta}_t\}_{t=0}^{\infty}, L, \|\cdot\|_q) := \inf\{t \in \mathbb{N} \mid \mathcal{C}_{\varepsilon}(t)\},$$

where the condition  $\mathcal{C}_\varepsilon(t)$  is defined as follows.

In the deterministic setting,  $\mathcal{C}_\varepsilon(t)$  is defined as:

$$\|\nabla L(\boldsymbol{\theta}_t)\|_q \leq P^{\frac{1}{q}} \varepsilon.$$

In the stochastic setting,  $\mathcal{C}_\varepsilon(t)$  is defined as:

$$\mathbb{P}(\forall s \leq t, \|\nabla L(\boldsymbol{\theta}_s)\|_q \geq P^{\frac{1}{q}} \varepsilon) \leq \frac{1}{2}.$$

Compared with the complexity definitions in previous studies, we introduce a distinction in the choice of norms and a normalization term  $P^{\frac{1}{q}}$  to ensure dimensional consistency across different norms.

Using this measure, we derive complexity bounds in deterministic (i.e., full-batch) settings. The parameter  $\zeta_0 \in (0, 1)$  controls the range of learning rates.

**Theorem 4.6** (Deterministic setting). *Assume  $\delta_D < \min(\Lambda_G, \Lambda_P)/3$ . Then, the iteration complexities in the deterministic setting are bounded as follows.*

For the gradient-based sequence, suppose that  $\varepsilon < \frac{\Lambda_G^2}{\rho_H \sqrt{P}}$  holds and that the learning rate at time  $t$  satisfies  $\eta_t = \zeta_t \min(\frac{1}{\Lambda_G}, \frac{1}{\sqrt{\rho_H \|\nabla L(\boldsymbol{\theta}_t^{Grad})\|_2}})$ , where  $\zeta_t \in [\zeta_0, 1]$ , we have

$$\mathcal{T}_\varepsilon(\{\boldsymbol{\theta}_t^{Grad}\}_{t=0}^\infty, L, \|\cdot\|_2) \leq \frac{6(L(\boldsymbol{\theta}_0) - L_*)}{P\varepsilon^2\zeta_0} \Lambda_G.$$

For the sign-based sequence, suppose that  $\varepsilon < \frac{\Lambda_P^2}{\rho_H \sqrt{P}}$  holds and that the learning rate at time  $t$  satisfies

$\eta_t = \zeta_t \min(\frac{\|\nabla L(\boldsymbol{\theta}_t^{Sign})\|_1}{\Lambda_P P}, \sqrt{\frac{\|\nabla L(\boldsymbol{\theta}_t^{Sign})\|_1}{\rho_H P^{3/2}}})$ , where  $\zeta_t \in [\zeta_0, 1]$ , we have

$$\mathcal{T}_\varepsilon(\{\boldsymbol{\theta}_t^{Sign}\}_{t=0}^\infty, L, \|\cdot\|_1) \leq \frac{6(L(\boldsymbol{\theta}_0) - L_*)}{P\varepsilon^2\zeta_0} \Lambda_P.$$

The iteration complexity of the gradient-based and sign-based sequences is evaluated using the norms  $\|\cdot\|_2$  and  $\|\cdot\|_1$ , respectively. This choice of norms is justified because they correspond to the dual norms that determine the steepest descent direction, as discussed in Section 3.2. We provide the proof and its intuition in Appendix A.

#### Gradient heterogeneity can increase the iteration complexity of the gradient-based sequence.

The theorem indicates that the iteration complexity of the gradient-based and sign-based sequences is characterized by  $\Lambda_G$  and  $\Lambda_P$ , respectively. As discussed earlier, under gradient–Hessian correlation, large gradient heterogeneity leads to a large  $\Lambda_G$ . Consequently, the iteration complexity of the gradient-based sequence can surpass that of the sign-based sequence under such conditions.

**Connection to Zhang et al. (2024a).** Zhang et al. (2024a) show that the Adam–SGD gap arises from Hessian heterogeneity. This finding is consistent with our theoretical results (Theorem 4.6), and our analysis further explains this gap by taking gradient heterogeneity into account.

#### 4.5 Stochastic setting

In practice, optimization is performed in a stochastic setting, where the gradient is estimated using a mini-batch. In this setting, we add the assumptions about noise, defined as the difference between the full-batch and mini-batch gradient.

**Assumption 4.7** (Noise). For all  $\boldsymbol{\theta} \in \mathcal{R}_{\text{FT}}$ , there exist constants  $\sigma_3, \sigma_2 \geq 0$  such that:

$$\mathbb{E}[\nabla \widehat{L}(\boldsymbol{\theta})] = \nabla L(\boldsymbol{\theta}), \tag{2}$$

$$\mathbb{E}[\|\nabla\hat{L}(\boldsymbol{\theta}) - \nabla L(\boldsymbol{\theta})\|_2^3] \leq \sigma_3 \|\nabla L(\boldsymbol{\theta})\|_2^3, \quad (3)$$

and for all  $i \in \{1, \dots, P\}$ ,

$$\mathbb{E}[|\nabla\hat{L}(\boldsymbol{\theta})_i - \nabla L(\boldsymbol{\theta})_i|^2] \leq \sigma_2 |\nabla L(\boldsymbol{\theta})_i|^2. \quad (4)$$

The assumption in Eq. (2) is standard in stochastic optimization (Bernstein et al., 2018). We introduce Eq. (3) to bound the third-order moment of the gradient noise norm and Eq. (4) to model its coordinate-wise correlation with the gradient. This correlation is supported by Figure 2 (additional settings in Appendix D.3). The coordinate-wise assumption is needed for analyzing errors in the gradient sign and block-wise gradient. Additionally, bounding the noise is a common practice in stochastic optimization (Crawshaw et al., 2022; Zhang et al., 2020a).

Using these assumptions, we establish the complexity bounds for the stochastic setting, where  $\zeta_0 \in (0, 1)$  controls the range of learning rates as in the deterministic setting.

**Theorem 4.8** (Stochastic setting). *Assume  $\delta_D < \min(\Lambda_G, \Lambda_P)/3$ . Then, the iteration complexities in the stochastic setting are bounded as follows.*

*For the gradient-based sequence, suppose that  $\varepsilon < \frac{(1+\sigma_2)^2 \Lambda_G^2}{4(1+\sigma_3)\rho_H \sqrt{P}}$  holds and that the learning rate at time  $t$  satisfies  $\eta_t = \zeta_t \min(\frac{1}{(1+\sigma_2)\Lambda_G}, \frac{1}{2\sqrt{(1+\sigma_3)\rho_H \|\nabla L(\boldsymbol{\theta}_t^{Grad})\|_2}})$ , where  $\zeta_t \in [\zeta_0, 1]$ , we have*

$$\mathcal{T}_\varepsilon(\{\boldsymbol{\theta}_t^{Grad}\}_{t=0}^\infty, L, \|\cdot\|_2) \leq \frac{12(1+\sigma_2)(L(\boldsymbol{\theta}_0) - L_*)}{P\varepsilon^2 \zeta_0} \Lambda_G.$$

*For the sign-based sequence, suppose that  $\varepsilon < \frac{\Lambda_P^2}{\rho_H \sqrt{P}}$  and  $\sigma_2 \leq \frac{1}{24}$  hold and that the learning rate at time  $t$  satisfies  $\eta_t = \zeta_t \min(\frac{\|\nabla L(\boldsymbol{\theta}_t^{Sign})\|_1}{\Lambda_P P}, \sqrt{\frac{\|\nabla L(\boldsymbol{\theta}_t^{Sign})\|_1}{\rho_H P^{3/2}}})$ , where  $\zeta_t \in [\zeta_0, 1]$ , we have*

$$\mathcal{T}_\varepsilon(\{\boldsymbol{\theta}_t^{Sign}\}_{t=0}^\infty, L, \|\cdot\|_1) \leq \frac{12(1+24\sigma_2)(L(\boldsymbol{\theta}_0) - L_*)}{P\varepsilon^2 \zeta_0} \Lambda_P.$$

This theorem shows that the dependence on the noise is the same for both sequences up to a constant, so the difference in noise dependence may be minor. Therefore, the performance gap is more likely due to the difference between  $\Lambda_G$  and  $\Lambda_P$ , as in the deterministic setting. We further analyze the setting with learning rates adapted to the noise in Appendix G.

#### 4.6 Implication for learning rates of SignSGD

Theorem 4.6 requires the learning rate to satisfy

$$\eta_t = \zeta_t \min\left(\frac{\|\nabla L(\boldsymbol{\theta}_t^{Sign})\|_1}{\Lambda_P P}, \sqrt{\frac{\|\nabla L(\boldsymbol{\theta}_t^{Sign})\|_1}{\rho_H P^{3/2}}}\right),$$

where both terms scale monotonically with  $\|\nabla L(\boldsymbol{\theta}_t^{Sign})\|_1$ . This implies that SignSGD should adapt its step size according to the  $\ell_1$ -norm of the gradient.

In the fine-tuning regime, where gradients are typically small, the linear term dominates whenever  $\|\nabla L(\boldsymbol{\theta}_t^{Sign})\|_1 \leq (\Lambda_P^2/\rho_H)\sqrt{P}$ . Thus, using the mini-batch loss  $\hat{L}$  in practice, the learning rate condition for SignSGD effectively reduces to  $\eta_t := \gamma_t \|\nabla \hat{L}(\boldsymbol{\theta}_t^{Sign})\|_1$ , where  $\gamma_t$  is a hyperparameter and we refer to this method, SignSGD with  $\ell_1$ -scaled learning rates, as SignSGD (S). This scaling corresponds to the steepest descent with respect to the  $\ell_\infty$ -norm (Appendix H.1) (Balles et al., 2020; Bernstein & Newhouse, 2024), and is also recovered as the optimal learning rate in our quadratic analysis (Appendix D.9).

## 4.7 Optimization of Transformers

Transformers show much greater parameter heterogeneity than other models (Zhang et al., 2024a; Cui & Wang, 2024), as confirmed by our experiments (Figure 3). On the basis of Theorems 4.6 and 4.8, we identify gradient heterogeneity as a key factor in the performance gap between Adam and SGD in Transformers. Here, we discuss the role of layer normalization in Transformers.

**Post-LN and Pre-LN.** In Transformers, residual connections and layer normalizations are combined with multi-head attention and feed-forward networks. The two main Transformer architectures are post-layer normalization (Post-LN), where the residual connection is followed by the layer normalization, and pre-layer normalization (Pre-LN), where the layer normalization precedes the residual connection. Pre-LN is known for greater stability (Wang et al., 2019b; Xiong et al., 2020; Takase et al., 2022).

**Jacobian of Transformers.** The Jacobians of a Transformer layer with Pre-LN and Post-LN are expressed as:

$$\mathbf{J}_{\text{Pre-LN}} = (\mathbf{J}_{\text{FFN}}\mathbf{J}_{\text{LN}} + \mathbf{I}_{nd}) (\mathbf{J}_{\text{ATT}}\mathbf{J}_{\text{LN}} + \mathbf{I}_{nd}) \quad (5)$$

$$\mathbf{J}_{\text{Post-LN}} = \mathbf{J}_{\text{LN}} (\mathbf{J}_{\text{FFN}} + \mathbf{I}_{nd}) \mathbf{J}_{\text{LN}} (\mathbf{J}_{\text{ATT}} + \mathbf{I}_{nd}), \quad (6)$$

where  $\mathbf{J}_{\text{ATT}}$  and  $\mathbf{J}_{\text{FFN}}$  denote the Jacobians of the self-attention and feed-forward network modules, respectively. For simplicity, the evaluation points of the Jacobians are omitted. The Jacobian of the layer normalization is represented by  $\mathbf{J}_{\text{LN}}$ , calculated for an input  $\mathbf{X} \in \mathbb{R}^{n \times d}$  as:

$$\mathbf{J}_{\text{LN}}(\mathbf{X}) = \text{blockdiag}(\{\mathbf{L}_i(\mathbf{X})\}_{i=1}^n), \quad (7)$$

where each block  $\mathbf{L}_i \in \mathbb{R}^{d \times d}$  is defined as:

$$\mathbf{L}_i(\mathbf{X}) := \frac{\sqrt{d}}{\|\widetilde{\mathbf{X}}_{i,:}\|_2} \left( \mathbf{I}_d - \frac{\widetilde{\mathbf{X}}_{i,:}\widetilde{\mathbf{X}}_{i,:}^\top}{\|\widetilde{\mathbf{X}}_{i,:}\|_2^2} \right) \left( \mathbf{I}_d - \frac{\mathbf{1}\mathbf{1}^\top}{d} \right),$$

and  $\widetilde{\mathbf{X}}_{i,:} := \mathbf{X}_{i,:}(\mathbf{I}_d - \frac{\mathbf{1}\mathbf{1}^\top}{d})$ . These derivations are provided in Appendix B.

**Greater gradient heterogeneity in Post-LN.** Equation (7) shows that the Jacobian of layer normalization,  $\mathbf{J}_{\text{LN}}$ , depends on the input, causing variations in its scale across layers. From Eqs. (5) and (6), we observe that in Post-LN,  $\mathbf{J}_{\text{LN}}$  appears in a multiplicative form and thus proportionally scales the entire Jacobian, leading to greater gradient heterogeneity across layers than in Pre-LN. Further discussion of gradient heterogeneity in Transformers, particularly in the attention mechanism, is provided in Appendix F.

## 5 Numerical evaluation

We numerically evaluate the following claims.

- Gradient heterogeneity is pronounced in Transformers and is influenced by the position of layer normalization (Section 5.2).
- SGD encounters greater difficulty in optimization under gradient heterogeneity compared with adaptive optimizers such as Adam (Section 5.3).

We provide details of the experimental setup and figures in Appendix C and additional results in Appendix D.

### 5.1 Experimental setup

**Datasets and models.** We used a total of nine datasets and three pre-trained models obtained from public sources. For NLP tasks, we used four datasets from SuperGLUE (Wang et al., 2019a) (BoolQ, CB, RTE, and WiC) and three datasets from GLUE (Wang et al., 2018) (CoLA, MRPC, and SST-2) with the RoBERTa-Base model (Liu et al., 2019). For vision tasks, we used the Flowers102 (Nilsback & Zisserman, 2008) and FGVC-Aircraft (Aircraft) (Maji et al., 2013) datasets with ViT-Base (Dosovitskiy et al., 2021) and ResNet18 (He et al., 2016).

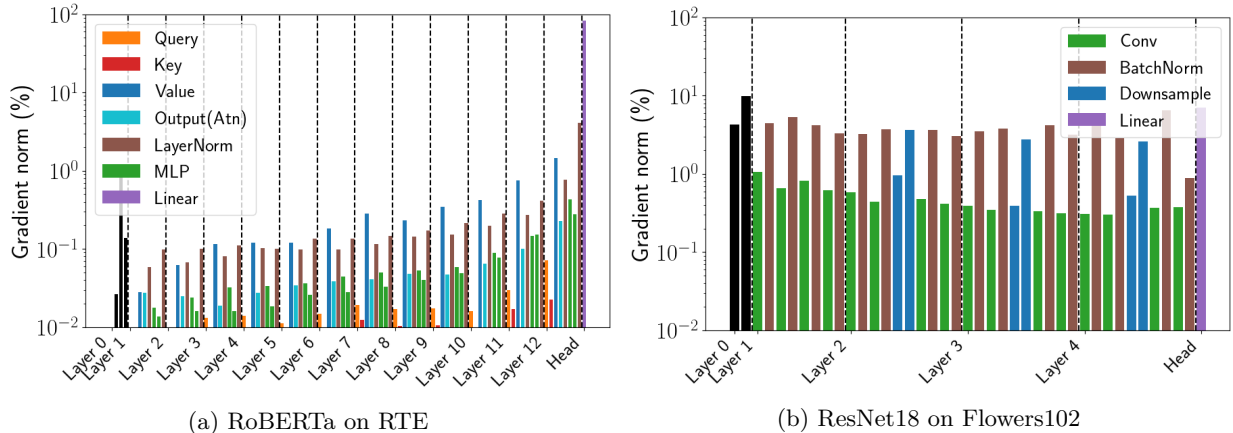


Figure 3: **Transformers exhibit large gradient heterogeneity.** Gradient norms of individual parameters in pre-trained models.

Table 2: **Post-LN increases gradient heterogeneity.** Higher Gini coefficient indicates greater heterogeneity.

| Norm Type | Model state  | Gini Coeff.       |
|-----------|--------------|-------------------|
| Pre-LN    | Random init. | $0.880 \pm 0.004$ |
| Post-LN   | Random init. | $0.941 \pm 0.012$ |
| Post-LN   | Pre-trained  | $0.944 \pm 0.005$ |

**Training.** We compared optimizers with momentum (default) and without momentum (w/o M), as well as SignSGD with  $\ell_1$ -scaled learning rates (SignSGD (S); Section 4.6). Gradient clipping was applied, and the learning-rate schedule was fixed within each domain. All models were fine-tuned from pre-trained weights.

## 5.2 Gradient heterogeneity

Figure 3 shows that RoBERTa exhibits the highest gradient heterogeneity among the models, followed by ViT and ResNet18, indicating that Transformers have more pronounced gradient heterogeneity. In RoBERTa, gradients are smaller near input layers compared to output layers, consistent with our analysis in Section 4.7.

**Effect of layer normalization.** Table 2 shows Gini coefficients for different normalization placements in RoBERTa on RTE. Post-LN shows higher heterogeneity than Pre-LN, consistent with our analysis (Section 4.7). Pre-trained weights are only available for Post-LN.

## 5.3 Training curves

**Limitations of SGD under gradient heterogeneity.** As shown in Figure 4, all optimizers successfully train ViT (and ResNet; see Figure S.7), whereas SGD fails to optimize RoBERTa, highlighting the challenge posed by gradient heterogeneity. In contrast, SignSGD (S) reliably trains both ViT and RoBERTa. These observations are consistent with our theoretical analysis in Theorems 4.6 and 4.8. Additionally, the final training losses are similar between momentum and no-momentum variants of SGD and SignSGD, and Adam performs similarly to RMSProp, suggesting that adaptive learning rates, rather than momentum, are the primary cause of the performance gap (Kunstner et al., 2023). Note that the RTE dataset, which has two almost balanced classes, rules out the heavy-tailed class imbalance as an explanation for the Adam–SGD gap (Kunstner et al., 2024).

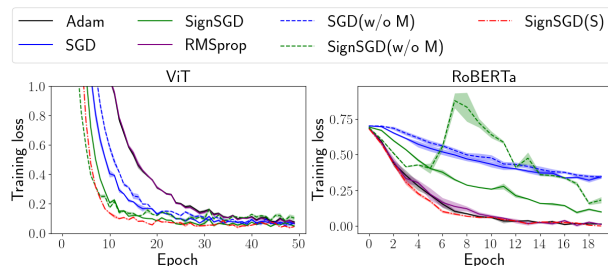


Figure 4: **RoBERTa is difficult to optimize with SGD.** Training loss curves for ViT on Flowers102 (left) and RoBERTa on RTE (right). Shaded regions denote interquartile ranges.

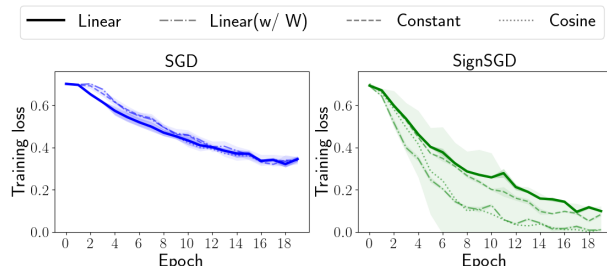


Figure 5: **Learning-rate schedulers improve SignSGD but have limited effect on SGD.** Training loss curves for RoBERTa on RTE. “w/ W” denotes with warmup.

**Effectiveness of learning-rate schedulers.** In NLP tasks, we use linear learning-rate scheduling by default. To assess whether SGD’s poor performance is due to its schedule, we train RoBERTa with various schedules. In Figure 5, learning-rate schedules do not improve SGD, while SignSGD benefits significantly from the appropriate schedules, achieving performance comparable to Adam.

## 6 Conclusion

We derive upper bounds on the iteration complexity of gradient-based and sign-based optimization methods. Our results suggest gradient and Hessian heterogeneity as key factors underlying the performance gap between Adam and SGD in Transformer models. Our analysis leverages the fact that SGD and SignSGD correspond to steepest descent under different norms, yielding implications for learning rate scaling of SignSGD. We further show that gradient heterogeneity is particularly pronounced in Post-LN Transformers. Empirically, SGD degrades under large gradient heterogeneity, whereas SignSGD, with appropriate learning-rate scheduling, achieves performance comparable to Adam.

## References

- Kwangjun Ahn, Xiang Cheng, Minhak Song, Chulhee Yun, Ali Jadbabaie, and Suvrit Sra. Linear attention is (maybe) all you need (to understand transformer optimization). In *International Conference on Learning Representations*, 2024.
- Joshua Ainslie, Santiago Ontanon, Chris Alberti, Vaclav Cvicek, Zachary Fisher, Philip Pham, Anirudh Ravula, Sumit Sanghai, Qifan Wang, and Li Yang. Etc: Encoding long and structured inputs in transformers. In *Empirical Methods in Natural Language Processing*, 2020.
- Lukas Balles and Philipp Hennig. Dissecting adam: The sign, magnitude and variance of stochastic gradients. In *International Conference on Machine Learning*, volume 80. PMLR, 2018.
- Lukas Balles, Fabian Pedregosa, and Nicolas Le Roux. The geometry of sign gradient descent. *arXiv preprint arXiv:2002.08056*, 2020.
- Han Bao, Ryuichiro Hataya, and Ryo Karakida. Self-attention networks localize when qk-eigenspectrum concentrates. In *International Conference on Machine Learning*, 2024.
- Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*, 2020.
- Jeremy Bernstein and Laker Newhouse. Old optimizer, new norm: An anthology. *arXiv preprint arXiv:2409.20325*, 2024.

- Jeremy Bernstein, Yu-Xiang Wang, Kamyar Azizzadenesheli, and Animashree Anandkumar. signsgd: Compressed optimisation for non-convex problems. In *International Conference on Machine Learning*. PMLR, 2018.
- Dimitri P. Bertsekas. *Nonlinear Programming*. Athena Scientific, Belmont, MA, 2 edition, 1999.
- Stephen Boyd and Lieven Vandenbergh. *Convex Optimization*. Cambridge University Press, 2004.
- David Carlson, Ya-Ping Hsieh, Edo Collins, Lawrence Carin, and Volkan Cevher. Stochastic spectral descent for discrete graphical models. *IEEE Journal of Selected Topics in Signal Processing*, 10(2):296–311, 2015.
- Yair Carmon, John C Duchi, Oliver Hinder, and Aaron Sidford. Lower bounds for finding stationary points i. *Mathematical Programming*, 184(1), 2020.
- Guanzheng Chen, Fangyu Liu, Zaiqiao Meng, and Shangsong Liang. Revisiting parameter-efficient tuning: Are we really there yet? In *Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 2022.
- Lizhang Chen, Bo Liu, Kaizhao Liang, and qiang liu. Lion secretly solves a constrained optimization: As lyapunov predicts. In *International Conference on Learning Representations*, 2024a.
- Xiangning Chen, Chen Liang, Da Huang, Esteban Real, Kaiyuan Wang, Hieu Pham, Xuanyi Dong, Thang Luong, Cho-Jui Hsieh, Yifeng Lu, et al. Symbolic discovery of optimization algorithms. *Advances in Neural Information Processing Systems*, 36, 2024b.
- D Choi. On empirical comparisons of optimizers for deep learning. *arXiv preprint arXiv:1910.05446*, 2019.
- Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. BoolQ: Exploring the surprising difficulty of natural yes/no questions. In *North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Association for Computational Linguistics, 2019.
- Kevin Clark. What does bert look at? an analysis of bert’s attention. *arXiv preprint arXiv:1906.04341*, 2019.
- Ronan Collobert. *Large scale machine learning*. PhD thesis, Université de Paris VI, 2004.
- Michael Crawshaw, Mingrui Liu, Francesco Orabona, Wei Zhang, and Zhenxun Zhuang. Robustness to unbounded smoothness of generalized signsgd. *Advances in Neural Information Processing Systems*, 35, 2022.
- Wanyun Cui and Qianle Wang. Cherry on top: Parameter heterogeneity and quantization in large language models. *arXiv preprint arXiv:2404.02837*, 2024.
- Marie-Catherine De Marneffe, Mandy Simons, and Judith Tonhauser. The commitmentbank: Investigating projection in naturally occurring discourse. In *Sinn und Bedeutung*, volume 23, 2019.
- Bill Dolan and Chris Brockett. Automatically constructing a corpus of sentential paraphrases. In *Third international workshop on paraphrasing (IWP2005)*, 2005.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *The IEEE conference on computer vision and pattern recognition*, 2016.

- Ruining He, Anirudh Ravula, Bhargav Kanagal, and Joshua Ainslie. Realformer: Transformer likes residual attention. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, 2021.
- Nam Hyeon-Woo, Kim Yu-Ji, Byeongho Heo, Dongyoon Han, Seong Joon Oh, and Tae-Hyun Oh. Scratching visual transformer’s back with uniform attention. In *The IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2023.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- Kaiqi Jiang, Dhruv Malik, and Yuanzhi Li. How does adaptive optimization impact local neural network geometry? *Advances in Neural Information Processing Systems*, 36, 2024.
- Sai Praneeth Karimireddy, Quentin Rebjock, Sebastian Stich, and Martin Jaggi. Error feedback fixes signsgd and other gradient compression schemes. In *International conference on machine learning*, pp. 3252–3261. PMLR, 2019.
- Jonathan A Kelner, Yin Tat Lee, Lorenzo Orecchia, and Aaron Sidford. An almost-linear-time algorithm for approximate max flow in undirected graphs, and its multicommodity generalizations. In *Proceedings of the twenty-fifth annual ACM-SIAM symposium on Discrete algorithms*, pp. 217–226. SIAM, 2014.
- Diederick P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015.
- Frederik Kunstner, Jacques Chen, Jonathan Wilder Lavington, and Mark Schmidt. Noise is not the main factor behind the gap between sgd and adam on transformers, but sign descent might be. In *International Conference on Learning Representations*, 2023.
- Frederik Kunstner, Robin Yadav, Alan Milligan, Mark Schmidt, and Alberto Bietti. Heavy-tailed class imbalance and why adam outperforms gradient descent on language models. In *Advances in Neural Information Processing Systems*, 2024.
- Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Bingrui Li, Wei Huang, Andi Han, Zhanpeng Zhou, Taiji Suzuki, Jun Zhu, and Jianfei Chen. On the optimization and generalization of two-layer transformers with sign gradient descent. In *International Conference on Learning Representations*, 2025.
- Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- Lucas Maes, Tianyue H Zhang, Alexia Jolicoeur-Martineau, Ioannis Mitliagkas, Damien Scieur, Simon Lacoste-Julien, and Charles Guille-Escuret. Understanding adam requires better rotation dependent assumptions. *arXiv preprint arXiv:2410.19964*, 2024.
- Subhansu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-grained visual classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013.
- James Martens and Roger Grosse. Optimizing neural networks with kronecker-factored approximate curvature. In *International conference on machine learning*. PMLR, 2015.
- Yurii Nesterov. *Introductory lectures on convex optimization: A basic course*, volume 87. Springer Science & Business Media, 2013.

- Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In *Indian conference on computer vision, graphics & image processing*. IEEE, 2008.
- Lorenzo Noci, Sotiris Anagnostidis, Luca Biggio, Antonio Orvieto, Sidak Pal Singh, and Aurelien Lucchi. Signal propagation in transformers: Theoretical perspectives and the role of rank collapse. *Advances in Neural Information Processing Systems*, 35, 2022.
- Weronika Ormaniec, Felix Dangel, and Sidak Pal Singh. What does it mean to be a transformer? insights from a theoretical hessian analysis. In *International Conference on Learning Representations*, 2025.
- Antonio Orvieto and Robert Gower. In search of adam’s secret sauce. *Advances in Neural Information Processing Systems*, 38, 2025.
- Yan Pan and Yuanzhi Li. Toward understanding why adam converges faster than SGD for transformers. In *OPT 2022: Optimization for Machine Learning (NeurIPS 2022 Workshop)*, 2022.
- Namuk Park and Songkuk Kim. How do vision transformers work? In *International Conference on Learning Representations*, 2022.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library, 2019.
- Mohammad Taher Pilehvar and Jose Camacho-Collados. Wic: the word-in-context dataset for evaluating context-sensitive meaning representations. In *North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2019.
- Andrew Rabinovich, Andrea Vedaldi, Carolina Galleguillos, Eric Wiewiora, and Serge Belongie. Objects in context. In *International Conference on Computer Vision*. IEEE, 2007.
- Elan Rosenfeld and Andrej Risteski. Outliers with opposing signals have an outsized effect on neural network optimization. In *International Conference on Learning Representations*, 2024.
- Robin M Schmidt, Frank Schneider, and Philipp Hennig. Descending through a crowded valley-benchmarking deep learning optimizers. In *International Conference on Machine Learning*. PMLR, 2021.
- Han Shi, JIAHUI GAO, Hang Xu, Xiaodan Liang, Zhenguo Li, Lingpeng Kong, Stephen MS Lee, and James Kwok. Revisiting over-smoothing in bert from the perspective of graph. In *International Conference on Learning Representations*, 2022.
- Jamie Shotton, John Winn, Carsten Rother, and Antonio Criminisi. Textonboost for image understanding: Multi-class object recognition and segmentation by jointly modeling texture, layout, and context. *International journal of computer vision*, 81, 2009.
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *Empirical Methods in Natural Language Processing*, 2013.
- Sho Takase, Shun Kiyono, Sosuke Kobayashi, and Jun Suzuki. On layer normalizations and residual connections in transformers. *arXiv preprint arXiv:2206.00330*, 2022.
- Tijmen Tieleman and G Hinton. Divide the gradient by a running average of its recent magnitude. coursera: Neural networks for machine learning. *Technical report*, 2017.
- Akiyoshi Tomihari and Issei Sato. Understanding linear probing then fine-tuning language models from ntk perspective. In *Advances in Neural Information Processing Systems*, volume 37, 2024.
- Antonio Torralba. Contextual priming for object detection. *International journal of computer vision*, 53, 2003.

- A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. Glue: A multi-task benchmark and analysis platform for natural language understanding. In *International Conference on Learning Representations*, 2018.
- Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. Superglue: A stickier benchmark for general-purpose language understanding systems. *Advances in Neural Information Processing Systems*, 32, 2019a.
- Qiang Wang, Bei Li, Tong Xiao, Jingbo Zhu, Changliang Li, Derek F Wong, and Lidia S Chao. Learning deep transformer models for machine translation. In *Annual Meeting of the Association for Computational Linguistics*, 2019b.
- Shulun Wang, Feng Liu, and Bin Liu. Escaping the gradient vanishing: Periodic alternatives of softmax in attention mechanism. *IEEE Access*, 9, 2021.
- Alex Warstadt, Amanpreet Singh, and Samuel R. Bowman. Neural network acceptability judgments. *Transactions of the Association for Computational Linguistics*, 7, 2019.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-the-art natural language processing. In *Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 2020.
- Xinyi Wu, Amir Ajorlou, Yifei Wang, Stefanie Jegelka, and Ali Jadbabaie. On the role of attention masks and layernorm in transformers. In *Advances in Neural Information Processing Systems*, 2024.
- Shuo Xie and Zhiyuan Li. Implicit bias of adamw:  $\ell_\infty$ -norm constrained optimization. In *International Conference on Machine Learning*. PMLR, 2024.
- Ruibin Xiong, Yunchang Yang, Di He, Kai Zheng, Shuxin Zheng, Chen Xing, Huishuai Zhang, Yanyan Lan, Liwei Wang, and Tiejun Liu. On layer normalization in the transformer architecture. In *International Conference on Machine Learning*. PMLR, 2020.
- Zhewei Yao, Amir Gholami, Kurt Keutzer, and Michael W Mahoney. Pyhessian: Neural networks through the lens of the hessian. In *IEEE international conference on big data (Big data)*. IEEE, 2020.
- Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. Big bird: Transformers for longer sequences. *Advances in Neural Information Processing Systems*, 33, 2020.
- Shuangfei Zhai, Tatiana Likhomanenko, Etai Littwin, Dan Busbridge, Jason Ramapuram, Yizhe Zhang, Jiatao Gu, and Joshua M. Susskind. Stabilizing transformer training by preventing attention entropy collapse. In *International Conference on Machine Learning*, volume 202. PMLR, 2023.
- Huishuai Zhang, Caiming Xiong, James Bradbury, and Richard Socher. Block-diagonal hessian-free optimization for training neural networks. *arXiv preprint arXiv:1712.07296*, 2017.
- Jingzhao Zhang, Tianxing He, Suvrit Sra, and Ali Jadbabaie. Why gradient clipping accelerates training: A theoretical justification for adaptivity. In *International Conference on Learning Representations*, 2020a.
- Jingzhao Zhang, Sai Praneeth Karimireddy, Andreas Veit, Seungyeon Kim, Sashank Reddi, Sanjiv Kumar, and Suvrit Sra. Why are adaptive methods good for attention models? *Advances in Neural Information Processing Systems*, 33, 2020b.
- Yushun Zhang, Congliang Chen, Naichen Shi, Ruoyu Sun, and Zhi-Quan Luo. Adam can converge without any modification on update rules. *Advances in Neural Information Processing Systems*, 35, 2022.

Yushun Zhang, Congliang Chen, Tian Ding, Ziniu Li, Ruoyu Sun, and Zhiquan Luo. Why transformers need adam: A hessian perspective. *Advances in Neural Information Processing Systems*, 37, 2024a.

Yushun Zhang, Congliang Chen, Ziniu Li, Tian Ding, Chenwei Wu, Diederik P. Kingma, Yinyu Ye, Zhi-Quan Luo, and Ruoyu Sun. Adam-mini: Use fewer learning rates to gain more, 2024b.

Rosie Zhao, Depen Morwani, David Brandfonbrener, Nikhil Vyas, and Sham M. Kakade. Deconstructing what makes a good optimizer for autoregressive language models. In *International Conference on Learning Representations*, 2025.

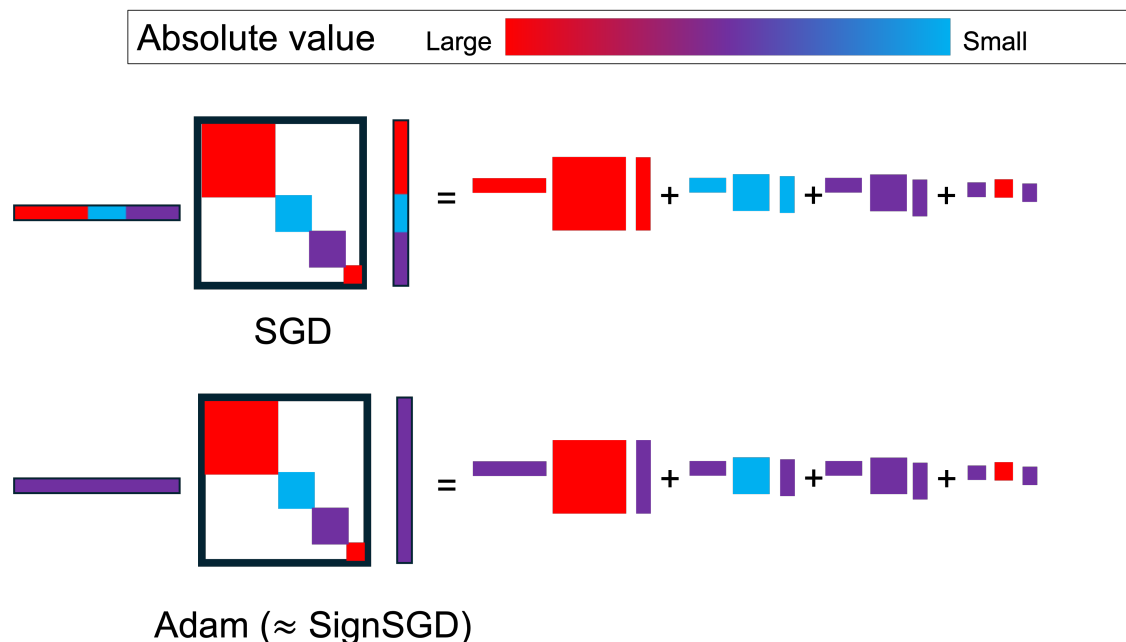


Figure S.1: Proof intuition for Theorem 4.6. This figure illustrates the key quantity  $|\Delta_t^\top \nabla^2 L_D(\theta_t) \Delta_t|$  discussed in Appendix A.1. Under gradient–Hessian correlation, gradient-based sequences (top) align large gradient norms with large Hessian operator norms, amplifying the block-diagonal quadratic form. Sign-based updates (bottom), which use unit-magnitude directions, mitigate this effect and yield more uniform contributions across blocks.

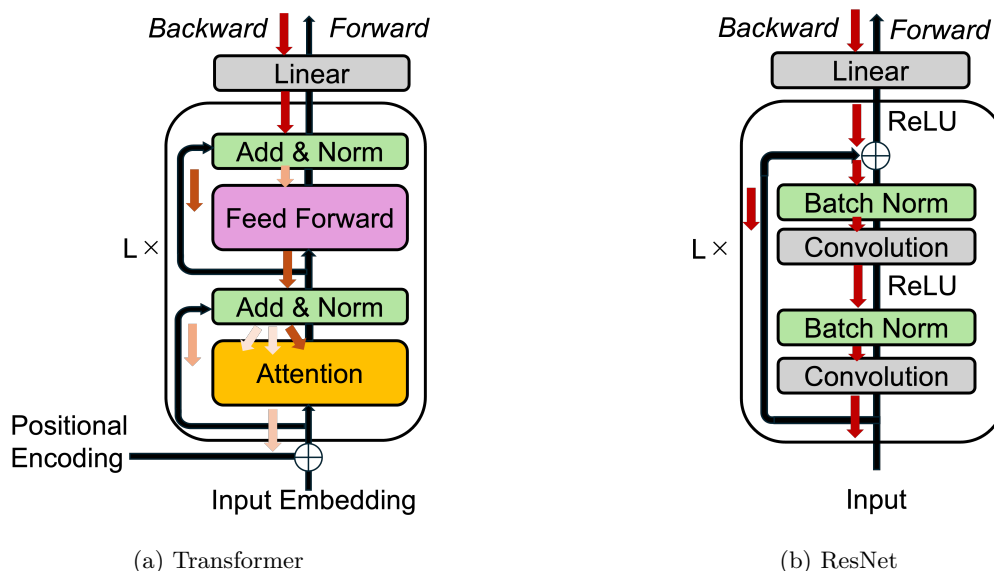


Figure S.2: Architecture and gradient heterogeneity across models. ResNets, as a representative CNN architecture, are constructed by the repetitive stacking of homogeneous parameter blocks (convolutional layers), which promotes relatively uniform gradient propagation. In contrast, Transformers involve stacking of heterogeneous parameter blocks, such as the Query, Key, Value, and output projection layers in attention as well as MLP layers, leading to uneven gradient propagation across modules and pronounced gradient heterogeneity.

## A Proof

### A.1 Proof intuition

Here, we outline the core analysis underlying Theorem 4.6. We apply the descent lemma (Bertsekas, 1999) in our setting and construct an upper bound for the term

$$L(\boldsymbol{\theta}_{t+1}) - L(\boldsymbol{\theta}_t).$$

Under Assumption 4.2, the key quantity to be controlled is

$$|\Delta_t^\top \nabla^2 L(\boldsymbol{\theta}_t) \Delta_t|,$$

where

$$\Delta_t := \begin{cases} \nabla L(\boldsymbol{\theta}_t^{\text{Grad}}) & \text{(gradient-based sequence)} \\ \text{sign}(\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})) & \text{(sign-based sequence)} \end{cases}$$

is the update term without learning rate.

With Assumption 4.3, the dominant contribution to this quadratic form can be characterized by the block-diagonal component

$$|\Delta_t^\top \nabla^2 L_D(\boldsymbol{\theta}_t) \Delta_t|,$$

while the contribution of the off-diagonal blocks is controlled by the approximation error  $\delta_D$ .

For the gradient-based sequence,  $\Delta_t$  coincides with the gradient itself, so parameter blocks with large gradient norms are naturally aligned with Hessian blocks that have large operator norms, resulting in a large contribution to the quadratic form. In contrast, for the sign-based sequence, each coordinate of  $\Delta_t$  has unit magnitude, which suppresses such alignment and prevents large contributions from individual parameter blocks. These effects are illustrated in Figure S.1. The complete proof is provided in Appendix A.3.

### A.2 Technical lemma

**Lemma A.1.** *Under assumption 4.2, for any  $\boldsymbol{\theta}, \boldsymbol{\theta}' \in \mathbb{R}^P$ , the following inequality holds:*

$$L(\boldsymbol{\theta}') - L(\boldsymbol{\theta}) \leq \nabla L(\boldsymbol{\theta})^\top (\boldsymbol{\theta}' - \boldsymbol{\theta}) + \frac{1}{2} (\boldsymbol{\theta}' - \boldsymbol{\theta})^\top \nabla^2 L(\boldsymbol{\theta}) (\boldsymbol{\theta}' - \boldsymbol{\theta}) + \frac{\rho_H}{6} \|\boldsymbol{\theta}' - \boldsymbol{\theta}\|_2^3.$$

*Proof.* Define  $\boldsymbol{\nu}(\alpha) := \boldsymbol{\theta} + \alpha(\boldsymbol{\theta}' - \boldsymbol{\theta})$ . By the  $\rho_H$ -Lipschitz continuity of the Hessian (assumption 4.2), we have:

$$\begin{aligned} & (\nabla L(\boldsymbol{\theta}') - \nabla L(\boldsymbol{\theta}))^\top (\boldsymbol{\theta}' - \boldsymbol{\theta}) \\ &= \int_0^1 (\boldsymbol{\theta}' - \boldsymbol{\theta})^\top \nabla^2 L(\boldsymbol{\nu}(\alpha)) (\boldsymbol{\theta}' - \boldsymbol{\theta}) d\alpha \\ &= (\boldsymbol{\theta}' - \boldsymbol{\theta})^\top \nabla^2 L(\boldsymbol{\theta}) (\boldsymbol{\theta}' - \boldsymbol{\theta}) + \int_0^1 (\boldsymbol{\theta}' - \boldsymbol{\theta})^\top (\nabla^2 L(\boldsymbol{\nu}(\alpha)) - \nabla^2 L(\boldsymbol{\theta})) (\boldsymbol{\theta}' - \boldsymbol{\theta}) d\alpha \\ &\leq (\boldsymbol{\theta}' - \boldsymbol{\theta})^\top \nabla^2 L(\boldsymbol{\theta}) (\boldsymbol{\theta}' - \boldsymbol{\theta}) + \int_0^1 \|\nabla^2 L(\boldsymbol{\nu}(\alpha)) - \nabla^2 L(\boldsymbol{\theta})\|_2 \|\boldsymbol{\theta}' - \boldsymbol{\theta}\|_2^2 d\alpha \\ &\leq (\boldsymbol{\theta}' - \boldsymbol{\theta})^\top \nabla^2 L(\boldsymbol{\theta}) (\boldsymbol{\theta}' - \boldsymbol{\theta}) + \int_0^1 \rho_H \alpha \|\boldsymbol{\theta}' - \boldsymbol{\theta}\|_2^3 d\alpha \\ &= (\boldsymbol{\theta}' - \boldsymbol{\theta})^\top \nabla^2 L(\boldsymbol{\theta}) (\boldsymbol{\theta}' - \boldsymbol{\theta}) + \frac{\rho_H}{2} \|\boldsymbol{\theta}' - \boldsymbol{\theta}\|_2^3. \end{aligned}$$

Using this inequality, we obtain:

$$L(\boldsymbol{\theta}') - L(\boldsymbol{\theta})$$

$$\begin{aligned}
&= \int_0^1 \nabla L(\boldsymbol{\nu}(\alpha))^\top (\boldsymbol{\theta}' - \boldsymbol{\theta}) d\alpha \\
&= \nabla L(\boldsymbol{\theta})^\top (\boldsymbol{\theta}' - \boldsymbol{\theta}) + \int_0^1 (\nabla L(\boldsymbol{\nu}(\alpha)) - \nabla L(\boldsymbol{\theta}))^\top (\boldsymbol{\theta}' - \boldsymbol{\theta}) d\alpha \\
&= \nabla L(\boldsymbol{\theta})^\top (\boldsymbol{\theta}' - \boldsymbol{\theta}) + \int_0^1 (\nabla L(\boldsymbol{\nu}(\alpha)) - \nabla L(\boldsymbol{\theta}))^\top \frac{1}{\alpha} (\boldsymbol{\nu}(\alpha) - \boldsymbol{\theta}) d\alpha \\
&\leq \nabla L(\boldsymbol{\theta})^\top (\boldsymbol{\theta}' - \boldsymbol{\theta}) + \int_0^1 \frac{1}{\alpha} \left( (\boldsymbol{\nu}(\alpha) - \boldsymbol{\theta})^\top \nabla^2 L(\boldsymbol{\theta}) (\boldsymbol{\nu}(\alpha) - \boldsymbol{\theta}) + \frac{\rho_H}{2} \|\boldsymbol{\nu}(\alpha) - \boldsymbol{\theta}\|_2^3 \right) d\alpha \\
&= \nabla L(\boldsymbol{\theta})^\top (\boldsymbol{\theta}' - \boldsymbol{\theta}) + \int_0^1 \left( (\boldsymbol{\theta}' - \boldsymbol{\theta})^\top \nabla^2 L(\boldsymbol{\theta}) (\boldsymbol{\theta}' - \boldsymbol{\theta}) \alpha + \frac{\rho_H}{2} \|\boldsymbol{\theta}' - \boldsymbol{\theta}\|_2^3 \alpha^2 \right) d\alpha \\
&= \nabla L(\boldsymbol{\theta})^\top (\boldsymbol{\theta}' - \boldsymbol{\theta}) + \frac{1}{2} (\boldsymbol{\theta}' - \boldsymbol{\theta})^\top \nabla^2 L(\boldsymbol{\theta}) (\boldsymbol{\theta}' - \boldsymbol{\theta}) + \frac{\rho_H}{6} \|\boldsymbol{\theta}' - \boldsymbol{\theta}\|_2^3.
\end{aligned}$$

□

**Lemma A.2.** For any  $a, b \geq 0$ , the following inequality holds:

$$(a + b)^3 \leq 4(a^3 + b^3).$$

*Proof.* Calculating the difference between the right-hand and left-hand side, we obtain:

$$\begin{aligned}
4(a^3 + b^3) - (a + b)^3 &= 4(a^3 + b^3) - (a^3 + 3a^2b + 3ab^2 + b^3) \\
&= 3(a^3 + b^3) - 3a^2b - 3ab^2 \\
&= 3(a + b)(a - b)^2 \geq 0.
\end{aligned}$$

□

### A.3 Proof of Theorem 4.6

**Theorem 4.6 is restated.** Assume  $\delta_D < \min(\Lambda_G, \Lambda_P)/3$ . Then, the iteration complexities in the deterministic setting are bounded as follows.

For the gradient-based sequence, suppose that  $\varepsilon < \frac{\Lambda_G^2}{\rho_H \sqrt{P}}$  holds and that the learning rate at time  $t$  satisfies  $\eta_t = \zeta_t \min(\frac{1}{\Lambda_G}, \frac{1}{\sqrt{\rho_H \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2}})$ , where  $\zeta_t \in [\zeta_0, 1]$ , we have

$$\mathcal{T}_\varepsilon(\{\boldsymbol{\theta}_t^{\text{Grad}}\}_{t=0}^\infty, L, \|\cdot\|_2) \leq \frac{6(L(\boldsymbol{\theta}_0) - L_*)}{P\varepsilon^2\zeta_0} \Lambda_G.$$

For the sign-based sequence, suppose that  $\varepsilon < \frac{\Lambda_P^2}{\rho_H \sqrt{P}}$  holds and that the learning rate at time  $t$  satisfies

$\eta_t = \zeta_t \min(\frac{\|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1}{\Lambda_P P}, \sqrt{\frac{\|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1}{\rho_H P^{3/2}}})$ , where  $\zeta_t \in [\zeta_0, 1]$ , we have

$$\mathcal{T}_\varepsilon(\{\boldsymbol{\theta}_t^{\text{Sign}}\}_{t=0}^\infty, L, \|\cdot\|_1) \leq \frac{6(L(\boldsymbol{\theta}_0) - L_*)}{P\varepsilon^2\zeta_0} \Lambda_P.$$

*Proof of gradient-based sequence.* The update rule of the gradient-based sequence in deterministic setting is  $\boldsymbol{\theta}_{t+1}^{\text{Grad}} = \boldsymbol{\theta}_t^{\text{Grad}} - \eta_t \nabla L(\boldsymbol{\theta}_t^{\text{Grad}})$ . Thus, we obtain:

$$\begin{aligned}
&L(\boldsymbol{\theta}_{t+1}^{\text{Grad}}) - L(\boldsymbol{\theta}_t^{\text{Grad}}) \\
&\leq \nabla L(\boldsymbol{\theta}_t^{\text{Grad}})^\top (\boldsymbol{\theta}_{t+1}^{\text{Grad}} - \boldsymbol{\theta}_t^{\text{Grad}}) + \frac{1}{2} (\boldsymbol{\theta}_{t+1}^{\text{Grad}} - \boldsymbol{\theta}_t^{\text{Grad}})^\top \nabla^2 L(\boldsymbol{\theta}_t^{\text{Grad}}) (\boldsymbol{\theta}_{t+1}^{\text{Grad}} - \boldsymbol{\theta}_t^{\text{Grad}}) + \frac{\rho_H}{6} \|\boldsymbol{\theta}_{t+1}^{\text{Grad}} - \boldsymbol{\theta}_t^{\text{Grad}}\|_2^3 \\
&= -\eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 + \frac{\eta_t^2}{2} \nabla L(\boldsymbol{\theta}_t^{\text{Grad}})^\top \nabla^2 L(\boldsymbol{\theta}_t^{\text{Grad}}) \nabla L(\boldsymbol{\theta}_t^{\text{Grad}}) + \eta_t^3 \frac{\rho_H}{6} \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^3
\end{aligned}$$

$$\begin{aligned}
&= -\eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 + \frac{\eta_t^2}{2} \nabla L(\boldsymbol{\theta}_t^{\text{Grad}})^\top \nabla^2 L_D(\boldsymbol{\theta}_t^{\text{Grad}}) \nabla L(\boldsymbol{\theta}_t^{\text{Grad}}) \\
&\quad + \frac{\eta_t^2}{2} \nabla L(\boldsymbol{\theta}_t^{\text{Grad}})^\top (\nabla^2 L(\boldsymbol{\theta}_t^{\text{Grad}}) - \nabla^2 L_D(\boldsymbol{\theta}_t^{\text{Grad}})) \nabla L(\boldsymbol{\theta}_t^{\text{Grad}}) + \eta_t^3 \frac{\rho_H}{6} \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^3 \\
&= -\eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 + \frac{\eta_t^2}{2} \sum_b [\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})]_b^\top [\nabla^2 L(\boldsymbol{\theta}_t^{\text{Grad}})]_b [\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})]_b \\
&\quad + \frac{\eta_t^2}{2} \nabla L(\boldsymbol{\theta}_t^{\text{Grad}})^\top (\nabla^2 L(\boldsymbol{\theta}_t^{\text{Grad}}) - \nabla^2 L_D(\boldsymbol{\theta}_t^{\text{Grad}})) \nabla L(\boldsymbol{\theta}_t^{\text{Grad}}) + \eta_t^3 \frac{\rho_H}{6} \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^3 \\
&\leq -\eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 + \frac{\eta_t^2}{2} \sum_b \|[\nabla^2 L(\boldsymbol{\theta}_t^{\text{Grad}})]_b\|_2 \|[\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})]_b\|_2^2 \\
&\quad + \frac{\eta_t^2}{2} \|\nabla^2 L(\boldsymbol{\theta}_t^{\text{Grad}}) - \nabla^2 L_D(\boldsymbol{\theta}_t^{\text{Grad}})\|_2 \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 + \eta_t^3 \frac{\rho_H}{6} \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^3 \\
&\leq -\eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 + \frac{\eta_t^2}{2} \Lambda_G \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 + \frac{\eta_t^2}{2} \delta_D \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 + \eta_t^3 \frac{\rho_H}{6} \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^3 \tag{8} \\
&\leq -\eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 + \frac{\eta_t}{2} \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 + \frac{\eta_t}{6} \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 + \frac{\eta_t}{6} \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 \\
&= -\frac{\eta_t}{6} \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2,
\end{aligned}$$

where we use Lemma A.1,  $\eta_t \leq \min\left(\frac{1}{\Lambda_G}, \frac{1}{\sqrt{\rho_H \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2}}\right)$ , and  $\delta_D < \Lambda_G/3$ .

Taking a telescoping sum and noting that  $\boldsymbol{\theta}_0 = \boldsymbol{\theta}_0^{\text{Grad}}$  and  $\eta_t \geq \zeta_0 \min\left(\frac{1}{\Lambda_G}, \frac{1}{\sqrt{\rho_H \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2}}\right)$ , we have:

$$\begin{aligned}
L(\boldsymbol{\theta}_T^{\text{Grad}}) - L(\boldsymbol{\theta}_0) &\leq -\frac{1}{6} \sum_{t=0}^{T-1} \eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 \\
&\leq -\frac{\zeta_0}{6} \sum_{t=0}^{T-1} \min\left(\frac{\|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2}{\Lambda_G}, \frac{\|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^{3/2}}{\sqrt{\rho_H}}\right).
\end{aligned}$$

Assume that  $\|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2 \geq \sqrt{P}\varepsilon$  holds for all  $0 \leq t < T$ . Then, using  $\varepsilon < \frac{\Lambda_G^2}{\rho_H \sqrt{P}}$ , we have:

$$\begin{aligned}
L(\boldsymbol{\theta}_T^{\text{Grad}}) - L(\boldsymbol{\theta}_0) &\leq -\frac{T\zeta_0}{6} \min\left(\frac{P\varepsilon^2}{\Lambda_G}, \frac{P^{3/4}\varepsilon^{3/2}}{\sqrt{\rho_H}}\right) \\
&= -\frac{TP\varepsilon^2\zeta_0}{6\Lambda_G}.
\end{aligned}$$

Therefore, we have

$$\begin{aligned}
T &\leq \frac{6(L(\boldsymbol{\theta}_0) - L(\boldsymbol{\theta}_T^{\text{Grad}}))}{P\varepsilon^2\zeta_0} \Lambda_G \\
&\leq \frac{6(L(\boldsymbol{\theta}_0) - L_*)}{P\varepsilon^2\zeta_0} \Lambda_G.
\end{aligned}$$

This means

$$\mathcal{T}_\varepsilon(\{\boldsymbol{\theta}_t^{\text{Grad}}\}_{t=0}^\infty, L, \|\cdot\|_2) \leq \frac{6(L(\boldsymbol{\theta}_0) - L_*)}{P\varepsilon^2\zeta_0} \Lambda_G.$$

□

*Proof of sign-based sequence.* The update rule of the sign-based sequence in deterministic setting is  $\boldsymbol{\theta}_{t+1}^{\text{Sign}} = \boldsymbol{\theta}_t^{\text{Sign}} - \eta_t \text{sign}(\nabla L(\boldsymbol{\theta}_t^{\text{Sign}}))$ . Thus, we obtain:

$$L(\boldsymbol{\theta}_{t+1}^{\text{Sign}}) - L(\boldsymbol{\theta}_t^{\text{Sign}})$$

$$\begin{aligned}
&\leq \nabla L(\boldsymbol{\theta}_t^{\text{Sign}})^\top (\boldsymbol{\theta}_{t+1}^{\text{Sign}} - \boldsymbol{\theta}_t^{\text{Sign}}) + \frac{1}{2} (\boldsymbol{\theta}_{t+1}^{\text{Sign}} - \boldsymbol{\theta}_t^{\text{Sign}})^\top \nabla^2 L(\boldsymbol{\theta}_t^{\text{Sign}}) (\boldsymbol{\theta}_{t+1}^{\text{Sign}} - \boldsymbol{\theta}_t^{\text{Sign}}) + \frac{\rho_H}{6} \|\boldsymbol{\theta}_{t+1}^{\text{Sign}} - \boldsymbol{\theta}_t^{\text{Sign}}\|_2^3 \\
&= -\eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 + \frac{\eta_t^2}{2} \text{sign}(\nabla L(\boldsymbol{\theta}_t^{\text{Sign}}))^\top \nabla^2 L(\boldsymbol{\theta}_t^{\text{Sign}}) \text{sign}(\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})) + \eta_t^3 \frac{\rho_H}{6} \|\text{sign}(\nabla L(\boldsymbol{\theta}_t^{\text{Sign}}))\|_2^3 \\
&= -\eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 + \frac{\eta_t^2}{2} \text{sign}(\nabla L(\boldsymbol{\theta}_t^{\text{Sign}}))^\top \nabla^2 L_D(\boldsymbol{\theta}_t^{\text{Sign}}) \text{sign}(\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})) \\
&\quad + \frac{\eta_t^2}{2} \text{sign}(\nabla L(\boldsymbol{\theta}_t^{\text{Sign}}))^\top (\nabla^2 L(\boldsymbol{\theta}_t^{\text{Sign}}) - \nabla^2 L_D(\boldsymbol{\theta}_t^{\text{Sign}})) \text{sign}(\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})) + \eta_t^3 \frac{\rho_H}{6} P^{3/2} \\
&= -\eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 + \frac{\eta_t^2}{2} \sum_b [\text{sign}(\nabla L(\boldsymbol{\theta}_t^{\text{Sign}}))]_b^\top [\nabla^2 L(\boldsymbol{\theta}_t^{\text{Sign}})]_b [\text{sign}(\nabla L(\boldsymbol{\theta}_t^{\text{Sign}}))]_b \\
&\quad + \frac{\eta_t^2}{2} \text{sign}(\nabla L(\boldsymbol{\theta}_t^{\text{Sign}}))^\top (\nabla^2 L(\boldsymbol{\theta}_t^{\text{Sign}}) - \nabla^2 L_D(\boldsymbol{\theta}_t^{\text{Sign}})) \text{sign}(\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})) + \eta_t^3 \frac{\rho_H}{6} P^{3/2} \\
&\leq -\eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 + \frac{\eta_t^2}{2} \sum_b \|[\nabla^2 L(\boldsymbol{\theta}_t^{\text{Sign}})]_b\|_2 P_b + \frac{\eta_t^2}{2} \|\nabla^2 L(\boldsymbol{\theta}_t^{\text{Sign}}) - \nabla^2 L_D(\boldsymbol{\theta}_t^{\text{Sign}})\|_2 P + \eta_t^3 \frac{\rho_H}{6} P^{3/2} \\
&\leq -\eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 + \frac{\eta_t^2}{2} \Lambda_P P + \frac{\eta_t^2}{2} \delta_D P + \eta_t^3 \frac{\rho_H}{6} P^{3/2} \tag{9} \\
&\leq -\eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 + \frac{\eta_t}{2} \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 + \frac{\eta_t}{6} \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 + \frac{\eta_t}{6} \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 \\
&= -\frac{\eta_t}{6} \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1,
\end{aligned}$$

where we used Lemma A.1,  $\eta_t \leq \min(\frac{\|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1}{\Lambda_P P}, \sqrt{\frac{\|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1}{\rho_H P^{3/2}}})$ , and  $\delta_D < \Lambda_P/3$ .

Taking the telescoping sum, and noting that  $\boldsymbol{\theta}_0 = \boldsymbol{\theta}_0^{\text{Sign}}$  and  $\eta_t \geq \zeta_0 \min(\frac{\|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1}{\Lambda_P P}, \sqrt{\frac{\|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1}{\rho_H P^{3/2}}})$ , we have:

$$\begin{aligned}
L(\boldsymbol{\theta}_T^{\text{Sign}}) - L(\boldsymbol{\theta}_0) &\leq -\frac{1}{6} \sum_{t=0}^{T-1} \eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 \\
&\leq -\frac{\zeta_0}{6} \sum_{t=0}^{T-1} \min(\frac{\|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1}{P \Lambda_P}, \sqrt{\frac{\|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1}{\rho_H P^{3/2}}}) \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1.
\end{aligned}$$

Assume that  $\|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 \geq P\varepsilon$  holds for all  $0 \leq t < T$ . Then, using  $\varepsilon < \frac{\Lambda_P^2}{\rho_H \sqrt{P}}$ , we have

$$\begin{aligned}
L(\boldsymbol{\theta}_T^{\text{Sign}}) - L(\boldsymbol{\theta}_0) &\leq -\frac{TP\varepsilon\zeta_0}{6} \min(\frac{\varepsilon}{\Lambda_P}, \sqrt{\frac{\varepsilon}{\rho_H P^{1/2}}}) \\
&= -\frac{TP\varepsilon^2\zeta_0}{6\Lambda_P}.
\end{aligned}$$

Therefore, we have:

$$\begin{aligned}
T &\leq \frac{6(L(\boldsymbol{\theta}_0) - L(\boldsymbol{\theta}_T^{\text{Sign}}))}{P\varepsilon^2\zeta_0} \Lambda_P \\
&\leq \frac{6(L(\boldsymbol{\theta}_0) - L_*)}{P\varepsilon^2\zeta_0} \Lambda_P.
\end{aligned}$$

This means:

$$\mathcal{T}_\varepsilon(\{\boldsymbol{\theta}_t^{\text{Sign}}\}_{t=0}^\infty, L, \|\cdot\|_1) \leq \frac{6(L(\boldsymbol{\theta}_0) - L_*)}{P\varepsilon^2\zeta_0} \Lambda_P.$$

□

#### A.4 Proof of Theorem 4.8

**Theorem 4.8 is restated.** Assume  $\delta_D < \min(\Lambda_G, \Lambda_P)/3$ . Then, the iteration complexities the stochastic setting are bounded as follows.

For the gradient-based sequence, suppose that  $\varepsilon < \frac{(1+\sigma_2)^2 \Lambda_G^2}{4(1+\sigma_3)\rho_H \sqrt{P}}$  holds and that the learning rate at time  $t$  satisfies  $\eta_t = \zeta_t \min(\frac{1}{(1+\sigma_2)\Lambda_G}, \frac{1}{2\sqrt{(1+\sigma_3)\rho_H \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2}})$ , where  $\zeta_t \in [\zeta_0, 1]$ , we have

$$\mathcal{T}_\varepsilon(\{\boldsymbol{\theta}_t^{\text{Grad}}\}_{t=0}^\infty, L, \|\cdot\|_2) \leq \frac{12(1+\sigma_2)(L(\boldsymbol{\theta}_0) - L_*)}{P\varepsilon^2\zeta_0} \Lambda_G.$$

For the sign-based sequence, suppose that  $\varepsilon < \frac{\Lambda_P^2}{\rho_H \sqrt{P}}$  and  $\sigma_2 \leq \frac{1}{24}$  hold and that the learning rate at time  $t$  satisfies  $\eta_t = \zeta_t \min(\frac{\|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1}{\Lambda_P P}, \sqrt{\frac{\|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1}{\rho_H P^{3/2}}})$ , where  $\zeta_t \in [\zeta_0, 1]$ , we have

$$\mathcal{T}_\varepsilon(\{\boldsymbol{\theta}_t^{\text{Sign}}\}_{t=0}^\infty, L, \|\cdot\|_1) \leq \frac{12(1+24\sigma_2)(L(\boldsymbol{\theta}_0) - L_*)}{P\varepsilon^2\zeta_0} \Lambda_P.$$

*Proof of gradient-based sequence.* The update rule of the gradient-based sequence in stochastic setting is  $\boldsymbol{\theta}_{t+1}^{\text{Grad}} = \boldsymbol{\theta}_t^{\text{Grad}} - \eta_t \nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Grad}})$ . Thus, using Lemma A.1, we obtain:

$$\begin{aligned} & \mathbb{E} [L(\boldsymbol{\theta}_{t+1}^{\text{Grad}}) - L(\boldsymbol{\theta}_t^{\text{Grad}}) \mid \boldsymbol{\theta}_t^{\text{Grad}}] \\ & \leq \mathbb{E} \left[ \nabla L(\boldsymbol{\theta}_t^{\text{Grad}})^\top (\boldsymbol{\theta}_{t+1}^{\text{Grad}} - \boldsymbol{\theta}_t^{\text{Grad}}) + \frac{1}{2} (\boldsymbol{\theta}_{t+1}^{\text{Grad}} - \boldsymbol{\theta}_t^{\text{Grad}})^\top \nabla^2 L(\boldsymbol{\theta}_t^{\text{Grad}}) (\boldsymbol{\theta}_{t+1}^{\text{Grad}} - \boldsymbol{\theta}_t^{\text{Grad}}) + \frac{\rho_H}{6} \|\boldsymbol{\theta}_{t+1}^{\text{Grad}} - \boldsymbol{\theta}_t^{\text{Grad}}\|_2^3 \mid \boldsymbol{\theta}_t^{\text{Grad}} \right] \\ & = -\eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 + \mathbb{E} \left[ \frac{\eta_t^2}{2} \nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Grad}})^\top \nabla^2 L(\boldsymbol{\theta}_t^{\text{Grad}}) \nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Grad}}) + \eta_t^3 \frac{\rho_H}{6} \|\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^3 \mid \boldsymbol{\theta}_t^{\text{Grad}} \right] \\ & = -\eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 + \mathbb{E} \left[ \frac{\eta_t^2}{2} \nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Grad}})^\top \nabla^2 L_D(\boldsymbol{\theta}_t^{\text{Grad}}) \nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Grad}}) \mid \boldsymbol{\theta}_t^{\text{Grad}} \right] \\ & \quad + \mathbb{E} \left[ \frac{\eta_t^2}{2} \nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Grad}})^\top (\nabla^2 L(\boldsymbol{\theta}_t^{\text{Grad}}) - \nabla^2 L_D(\boldsymbol{\theta}_t^{\text{Grad}})) \nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Grad}}) + \eta_t^3 \frac{\rho_H}{6} \|\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^3 \mid \boldsymbol{\theta}_t^{\text{Grad}} \right] \\ & = -\eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 + \mathbb{E} \left[ \frac{\eta_t^2}{2} \sum_b [\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Grad}})]_b^\top [\nabla^2 L(\boldsymbol{\theta}_t^{\text{Grad}})]_b [\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Grad}})]_b \mid \boldsymbol{\theta}_t^{\text{Grad}} \right] \\ & \quad + \mathbb{E} \left[ \frac{\eta_t^2}{2} \nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Grad}})^\top (\nabla^2 L(\boldsymbol{\theta}_t^{\text{Grad}}) - \nabla^2 L_D(\boldsymbol{\theta}_t^{\text{Grad}})) \nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Grad}}) + \eta_t^3 \frac{\rho_H}{6} \|\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^3 \mid \boldsymbol{\theta}_t^{\text{Grad}} \right] \\ & \leq -\eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 + \mathbb{E} \left[ \frac{\eta_t^2}{2} \sum_b \|\nabla^2 L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2 \|\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 \mid \boldsymbol{\theta}_t^{\text{Grad}} \right] \\ & \quad + \mathbb{E} \left[ \frac{\eta_t^2}{2} \|\nabla^2 L(\boldsymbol{\theta}_t^{\text{Grad}}) - \nabla^2 L_D(\boldsymbol{\theta}_t^{\text{Grad}})\|_2 \|\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 \mid \boldsymbol{\theta}_t^{\text{Grad}} \right] + \mathbb{E} \left[ \eta_t^3 \frac{\rho_H}{6} \|\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^3 \mid \boldsymbol{\theta}_t^{\text{Grad}} \right]. \quad (10) \end{aligned}$$

For the second and third term, using Eqs.(2)(4), and  $\delta_D < \Lambda_G/3$ , we can derive an upper bound as follows:

$$\begin{aligned} & \mathbb{E} \left[ \frac{\eta_t^2}{2} \sum_b \|\nabla^2 L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2 \|\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 \mid \boldsymbol{\theta}_t^{\text{Grad}} \right] + \mathbb{E} \left[ \frac{\eta_t^2}{2} \|\nabla^2 L(\boldsymbol{\theta}_t^{\text{Grad}}) - \nabla^2 L_D(\boldsymbol{\theta}_t^{\text{Grad}})\|_2 \|\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 \mid \boldsymbol{\theta}_t^{\text{Grad}} \right] \\ & \leq \mathbb{E} \left[ \frac{\eta_t^2}{2} \sum_b \|\nabla^2 L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2 \|\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 \mid \boldsymbol{\theta}_t^{\text{Grad}} \right] + \mathbb{E} \left[ \frac{\eta_t^2}{2} \delta_D \|\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 \mid \boldsymbol{\theta}_t^{\text{Grad}} \right] \\ & = \frac{\eta_t^2}{2} \sum_b \|\nabla^2 L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2 \sum_i \mathbb{E} \left[ (([\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})]_b)_i + ([\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Grad}})]_b)_i - ([\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})]_b)_i)^2 \mid \boldsymbol{\theta}_t^{\text{Grad}} \right] \\ & \quad + \frac{\eta_t^2}{2} \delta_D \sum_i \mathbb{E} \left[ (\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})_i + \nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Grad}})_i - \nabla L(\boldsymbol{\theta}_t^{\text{Grad}})_i)^2 \mid \boldsymbol{\theta}_t^{\text{Grad}} \right] \end{aligned}$$

$$\begin{aligned}
&\leq \frac{\eta_t^2}{2} \sum_b \|\nabla^2 L(\boldsymbol{\theta}_t^{\text{Grad}})\|_b (1 + \sigma_2) (\|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_b)_i^2 + \frac{\eta_t^2}{2} \delta_D \sum_i (1 + \sigma_2) \nabla L(\boldsymbol{\theta}_t^{\text{Grad}})_i^2 \\
&\leq \frac{\eta_t^2}{2} (1 + \sigma_2) \Lambda_G \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 + \frac{\eta_t^2}{2} (1 + \sigma_2) \delta_D \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 \\
&\leq \frac{2\eta_t^2}{3} (1 + \sigma_2) \Lambda_G \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 .
\end{aligned} \tag{11}$$

For the fourth term, using Lemma A.2 and Eq.(3), we can derive an upper bound as follows:

$$\begin{aligned}
&\mathbb{E} \left[ \eta_t^3 \frac{\rho_H}{6} \|\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^3 \mid \boldsymbol{\theta}_t^{\text{Grad}} \right] \\
&\leq \eta_t^3 \frac{\rho_H}{6} \mathbb{E} \left[ (\|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2 + \|\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Grad}}) - \nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2)^3 \mid \boldsymbol{\theta}_t^{\text{Grad}} \right] \\
&\leq \frac{2\eta_t^3 \rho_H}{3} \mathbb{E} \left[ \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^3 + \|\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Grad}}) - \nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^3 \mid \boldsymbol{\theta}_t^{\text{Grad}} \right] \\
&\leq \frac{2\eta_t^3 \rho_H}{3} (1 + \sigma_3) \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^3 .
\end{aligned} \tag{12}$$

Combining Eqs.(10)(11) (12) and  $\eta_t \leq \min(\frac{1}{(1+\sigma_2)\Lambda_G}, \frac{1}{2\sqrt{(1+\sigma_3)\rho_H\|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2}})$ , we have:

$$\begin{aligned}
&\mathbb{E} [L(\boldsymbol{\theta}_{t+1}^{\text{Grad}}) - L(\boldsymbol{\theta}_t^{\text{Grad}}) \mid \boldsymbol{\theta}_t^{\text{Grad}}] \\
&\leq -\eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 + \frac{2\eta_t^2}{3} (1 + \sigma_2) \Lambda_G \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 + \frac{2\eta_t^3 \rho_H}{3} (1 + \sigma_3) \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^3 \\
&\leq -\frac{\eta_t}{6} \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2
\end{aligned}$$

Assume that the probability of the event  $\mathcal{E}(T) = \{\forall s \leq T, \|\nabla L(\boldsymbol{\theta}_s^{\text{Grad}})\|_2 \geq \sqrt{P}\varepsilon\}$  satisfies  $\mathbb{P}(\mathcal{E}(T)) \geq \frac{1}{2}$ . By applying the telescoping sum and taking expectations, and noting that  $\boldsymbol{\theta}_0 = \boldsymbol{\theta}_0^{\text{Grad}}$ ,  $\eta_t \geq \zeta_0 \min(\frac{1}{(1+\sigma_2)\Lambda_G}, \frac{1}{2\sqrt{(1+\sigma_3)\rho_H\|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2}})$ , and  $\varepsilon < \frac{(1+\sigma_2)^2 \Lambda_G^2}{4(1+\sigma_3)\rho_H\sqrt{P}}$ , we have:

$$\begin{aligned}
&\mathbb{E} [L(\boldsymbol{\theta}_T^{\text{Grad}})] - L(\boldsymbol{\theta}_0) \\
&\leq -\frac{1}{6} \sum_{t=0}^{T-1} \mathbb{E} [\eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2] \\
&= -\frac{1}{6} \sum_{t=0}^{T-1} \left( \mathbb{E} [\eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 \mid \mathcal{E}(T)] \mathbb{P}(\mathcal{E}(T)) + \mathbb{E} [\eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 \mid \overline{\mathcal{E}(T)}] \mathbb{P}(\overline{\mathcal{E}(T)}) \right) \\
&\leq -\frac{1}{6} \sum_{t=0}^{T-1} \mathbb{E} [\eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 \mid \mathcal{E}(T)] \mathbb{P}(\mathcal{E}(T)) \\
&\leq -\frac{1}{12} \sum_{t=0}^{T-1} \mathbb{E} [\eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 \mid \mathcal{E}(T)] \\
&\leq -\frac{\zeta_0}{12} \sum_{t=0}^{T-1} \mathbb{E} \left[ \min\left(\frac{\|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2}{(1 + \sigma_2)\Lambda_G}, \frac{\|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^{3/2}}{2\sqrt{(1 + \sigma_3)\rho_H}}\right) \mid \mathcal{E}(T) \right] \\
&\leq -\frac{T\zeta_0}{12} \min\left(\frac{P\varepsilon^2}{(1 + \sigma_2)\Lambda_G}, \frac{P^{3/4}\varepsilon^{3/2}}{2\sqrt{(1 + \sigma_3)\rho_H}}\right) \\
&= -\frac{TP\varepsilon^2\zeta_0}{12(1 + \sigma_2)\Lambda_G}.
\end{aligned}$$

Therefore, we have

$$T \leq \frac{12(1 + \sigma_2)(L(\boldsymbol{\theta}_0) - \mathbb{E} [L(\boldsymbol{\theta}_T^{\text{Grad}})])}{P\varepsilon^2\zeta_0} \Lambda_G$$

$$\leq \frac{12(1 + \sigma_2)(L(\boldsymbol{\theta}_0) - L_*)}{P\varepsilon^2\zeta_0} \Lambda_G.$$

This means that when we take  $T > \frac{12(1+\sigma_2)(L(\boldsymbol{\theta}_0)-L_*)}{P\varepsilon^2\zeta_0} \Lambda_G$ , we have  $\mathbb{P}(\mathcal{E}(T)) < \frac{1}{2}$ . Therefore, we have

$$\mathcal{T}_\varepsilon(\{\boldsymbol{\theta}_t^{\text{Grad}}\}_{t=0}^\infty, L, \|\cdot\|_2) \leq \frac{12(1 + \sigma_2)(L(\boldsymbol{\theta}_0) - L_*)}{P\varepsilon^2\zeta_0} \Lambda_G.$$

□

*Proof of sign-based sequence.* The update rule of the sign-based sequence in stochastic setting is  $\boldsymbol{\theta}_{t+1}^{\text{Sign}} = \boldsymbol{\theta}_t^{\text{Sign}} - \eta_t \text{sign}(\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Sign}}))$ . Thus, using Lemma A.1, we obtain:

$$\begin{aligned} & \mathbb{E} \left[ L(\boldsymbol{\theta}_{t+1}^{\text{Sign}}) - L(\boldsymbol{\theta}_t^{\text{Sign}}) \mid \boldsymbol{\theta}_t^{\text{Sign}} \right] \\ & \leq \mathbb{E} \left[ \nabla L(\boldsymbol{\theta}_t^{\text{Sign}})^\top (\boldsymbol{\theta}_{t+1}^{\text{Sign}} - \boldsymbol{\theta}_t^{\text{Sign}}) + \frac{1}{2} (\boldsymbol{\theta}_{t+1}^{\text{Sign}} - \boldsymbol{\theta}_t^{\text{Sign}})^\top \nabla^2 L(\boldsymbol{\theta}_t^{\text{Sign}}) (\boldsymbol{\theta}_{t+1}^{\text{Sign}} - \boldsymbol{\theta}_t^{\text{Sign}}) + \frac{\rho_H}{6} \|\boldsymbol{\theta}_{t+1}^{\text{Sign}} - \boldsymbol{\theta}_t^{\text{Sign}}\|_2^3 \mid \boldsymbol{\theta}_t^{\text{Sign}} \right] \\ & = -\eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 + \mathbb{E} \left[ \frac{\eta_t^2}{2} \text{sign}(\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Sign}}))^\top \nabla^2 L(\boldsymbol{\theta}_t^{\text{Sign}}) \text{sign}(\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Sign}})) + \eta_t^3 \frac{\rho_H}{6} \|\text{sign}(\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Sign}}))\|_2^3 \mid \boldsymbol{\theta}_t^{\text{Sign}} \right] \\ & \quad + \mathbb{E} \left[ -\eta_t \nabla L(\boldsymbol{\theta}_t^{\text{Sign}})^\top (\text{sign}(\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Sign}})) - \text{sign}(\nabla L(\boldsymbol{\theta}_t^{\text{Sign}}))) \mid \boldsymbol{\theta}_t^{\text{Sign}} \right]. \end{aligned} \quad (13)$$

For the second term, we can derive an upper bound in the same way as in the deterministic case:

$$\begin{aligned} & \mathbb{E} \left[ \frac{\eta_t^2}{2} \text{sign}(\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Sign}}))^\top \nabla^2 L(\boldsymbol{\theta}_t^{\text{Sign}}) \text{sign}(\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Sign}})) + \eta_t^3 \frac{\rho_H}{6} \|\text{sign}(\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Sign}}))\|_2^3 \mid \boldsymbol{\theta}_t^{\text{Sign}} \right] \\ & = \mathbb{E} \left[ \frac{\eta_t^2}{2} \text{sign}(\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Sign}}))^\top \nabla^2 L_D(\boldsymbol{\theta}_t^{\text{Sign}}) \text{sign}(\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Sign}})) \mid \boldsymbol{\theta}_t^{\text{Sign}} \right] \\ & \quad + \mathbb{E} \left[ \frac{\eta_t^2}{2} \text{sign}(\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Sign}}))^\top (\nabla^2 L(\boldsymbol{\theta}_t^{\text{Sign}}) - \nabla^2 L_D(\boldsymbol{\theta}_t^{\text{Sign}})) \text{sign}(\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Sign}})) \mid \boldsymbol{\theta}_t^{\text{Sign}} \right] + \eta_t^3 \frac{\rho_H}{6} P^{3/2} \\ & = \mathbb{E} \left[ \frac{\eta_t^2}{2} \sum_b [\text{sign}(\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Sign}}))]_b^\top [\nabla^2 L(\boldsymbol{\theta}_t^{\text{Sign}})]_b [\text{sign}(\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Sign}}))]_b \mid \boldsymbol{\theta}_t^{\text{Sign}} \right] \\ & \quad + \mathbb{E} \left[ \frac{\eta_t^2}{2} \text{sign}(\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Sign}}))^\top (\nabla^2 L(\boldsymbol{\theta}_t^{\text{Sign}}) - \nabla^2 L_D(\boldsymbol{\theta}_t^{\text{Sign}})) \text{sign}(\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Sign}})) \mid \boldsymbol{\theta}_t^{\text{Sign}} \right] + \eta_t^3 \frac{\rho_H}{6} P^{3/2} \\ & \leq \frac{\eta_t^2}{2} \sum_b \|\nabla^2 L(\boldsymbol{\theta}_t^{\text{Sign}})]_b\|_2 P_b + \frac{\eta_t^2}{2} \|\nabla^2 L(\boldsymbol{\theta}_t^{\text{Sign}}) - \nabla^2 L_D(\boldsymbol{\theta}_t^{\text{Sign}})\|_2 P + \eta_t^3 \frac{\rho_H}{6} P^{3/2} \\ & \leq \frac{\eta_t^2}{2} \Lambda_P P + \frac{\eta_t^2}{2} \delta_D P + \eta_t^3 \frac{\rho_H}{6} P^{3/2}. \end{aligned} \quad (14)$$

For the third term, we can derive an upper bound as follows:

$$\begin{aligned} & \mathbb{E} \left[ -\eta_t \nabla L(\boldsymbol{\theta}_t^{\text{Sign}})^\top (\text{sign}(\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Sign}})) - \text{sign}(\nabla L(\boldsymbol{\theta}_t^{\text{Sign}}))) \mid \boldsymbol{\theta}_t^{\text{Sign}} \right] \\ & = \eta_t \sum_{i=1}^P \nabla L(\boldsymbol{\theta}_t^{\text{Sign}})_i \mathbb{E} \left[ \text{sign}(\nabla L(\boldsymbol{\theta}_t^{\text{Sign}}))_i - \text{sign}(\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Sign}}))_i \mid \boldsymbol{\theta}_t^{\text{Sign}} \right] \\ & = \eta_t \sum_{i=1}^P |\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})_i| 2 \mathbb{E} \left[ \mathbb{1}[\text{sign}(\nabla L(\boldsymbol{\theta}_t^{\text{Sign}}))_i \neq \text{sign}(\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Sign}}))_i] \mid \boldsymbol{\theta}_t^{\text{Sign}} \right] \\ & = \eta_t \sum_{i=1}^P |\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})_i| 2 \mathbb{P} \left( \text{sign}(\nabla L(\boldsymbol{\theta}_t^{\text{Sign}}))_i \neq \text{sign}(\nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Sign}}))_i \mid \boldsymbol{\theta}_t^{\text{Sign}} \right) \end{aligned}$$

$$\begin{aligned}
&\leq \eta_t \sum_{i=1}^P |\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})_i| 2\mathbb{P} \left( |\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})_i - \nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Sign}})_i| \geq |\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})_i| \mid \boldsymbol{\theta}_t^{\text{Sign}} \right) \\
&\leq \eta_t \sum_{i=1}^P |\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})_i| 2 \frac{\mathbb{E} \left[ |\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})_i - \nabla \widehat{L}(\boldsymbol{\theta}_t^{\text{Sign}})_i|^2 \mid \boldsymbol{\theta}_t^{\text{Sign}} \right]}{|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})_i|^2} \\
&\leq \eta_t \sum_{i=1}^P |\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})_i| 2\sigma_2 \\
&= 2\sigma_2 \eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1,
\end{aligned} \tag{15}$$

where the second inequality follows from Chebyshev's inequality, and the last inequality uses Eq. (4).

Combining Eqs.(13)(14)(15),  $\eta_t \leq \min(\frac{\|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1}{\Lambda_P P}, \sqrt{\frac{\|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1}{\rho_H P^{3/2}}})$ ,  $\delta_D < \Lambda_P/3$ , and  $\sigma_2 \leq \frac{1}{24}$ , we have:

$$\begin{aligned}
&\mathbb{E} \left[ L(\boldsymbol{\theta}_{t+1}^{\text{Sign}}) - L(\boldsymbol{\theta}_t^{\text{Sign}}) \mid \boldsymbol{\theta}_t^{\text{Sign}} \right] \\
&\leq -\eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 + \frac{\eta_t^2}{2} \Lambda_P P + \frac{\eta_t^2}{2} \delta_D P + \eta_t^3 \frac{\rho_H}{6} P^{3/2} + 2\sigma_2 \eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 \\
&\leq -\eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 + \frac{\eta_t}{2} \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 + \frac{\eta_t}{6} \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 + \frac{\eta_t}{6} \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 + 2\sigma_2 \eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 \\
&= -\frac{(1 - 12\sigma_2)\eta_t}{6} \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 \\
&\leq -\frac{\eta_t}{6(1 + 24\sigma_2)} \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1.
\end{aligned} \tag{16}$$

Assume that the probability of the event  $\mathcal{E}(T) = \{\forall s \leq T, \|\nabla L(\boldsymbol{\theta}_s^{\text{Sign}})\|_1 \geq P\varepsilon\}$  satisfies  $\mathbb{P}(\mathcal{E}(T)) \geq \frac{1}{2}$ . By applying the telescoping sum and taking expectations, and noting that  $\boldsymbol{\theta}_0 = \boldsymbol{\theta}_0^{\text{Sign}}$ ,  $\eta_t \geq \zeta_0 \min(\frac{\|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1}{\Lambda_P P}, \sqrt{\frac{\|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1}{\rho_H P^{3/2}}})$ , and  $\varepsilon < \frac{\Lambda_P^2}{\rho_H \sqrt{P}}$  we have:

$$\begin{aligned}
&\mathbb{E} \left[ L(\boldsymbol{\theta}_T^{\text{Sign}}) \right] - L(\boldsymbol{\theta}_0) \\
&\leq -\frac{1}{6(1 + 24\sigma_2)} \sum_{t=0}^{T-1} \mathbb{E} \left[ \eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 \right] \\
&= -\frac{1}{6(1 + 24\sigma_2)} \sum_{t=0}^{T-1} \left( \mathbb{E} \left[ \eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 \mid \mathcal{E}(T) \right] \mathbb{P}(\mathcal{E}(T)) + \mathbb{E} \left[ \bar{\eta}_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 \mid \overline{\mathcal{E}(T)} \right] \mathbb{P}(\overline{\mathcal{E}(T)}) \right) \\
&\leq -\frac{1}{6(1 + 24\sigma_2)} \sum_{t=0}^{T-1} \mathbb{E} \left[ \eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 \mid \mathcal{E}(T) \right] \mathbb{P}(\mathcal{E}(T)) \\
&\leq -\frac{1}{12(1 + 24\sigma_2)} \sum_{t=0}^{T-1} \mathbb{E} \left[ \eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 \mid \mathcal{E}(T) \right] \\
&\leq -\frac{\zeta_0}{12(1 + 24\sigma_2)} \sum_{t=0}^{T-1} \mathbb{E} \left[ \min\left( \frac{\|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1^2}{\Lambda_P P}, \frac{\|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1^{3/2}}{\sqrt{\rho_H P^{3/2}}} \right) \mid \mathcal{E}(T) \right] \\
&\leq -\frac{\zeta_0}{12(1 + 24\sigma_2)} \sum_{t=0}^{T-1} \min\left( \frac{P\varepsilon^2}{\Lambda_P}, P\varepsilon \sqrt{\frac{\varepsilon}{\rho_H P^{1/2}}} \right) \\
&= -\frac{TP\varepsilon^2 \zeta_0}{12(1 + 24\sigma_2)\Lambda_P}.
\end{aligned}$$

Therefore, we have:

$$\begin{aligned} T &\leq \frac{12(1 + 24\sigma_2)(L(\boldsymbol{\theta}_0) - \mathbb{E}[L(\boldsymbol{\theta}_T^{\text{Sign}})])}{P\varepsilon^2\zeta_0} \Lambda_P \\ &\leq \frac{12(1 + 24\sigma_2)(L(\boldsymbol{\theta}_0) - L_*)}{P\varepsilon^2\zeta_0} \Lambda_P. \end{aligned}$$

This means that when we take  $T > \frac{12(1+24\sigma_2)(L(\boldsymbol{\theta}_0)-L_*)}{P\varepsilon^2\zeta_0} \Lambda_P$ , we have  $\mathbb{P}(\mathcal{E}(T)) < \frac{1}{2}$ . Therefore, we have

$$\mathcal{T}_\varepsilon(\{\boldsymbol{\theta}_t^{\text{Sign}}\}_{t=0}^\infty, L, \|\cdot\|_1) \leq \frac{12(1 + 24\sigma_2)(L(\boldsymbol{\theta}_0) - L_*)}{P\varepsilon^2\zeta_0} \Lambda_P.$$

□

## B Derivation of Jacobian matrix in Section 4.7

### B.1 Jacobian of Transformer layer

The output of a Transformer layer for an input  $\mathbf{X} \in \mathbb{R}^{n \times d}$  is given by  $\mathcal{M}(\mathcal{A}(\mathbf{X}))$ , where  $\mathcal{A}(\cdot)$  is the attention layer and  $\mathcal{M}(\cdot)$  is the feed-forward layer. In the following, we denote the Jacobian of the self-attention module, the feed-forward module, and the layer normalization as  $\mathbf{J}_{\text{ATT}}$ ,  $\mathbf{J}_{\text{FFN}}$ , and  $\mathbf{J}_{\text{LN}}$ , respectively.

**In Pre-LN.** The self-attention and feed-forward layers in the Pre-LN architecture are given by

$$\begin{aligned} \mathcal{A}(\mathbf{X}) &= \text{ATT}(\text{LN}(\mathbf{X})) + \mathbf{X}, \\ \mathcal{M}(\mathbf{Y}) &= \text{FFN}(\text{LN}(\mathbf{Y})) + \mathbf{Y}. \end{aligned}$$

The Jacobian of these modules are as follows:

$$\begin{aligned} \frac{\partial \mathcal{A}(\mathbf{X})}{\partial \mathbf{X}} &= \left. \frac{\partial \text{ATT}(\mathbf{Z})}{\partial \mathbf{Z}} \right|_{\mathbf{Z}=\text{LN}(\mathbf{X})} \frac{\partial \text{LN}(\mathbf{X})}{\partial \mathbf{X}} + \frac{\partial \mathbf{X}}{\partial \mathbf{X}} \\ &= \mathbf{J}_{\text{ATT}}(\text{LN}(\mathbf{X})) \mathbf{J}_{\text{LN}}(\mathbf{X}) + \mathbf{I}_{nd}, \\ \frac{\partial \mathcal{M}(\mathbf{Y})}{\partial \mathbf{Y}} &= \left. \frac{\partial \text{FFN}(\mathbf{Y})}{\partial \mathbf{Y}} \right|_{\mathbf{Y}=\text{LN}(\mathbf{Y})} \frac{\partial \text{LN}(\mathbf{Y})}{\partial \mathbf{Y}} + \frac{\partial \mathbf{Y}}{\partial \mathbf{Y}} \\ &= \mathbf{J}_{\text{FFN}}(\text{LN}(\mathbf{Y})) \mathbf{J}_{\text{LN}}(\mathbf{Y}) + \mathbf{I}_{nd}. \end{aligned}$$

Therefore, the Jacobian of the Pre-LN layer is given by

$$\begin{aligned} \mathbf{J}_{\text{Pre-LN}}(\mathbf{X}) &= \left. \frac{\partial \mathcal{M}(\mathbf{Y})}{\partial \mathbf{Y}} \right|_{\mathbf{Y}=\mathcal{A}(\mathbf{X})} \frac{\partial \mathcal{A}(\mathbf{X})}{\partial \mathbf{X}} \\ &= (\mathbf{J}_{\text{FFN}}(\text{LN}(\mathcal{A}(\mathbf{X}))) \mathbf{J}_{\text{LN}}(\mathcal{A}(\mathbf{X})) + \mathbf{I}_{nd}) (\mathbf{J}_{\text{ATT}}(\text{LN}(\mathbf{X})) \mathbf{J}_{\text{LN}}(\mathbf{X}) + \mathbf{I}_{nd}) \end{aligned}$$

and with omitting the evaluation point, we can write the Jacobian as

$$\mathbf{J}_{\text{Pre-LN}} = (\mathbf{J}_{\text{FFN}} \mathbf{J}_{\text{LN}} + \mathbf{I}_{nd}) (\mathbf{J}_{\text{ATT}} \mathbf{J}_{\text{LN}} + \mathbf{I}_{nd}).$$

**In Post-LN.** The self-attention and feed-forward layers in the Post-LN layer are given by

$$\begin{aligned} \mathcal{A}(\mathbf{X}) &= \text{LN}(\text{ATT}(\mathbf{X}) + \mathbf{X}), \\ \mathcal{M}(\mathbf{Y}) &= \text{LN}(\text{FFN}(\mathbf{Y}) + \mathbf{Y}). \end{aligned}$$

The Jacobian of these modules are as follows:

$$\frac{\partial \mathcal{A}(\mathbf{X})}{\partial \mathbf{X}} = \left. \frac{\partial \text{LN}(\mathbf{Z})}{\partial \mathbf{Z}} \right|_{\mathbf{Z}=\text{ATT}(\mathbf{X})+\mathbf{X}} \left( \frac{\partial \text{ATT}(\mathbf{X})}{\partial \mathbf{X}} + \frac{\partial \mathbf{X}}{\partial \mathbf{X}} \right)$$

$$\begin{aligned}
&= \mathbf{J}_{\text{LN}}(\text{ATT}(\mathbf{X}) + \mathbf{X}) (\mathbf{J}_{\text{ATT}}(\mathbf{X}) + \mathbf{I}_{nd}), \\
\frac{\partial \mathcal{M}(\mathbf{Y})}{\partial \mathbf{Y}} &= \frac{\partial \text{LN}(\mathbf{Z})}{\partial \mathbf{Z}} \Big|_{\mathbf{Z}=\text{FFN}(\mathbf{Y})+\mathbf{Y}} \left( \frac{\partial \text{FFN}(\mathbf{Y})}{\partial \mathbf{Y}} + \frac{\partial \mathbf{Y}}{\partial \mathbf{Y}} \right) \\
&= \mathbf{J}_{\text{LN}}(\text{FFN}(\mathbf{Y}) + \mathbf{Y}) (\mathbf{J}_{\text{FFN}}(\mathbf{Y}) + \mathbf{I}_{nd}).
\end{aligned}$$

Therefore, the Jacobian of the Post-LN layer is given by

$$\begin{aligned}
\mathbf{J}_{\text{Post-LN}}(\mathbf{X}) &= \frac{\partial \mathcal{M}(\mathbf{Y})}{\partial \mathbf{Y}} \Big|_{\mathbf{Y}=\mathcal{A}(\mathbf{X})} \frac{\partial \mathcal{A}(\mathbf{X})}{\partial \mathbf{X}} \\
&= \mathbf{J}_{\text{LN}}(\text{FFN}(\mathcal{A}(\mathbf{X})) + \mathcal{A}(\mathbf{X})) (\mathbf{J}_{\text{FFN}}(\mathcal{A}(\mathbf{X})) + \mathbf{I}_{nd}) \mathbf{J}_{\text{LN}}(\text{ATT}(\mathbf{X}) + \mathbf{X}) (\mathbf{J}_{\text{ATT}}(\mathbf{X}) + \mathbf{I}_{nd})
\end{aligned}$$

and with omitting the evaluation point, we can write the Jacobian as

$$\mathbf{J}_{\text{Post-LN}} = \mathbf{J}_{\text{LN}} (\mathbf{J}_{\text{FFN}} + \mathbf{I}_{nd}) \mathbf{J}_{\text{LN}} (\mathbf{J}_{\text{ATT}} + \mathbf{I}_{nd}).$$

## B.2 Jacobian of layer normalization

Since the layer normalization is a row-wise operation, the Jacobian of the layer normalization for the input matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$  is given by

$$\mathbf{J}_{\text{LN}}(\mathbf{X}) = \text{blockdiag}(\left\{ \frac{\partial \text{LN}(\mathbf{X})_{i,:}}{\partial \mathbf{X}_{i,:}} \right\}_{i=1}^n).$$

where  $\frac{\partial \text{LN}(\mathbf{X})_{i,:}}{\partial \mathbf{X}_{i,:}}$  is the Jacobian of the layer normalization for the  $i$ -th row of the input matrix  $\mathbf{X}$ . The layer normalization for the  $i$ -th row of the input matrix  $\mathbf{X}$  is given by

$$\text{LN}(\mathbf{X})_{i,:} = \frac{\sqrt{d} \widetilde{\mathbf{X}}_{i,:}}{\|\widetilde{\mathbf{X}}_{i,:}\|},$$

where  $\widetilde{\mathbf{X}}_{i,:} := \mathbf{X}_{i,:} (\mathbf{I}_d - \frac{1}{d} \mathbf{1}\mathbf{1}^\top)$ . Therefore, the  $i$ -th block of the Jacobian of the layer normalization is given by

$$\begin{aligned}
\frac{\partial \text{LN}(\mathbf{X})_{i,:}}{\partial \mathbf{X}_{i,:}} &= \frac{\partial \text{LN}(\mathbf{X})_{i,:}}{\partial \widetilde{\mathbf{X}}_{i,:}} \frac{\partial \widetilde{\mathbf{X}}_{i,:}}{\partial \mathbf{X}_{i,:}} \\
&= \sqrt{d} \left( \frac{1}{\|\widetilde{\mathbf{X}}_{i,:}\|} \mathbf{I}_d - \widetilde{\mathbf{X}}_{i,:} \frac{\widetilde{\mathbf{X}}_{i,:}^\top}{\|\widetilde{\mathbf{X}}_{i,:}\|^3} \right) (\mathbf{I}_d - \frac{1}{d} \mathbf{1}\mathbf{1}^\top) \\
&= \frac{\sqrt{d}}{\|\widetilde{\mathbf{X}}_{i,:}\|_2} \left( \mathbf{I}_d - \frac{\widetilde{\mathbf{X}}_{i,:} \widetilde{\mathbf{X}}_{i,:}^\top}{\|\widetilde{\mathbf{X}}_{i,:}\|_2^2} \right) \left( \mathbf{I}_d - \frac{\mathbf{1}\mathbf{1}^\top}{d} \right).
\end{aligned}$$

Therefore, we can write the Jacobian of the layer normalization as

$$\mathbf{J}_{\text{LN}}(\mathbf{X}) = \text{blockdiag}(\{\mathbf{L}_i(\mathbf{X})\}_{i=1}^n),$$

where

$$\mathbf{L}_i(\mathbf{X}) = \frac{\sqrt{d}}{\|\widetilde{\mathbf{X}}_{i,:}\|_2} \left( \mathbf{I}_d - \frac{\widetilde{\mathbf{X}}_{i,:} \widetilde{\mathbf{X}}_{i,:}^\top}{\|\widetilde{\mathbf{X}}_{i,:}\|_2^2} \right) \left( \mathbf{I}_d - \frac{\mathbf{1}\mathbf{1}^\top}{d} \right).$$

## B.3 Jacobian of RMS normalization

Since the RMS normalization is a row-wise operation, the Jacobian of the RMS normalization for the input matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$  is given by

$$\mathbf{J}_{\text{RMS}}(\mathbf{X}) = \text{blockdiag} \left( \left\{ \frac{\partial \text{RMS}(\mathbf{X})_{i,:}}{\partial \mathbf{X}_{i,:}} \right\}_{i=1}^n \right).$$

where  $\frac{\partial \text{RMS}(\mathbf{X})_{i,:}}{\partial \mathbf{X}_{i,:}}$  is the Jacobian of the RMS normalization for the  $i$ -th row of the input matrix  $\mathbf{X}$ . The RMS normalization for the  $i$ -th row of the input matrix  $\mathbf{X}$  is given by

$$\text{RMS}(\mathbf{X})_{i,:} = \text{Diag}(\boldsymbol{\gamma}) \frac{\mathbf{X}_{i,:}}{r_i}, \quad r_i := \sqrt{\frac{1}{d} \|\mathbf{X}_{i,:}\|_2^2},$$

where  $\boldsymbol{\gamma} \in \mathbb{R}^d$  is a learnable scale parameter. Let  $\mathbf{x} := \mathbf{X}_{i,:} \in \mathbb{R}^{d \times 1}$  and  $\mathbf{y} := \text{RMS}(\mathbf{X})_{i,:} \in \mathbb{R}^{d \times 1}$ . Then  $\mathbf{y} = \text{Diag}(\boldsymbol{\gamma})\mathbf{x}/r$  with  $r = \sqrt{\frac{1}{d} \|\mathbf{x}\|_2^2}$ .

First, we compute the derivative of  $r$  with respect to  $\mathbf{x}$ . Since  $r = (\frac{1}{d} \mathbf{x}^\top \mathbf{x})^{1/2}$ , we have

$$\frac{\partial r}{\partial \mathbf{x}} = \frac{1}{2} \left( \frac{1}{d} \mathbf{x}^\top \mathbf{x} \right)^{-1/2} \cdot \frac{2}{d} \mathbf{x} = \frac{1}{dr} \mathbf{x}.$$

Hence,

$$\frac{\partial(r^{-1})}{\partial \mathbf{x}} = -r^{-2} \frac{\partial r}{\partial \mathbf{x}} = -\frac{1}{dr^3} \mathbf{x}.$$

Therefore, the  $i$ -th block of the Jacobian of the RMS normalization is given by

$$\begin{aligned} \frac{\partial \text{RMS}(\mathbf{X})_{i,:}}{\partial \mathbf{X}_{i,:}} &= \text{Diag}(\boldsymbol{\gamma}) \frac{\partial(\mathbf{x}r^{-1})}{\partial \mathbf{x}} \\ &= \text{Diag}(\boldsymbol{\gamma}) \left( \frac{1}{r} \mathbf{I}_d + \mathbf{x} \frac{\partial(r^{-1})}{\partial \mathbf{x}^\top} \right) \\ &= \text{Diag}(\boldsymbol{\gamma}) \left( \frac{1}{r} \mathbf{I}_d - \frac{1}{dr^3} \mathbf{x} \mathbf{x}^\top \right) \\ &= \text{Diag}(\boldsymbol{\gamma}) \left( \frac{1}{r_i} \mathbf{I}_d - \frac{1}{dr_i^3} \mathbf{X}_{i,:} \mathbf{X}_{i,:}^\top \right) \\ &= \text{Diag}(\boldsymbol{\gamma}) \frac{\sqrt{d}}{\|\mathbf{X}_{i,:}\|_2} \left( \mathbf{I}_d - \frac{\mathbf{X}_{i,:} \mathbf{X}_{i,:}^\top}{\|\mathbf{X}_{i,:}\|_2^2} \right). \end{aligned}$$

Therefore, we can write the Jacobian of the RMS normalization as

$$\mathbf{J}_{\text{RMS}}(\mathbf{X}) = \text{blockdiag}(\{\mathbf{R}_i(\mathbf{X})\}_{i=1}^n),$$

where

$$\mathbf{R}_i(\mathbf{X}) = \text{Diag}(\boldsymbol{\gamma}) \frac{\sqrt{d}}{\|\mathbf{X}_{i,:}\|_2} \left( \mathbf{I}_d - \frac{\mathbf{X}_{i,:} \mathbf{X}_{i,:}^\top}{\|\mathbf{X}_{i,:}\|_2^2} \right), \quad r_i = \sqrt{\frac{1}{d} \|\mathbf{X}_{i,:}\|_2^2}.$$

We note that all the above derivations for Pre-LN and Post-LN architectures remain valid when layer normalization is replaced with RMS normalization, by simply substituting  $\mathbf{J}_{\text{LN}}$  with  $\mathbf{J}_{\text{RMS}}$ .

## C Experimental details

### C.1 Implementation and training details

Our implementation, based on PyTorch (Paszke et al., 2019), uses the HuggingFace Transformers library (Wolf et al., 2020) for NLP tasks and primarily follows Tomihari & Sato (2024). All experiments were conducted on a single NVIDIA A100 GPU. The reported results are averages over five training seeds. We used the cross-entropy loss, defined as  $\ell(\mathbf{f}(\mathbf{x}), y) := -\log(\sigma_{\text{SM}}(\mathbf{f}(\mathbf{x}))_y)$ , where the function  $\sigma_{\text{SM}} : \mathbb{R}^C \rightarrow \mathbb{R}^C$  represents the softmax operation.

Following the methodology of Kunstner et al. (2023), we optimized the learning rate via grid search based on the training loss, while keeping other hyperparameters, such as batch size and the number of epochs, fixed. Momentum was set to 0.9 as the default configuration for both SGD and SignSGD. Gradient clipping with a threshold of 1.0 was applied to SGD, which is equivalent to normalized gradient descent up to a constant factor in the learning rate (Zhang et al., 2020a). For NLP tasks, we used linear learning rate scheduling, whereas for vision tasks, a warmup schedule was applied.

Other hyperparameters followed the default values provided by PyTorch, including Adam ( $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ,  $\epsilon = 1e - 8$ ) and RMSProp ( $\alpha = 0.99$ ,  $\epsilon = 1e - 8$ ). For NLP tasks, the original training set was split into a 9:1 training-to-validation ratio, with the original validation set used as the test set, following Chen et al. (2022); Tomihari & Sato (2024).

We provide dataset statistics and hyperparameter configurations in Table S.1 and Tables S.2–S.4, respectively.

### C.2 Details of each experiment and figure

**Correlation between Hessian and gradient.** In Figure 1, we show the correlation between the Hessian and the gradient. The maximum eigenvalue of the Hessian was computed using power iteration, as described in Park & Kim (2022), with the PyHessian implementation (Yao et al., 2020). To estimate the maximum eigenvectors of the block-diagonal elements of the Hessian, we calculated the product of the Hessian and a random vector for each parameter. The batch size used for these computations was the same as the training batch size. The maximum eigenvalue and the gradient were computed for each batch across all training data.

**Correlation between full-batch gradient and gradient error.** In Figure 2, we show the correlation between the full-batch gradient and the gradient error in a coordinate-wise manner. We randomly sampled 1,000 coordinates from the parameters and computed the squared norm of the full-batch gradient and the gradient error for each coordinate. The gradient error is defined as the difference between the full-batch gradient and the gradient computed with a mini-batch. The batch size was the same as the training batch size. The gradient error was computed for each batch across all training data.

**Gradient heterogeneity.** In Figure 3, we show the ratio of the gradient norm for each parameter relative to the sum of the gradient norms. Specifically, we plot:

$$\frac{G_{\theta}/\sqrt{P_{\theta}}}{\sum_{\theta'} G_{\theta'}/\sqrt{P_{\theta'}}},$$

for each parameter  $\theta$ , where  $G_{\theta}$  is the full-batch gradient norm of parameter  $\theta$ , and  $P_{\theta}$  is its dimension. To compare gradient norms across different parameters, we normalize each gradient norm by the square root of its parameter dimension. Bias parameters are omitted in these plots.

**Effect of layer normalization.** In Tables 2 and S.8, all models share the same RoBERTa backbone and differ only in the placement of the normalization layer. To minimize the effect of initialization, we trained scratch-initialized models for 1000 iterations. Note that pre-trained weights are available only for the Post-LN variant.

**Training Curve.** In Figure 4, we show training runs with the median final loss value among the five training seeds. The shaded area represents the interquartile range across the five seeds. This approach is used to reduce the influence of outliers on the reported results.

Table S.1: Dataset statistics, including the number of classes and counts of training (Train), validation (Val), and test samples for each dataset.

| Domain | Dataset                                 | Classes | Train  | Val   | Test  |
|--------|---|---------|--------|-------|-------|
| NLP    | CB (De Marneffe et al., 2019)           | 3       | 225    | 25    | 57    |
|        | RTE (Wang et al., 2018)                 | 2       | 2,241  | 249   | 277   |
|        | BoolQ (Clark et al., 2019)              | 2       | 8,484  | 943   | 3,270 |
|        | WiC (Pilehvar & Camacho-Collados, 2019) | 2       | 5,400  | 600   | 638   |
|        | CoLA (Warstadt et al., 2019)            | 2       | 7,695  | 855   | 1,040 |
|        | SST-2 (Socher et al., 2013)             | 2       | 60,614 | 6,735 | 872   |
|        | MRPC (Dolan & Brockett, 2005)           | 2       | 3,301  | 367   | 408   |
| Vision | Flowers102 (Nilsback & Zisserman, 2008) | 102     | 1,632  | 408   | 6,149 |
|        | Aircraft (Maji et al., 2013)            | 100     | 5,334  | 1,333 | 3,333 |

Table S.2: Hyperparameter configurations for RoBERTa-Base. The settings include batch size (bs), learning rate (lr), and the number of epochs (epochs). “w/o M” denotes optimizers without momentum and “Const”, “Cos”, and “Lin-W” denote constant, cosine, and linear with warm-up learning rate schedules, respectively.

| Optimizer       | Param  | CB     | RTE    | BoolQ  | WiC    | CoLA   | SST-2  | MRPC   |
|-----------------|--------|--------|--------|--------|--------|--------|--------|--------|
| Common          | bs     | 8      | 8      | 32     | 32     | 32     | 32     | 32     |
|                 | epochs | 20     | 20     | 20     | 20     | 20     | 10     | 20     |
| Adam            |        | $1e-4$ | $1e-5$ | $1e-5$ | $1e-5$ | $1e-5$ | $1e-5$ | $1e-5$ |
| SGD             |        | $1e-2$ | $1e-3$ | $1e-2$ | $1e-3$ | $1e-3$ | $1e-2$ | $1e-2$ |
| SGD (w/o M)     |        | $1e-1$ | $1e-2$ | $1e-1$ | $1e-2$ | $1e-2$ | $1e-1$ | $1e-1$ |
| SignSGD         |        | $1e-5$ | $1e-6$ | $1e-5$ | $1e-5$ | $1e-5$ | $1e-5$ | $1e-5$ |
| SignSGD (w/o M) |        | $1e-4$ | $1e-5$ | $1e-5$ | $1e-5$ | $1e-4$ | $1e-5$ | $1e-5$ |
| RMSProp         | lr     | $1e-5$ | $1e-5$ | $1e-5$ | $1e-5$ | $1e-5$ | $1e-5$ | $1e-5$ |
| SignSGD (S)     |        | $1e-4$ | $1e-4$ | $1e-4$ | $1e-4$ | $5e-4$ | $1e-4$ | $5e-4$ |
| SGD (Const)     |        | $1e-2$ | $1e-3$ | -      | -      | -      | -      | -      |
| SGD (Cos)       |        | $1e-2$ | $1e-3$ | -      | -      | -      | -      | -      |
| SGD (Lin-W)     |        | $1e-2$ | $1e-3$ | -      | -      | -      | -      | -      |
| SignSGD (Const) |        | $1e-6$ | $1e-6$ | -      | -      | -      | -      | -      |
| SignSGD (Cos)   |        | $1e-5$ | $1e-5$ | -      | -      | -      | -      | -      |
| SignSGD (Lin-W) |        | $1e-5$ | $1e-5$ | -      | -      | -      | -      | -      |

Table S.3: Hyperparameter configurations for ResNet18. The settings include batch size (bs), learning rate (lr), and the number of epochs (epochs). “w/o M” denotes optimizers without momentum.

| Optimizer       | Param  | Flowers102 | Aircraft |
|-----------------|--------|------------|----------|
| Common          | bs     | 32         | 32       |
|                 | epochs | 50         | 100      |
| Adam            |        | $1e-4$     | $1e-4$   |
| SGD             |        | $1e-2$     | $1e-2$   |
| SGD (w/o M)     | lr     | $1e-1$     | $1e-1$   |
| SignSGD         |        | $1e-5$     | $1e-5$   |
| SignSGD (w/o M) |        | $1e-4$     | $1e-4$   |
| RMSProp         |        | $1e-4$     | $1e-4$   |
| SignSGD (S)     |        | $5e-4$     | $5e-4$   |

Table S.4: Hyperparameter configurations for ViT-Base. The settings include batch size (bs), learning rate (lr), and the number of epochs (epochs). “w/o M” denotes optimizers without momentum.

| Optimizer       | Param  | Flowers102 | Aircraft |
|-----------------|--------|------------|----------|
| Common          | bs     | 32         | 32       |
|                 | epochs | 50         | 100      |
| Adam            |        | $1e-5$     | $1e-5$   |
| SGD             |        | $1e-2$     | $1e-2$   |
| SGD (w/o M)     | lr     | $1e-1$     | $5e-1$   |
| SignSGD         |        | $1e-5$     | $1e-5$   |
| SignSGD (w/o M) |        | $1e-4$     | $1e-5$   |
| RMSProp         |        | $1e-5$     | $1e-5$   |
| SignSGD (S)     |        | $5e-4$     | $1e-4$   |

## D Additional experimental results

### D.1 Correlation between Hessian and gradient

We show the correlation between the Hessian and the gradient in Figure S.3. The Hessian and gradient are computed using the pre-trained models or the trained models corresponding to the median final loss value among the five training seeds shown in Figures 4 and S.7 and Appendix D.6.

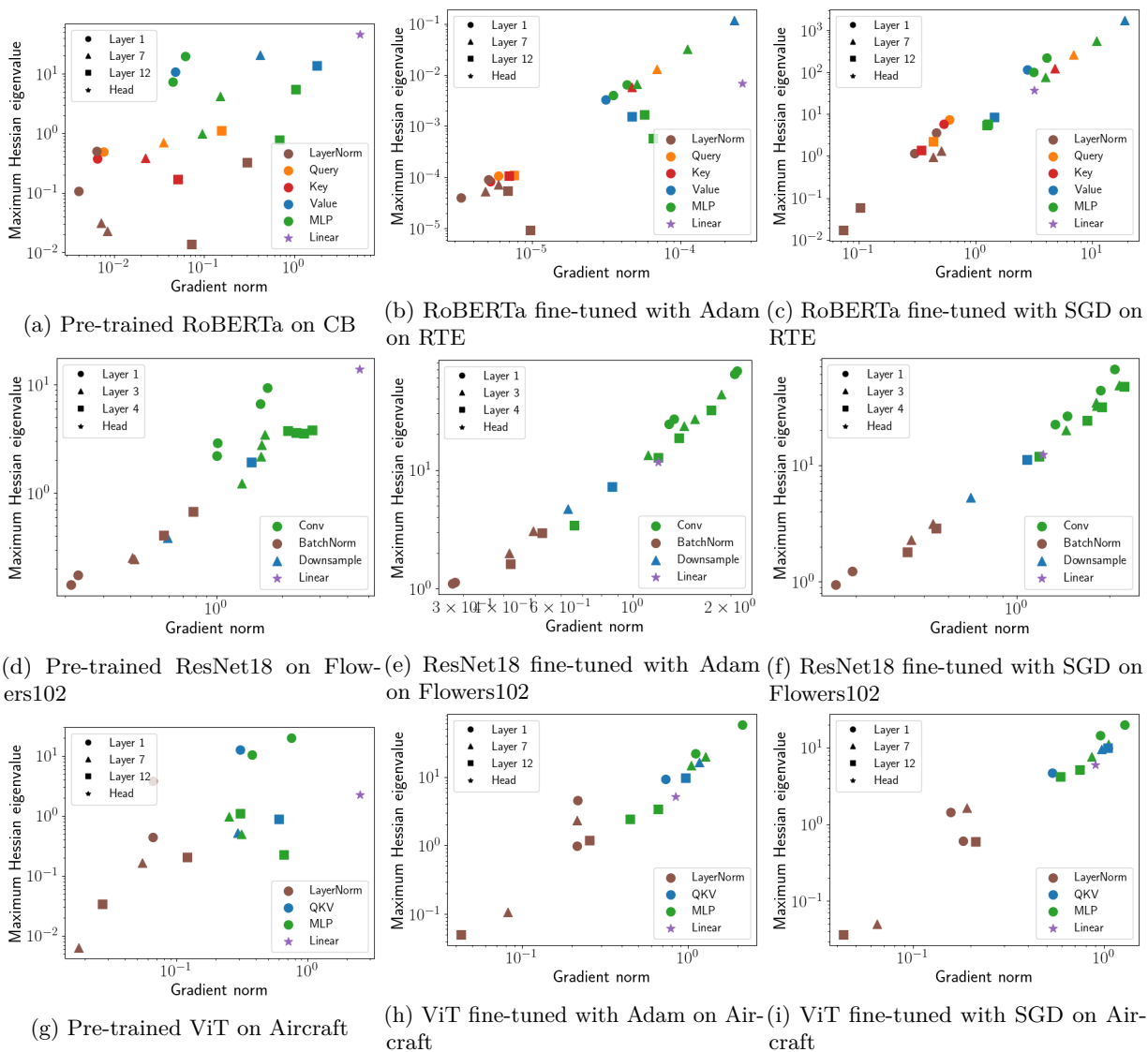


Figure S.3: Gradient vs. Hessian matrix.

## D.2 Correlation between Hessian and parameter dimension

We show the correlation between the Hessian and the parameter in Figure S.4. The Hessian and parameter dimension do not show a clear correlation.

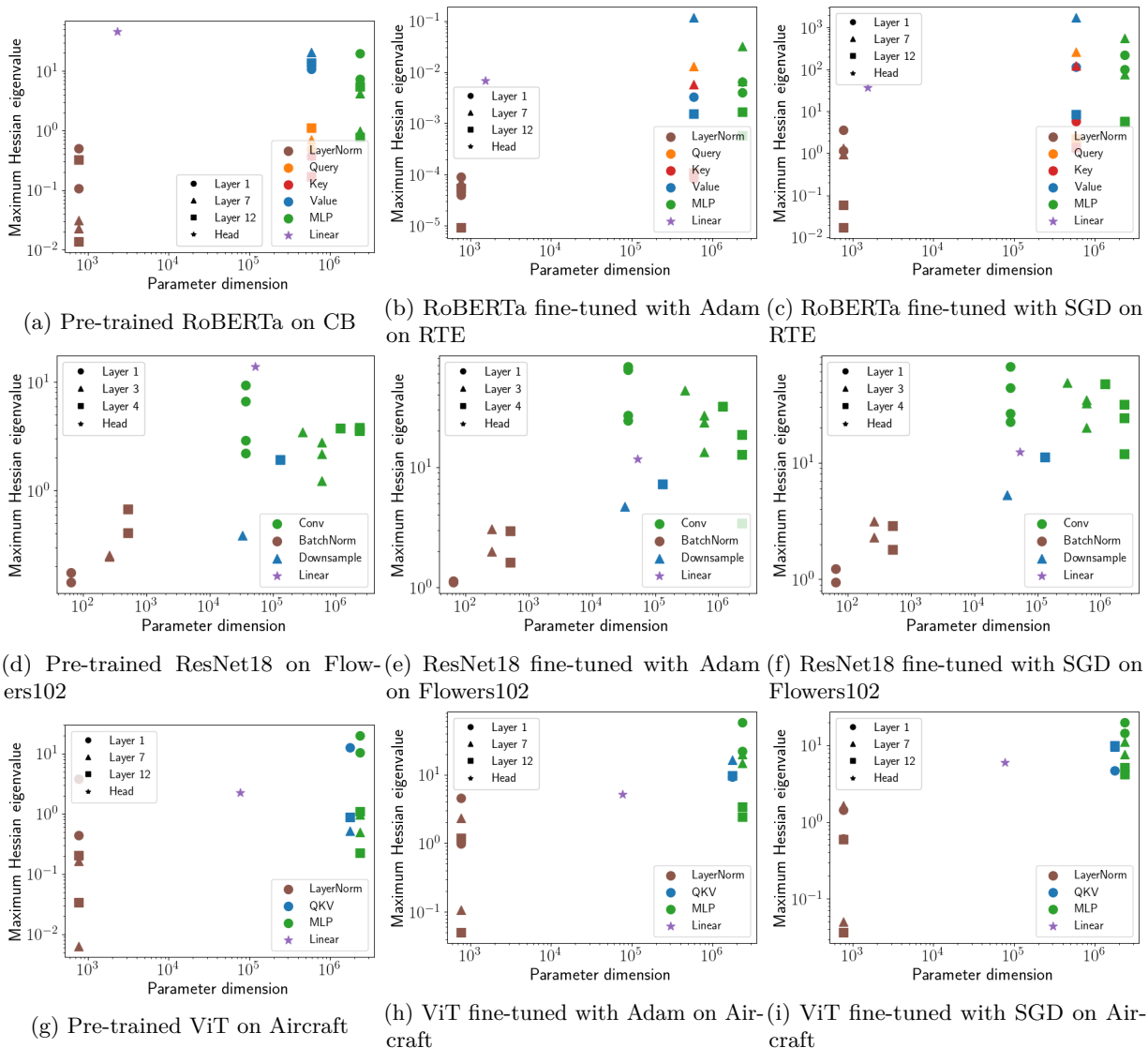


Figure S.4: Parameter dimension vs. Hessian matrix.

### D.3 Correlation between full-batch gradient and gradient error

We show the correlation between the full-batch gradient and the gradient error in Figure S.5.

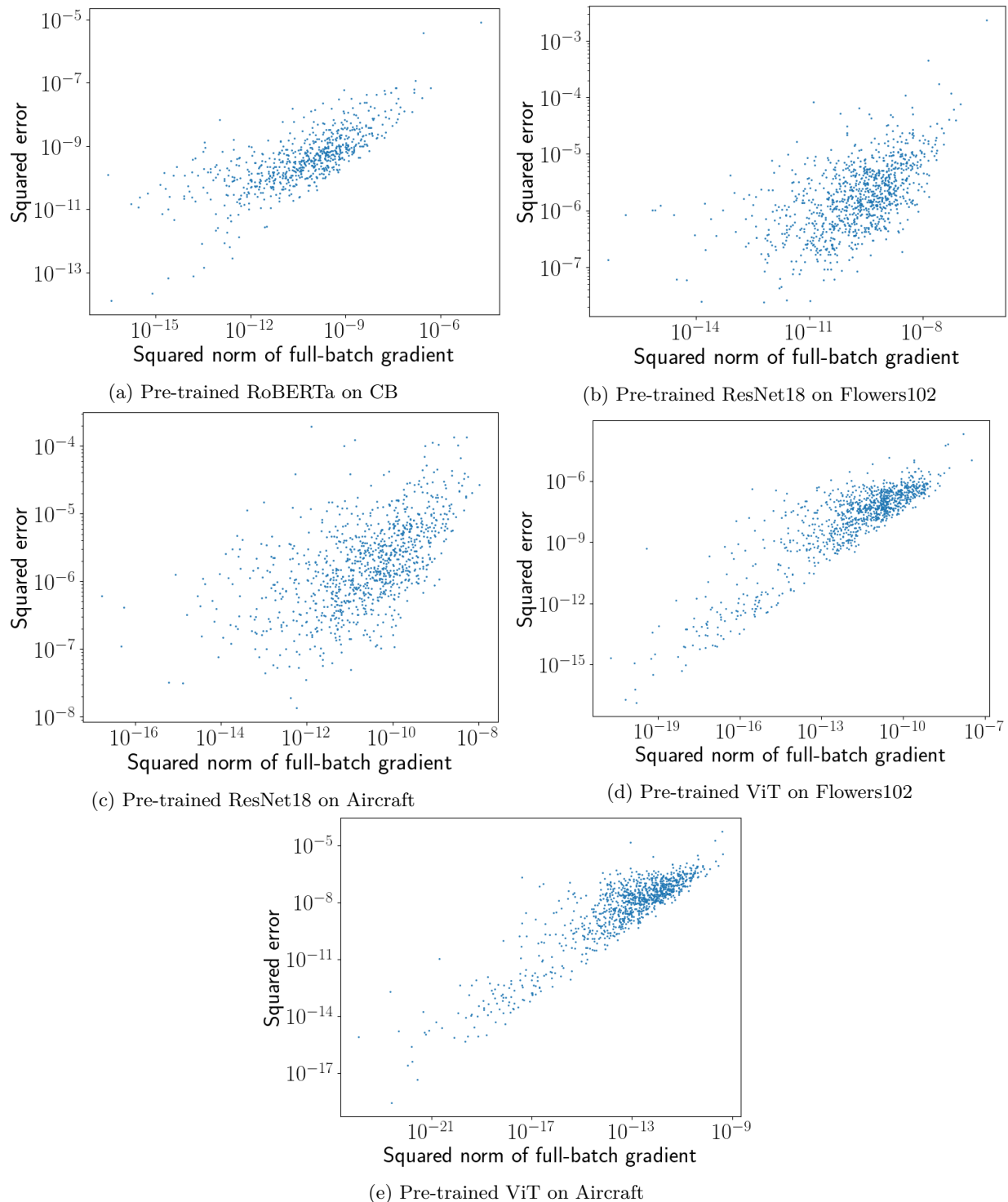


Figure S.5: coordinate-wise full-batch gradient vs. gradient error.

## D.4 Gradient per parameter

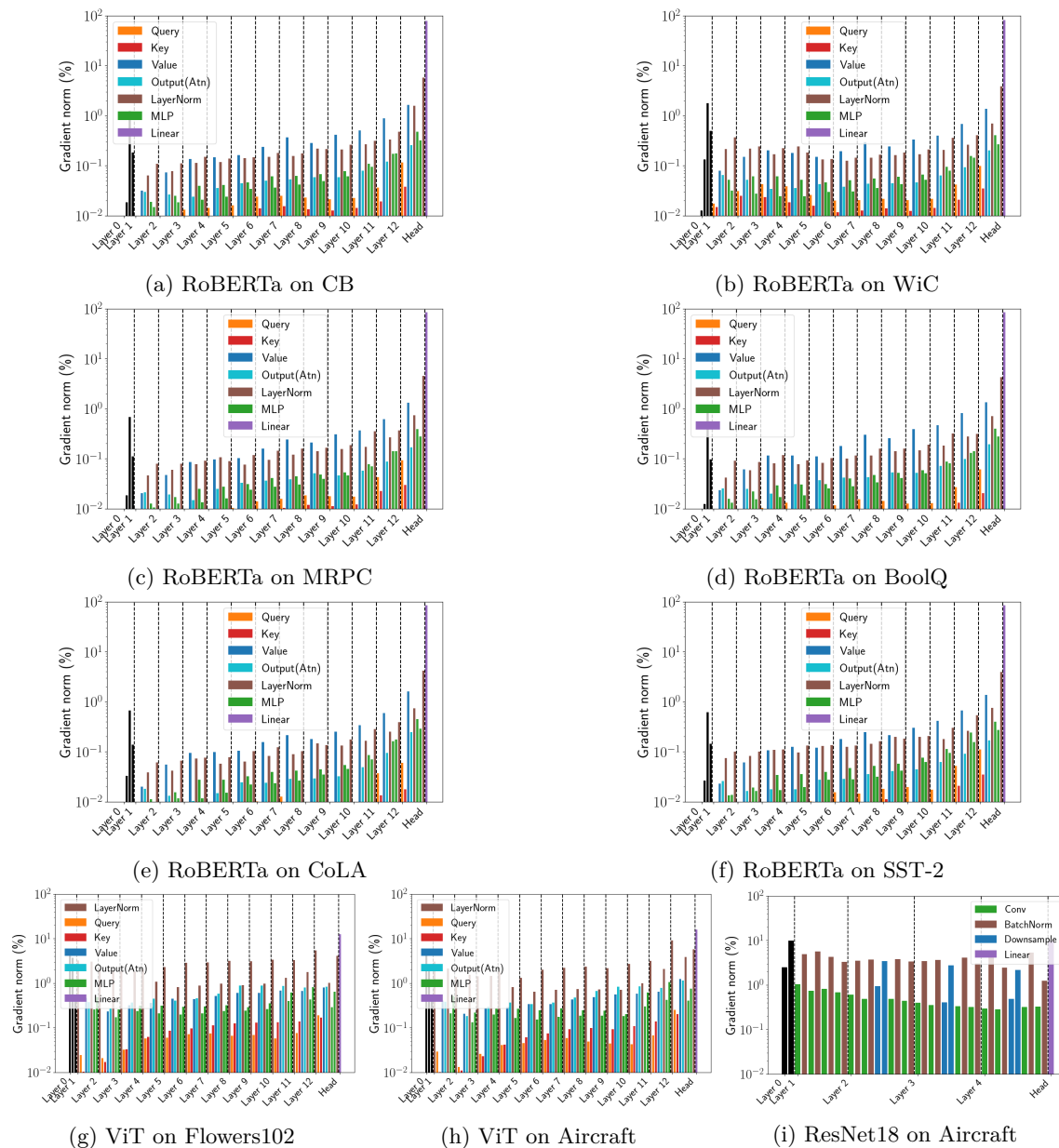


Figure S.6: Gradient norm of each parameter of pre-trained model.

## D.5 Quantitative measures of gradient heterogeneity

**Gini coefficient.** In Table S.5, we provide the Gini coefficient of the normalized gradients.

Gini coefficient is a measure of statistical dispersion intended to represent the inequality of a distribution, which ranges from 0 to 1 and the higher value indicates more heterogeneity.

Given a set of values  $\{x_1, x_2, \dots, x_n\}$  sorted in non-decreasing order, the Gini coefficient is defined as:

$$G = \frac{\sum_{i=1}^n \sum_{j=1}^n |x_i - x_j|}{2n^2 \bar{x}},$$

where  $\bar{x}$  is the mean of the values.

**Layer-wise gradient norm ratio.** In Table S.6, we present the ratio of the gradient norm for each layer, computed as:

$$\frac{G_l}{\sum_{l'} G_{l'}},$$

where  $G_l$  represents the sum of the normalized full-batch gradient norms of the parameters in layer  $l$ . Since all layers contain the same number of parameters, this comparison is valid.

| Model (Dataset)        | Gini coefficient |
|------------------------|------------------|
| RoBERTa-Base (CB)      | 0.932 ± 0.006    |
| RoBERTa-Base (RTE)     | 0.944 ± 0.005    |
| RoBERTa-Base (WiC)     | 0.931 ± 0.004    |
| RoBERTa-Base (BoolQ)   | 0.944 ± 0.001    |
| RoBERTa-Base (CoLA)    | 0.954 ± 0.003    |
| RoBERTa-Base (MRPC)    | 0.951 ± 0.001    |
| RoBERTa-Base (SST-2)   | 0.930 ± 0.032    |
| ResNet-18 (Flowers102) | 0.407 ± 0.013    |
| ResNet-18 (Aircraft)   | 0.433 ± 0.005    |
| ViT-Base (Flowers102)  | 0.539 ± 0.004    |
| ViT-Base (Aircraft)    | 0.598 ± 0.009    |

Table S.5: Gini coefficient of normalized gradients. ± represents standard deviation.

| Layer                 | 1             | 2             | 3             | 4             | 5             | 6             | 7             | 8             | 9             | 10            | 11            | 12            |
|-----------------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
| RoBERTa-Base (CB)     | 0.021 ± 0.001 | 0.022 ± 0.001 | 0.027 ± 0.002 | 0.031 ± 0.002 | 0.036 ± 0.002 | 0.045 ± 0.002 | 0.054 ± 0.002 | 0.060 ± 0.003 | 0.070 ± 0.004 | 0.092 ± 0.005 | 0.156 ± 0.015 | 0.387 ± 0.027 |
| RoBERTa-Base (RTE)    | 0.023 ± 0.003 | 0.024 ± 0.003 | 0.028 ± 0.003 | 0.030 ± 0.003 | 0.034 ± 0.002 | 0.042 ± 0.002 | 0.051 ± 0.004 | 0.058 ± 0.003 | 0.068 ± 0.003 | 0.093 ± 0.008 | 0.163 ± 0.014 | 0.387 ± 0.023 |
| RoBERTa-Base (WiC)    | 0.047 ± 0.014 | 0.042 ± 0.010 | 0.041 ± 0.005 | 0.040 ± 0.003 | 0.036 ± 0.002 | 0.040 ± 0.003 | 0.049 ± 0.004 | 0.055 ± 0.004 | 0.063 ± 0.003 | 0.086 ± 0.006 | 0.145 ± 0.009 | 0.355 ± 0.035 |
| RoBERTa-Base (BoolQ)  | 0.023 ± 0.001 | 0.024 ± 0.001 | 0.028 ± 0.001 | 0.031 ± 0.002 | 0.034 ± 0.002 | 0.043 ± 0.002 | 0.055 ± 0.003 | 0.062 ± 0.004 | 0.073 ± 0.004 | 0.098 ± 0.007 | 0.157 ± 0.010 | 0.370 ± 0.034 |
| RoBERTa-Base (CoLA)   | 0.017 ± 0.001 | 0.018 ± 0.001 | 0.023 ± 0.003 | 0.025 ± 0.002 | 0.029 ± 0.002 | 0.037 ± 0.003 | 0.042 ± 0.002 | 0.048 ± 0.002 | 0.058 ± 0.003 | 0.083 ± 0.006 | 0.169 ± 0.013 | 0.451 ± 0.027 |
| RoBERTa-Base (MRPC)   | 0.019 ± 0.002 | 0.020 ± 0.002 | 0.024 ± 0.002 | 0.028 ± 0.002 | 0.032 ± 0.002 | 0.040 ± 0.002 | 0.049 ± 0.003 | 0.057 ± 0.004 | 0.067 ± 0.004 | 0.089 ± 0.007 | 0.155 ± 0.010 | 0.421 ± 0.037 |
| RoBERTa-Base (SST-2)  | 0.025 ± 0.010 | 0.026 ± 0.010 | 0.032 ± 0.012 | 0.036 ± 0.012 | 0.040 ± 0.013 | 0.046 ± 0.012 | 0.054 ± 0.014 | 0.061 ± 0.014 | 0.070 ± 0.009 | 0.087 ± 0.008 | 0.148 ± 0.022 | 0.373 ± 0.086 |
| ViT-Base (Flowers102) | 0.093 ± 0.004 | 0.065 ± 0.002 | 0.073 ± 0.002 | 0.071 ± 0.004 | 0.069 ± 0.003 | 0.071 ± 0.005 | 0.075 ± 0.005 | 0.079 ± 0.003 | 0.083 ± 0.005 | 0.094 ± 0.002 | 0.105 ± 0.005 | 0.122 ± 0.004 |
| ViT-Base (Aircraft)   | 0.083 ± 0.005 | 0.058 ± 0.003 | 0.067 ± 0.003 | 0.063 ± 0.003 | 0.058 ± 0.002 | 0.063 ± 0.003 | 0.068 ± 0.001 | 0.073 ± 0.002 | 0.077 ± 0.003 | 0.090 ± 0.001 | 0.119 ± 0.005 | 0.181 ± 0.011 |

Table S.6: Layer-wise ratio of gradient norms in Transformers. ± represents standard deviation.

## D.6 Train curves

We show the training curves on different datasets from that in the main text. On the CB dataset, the final train loss is similar among all optimizers, but the convergence speed of SGD is slower than other optimizers. This is consistent with our analysis suggesting the difficulty of training of RoBERTa with SGD.

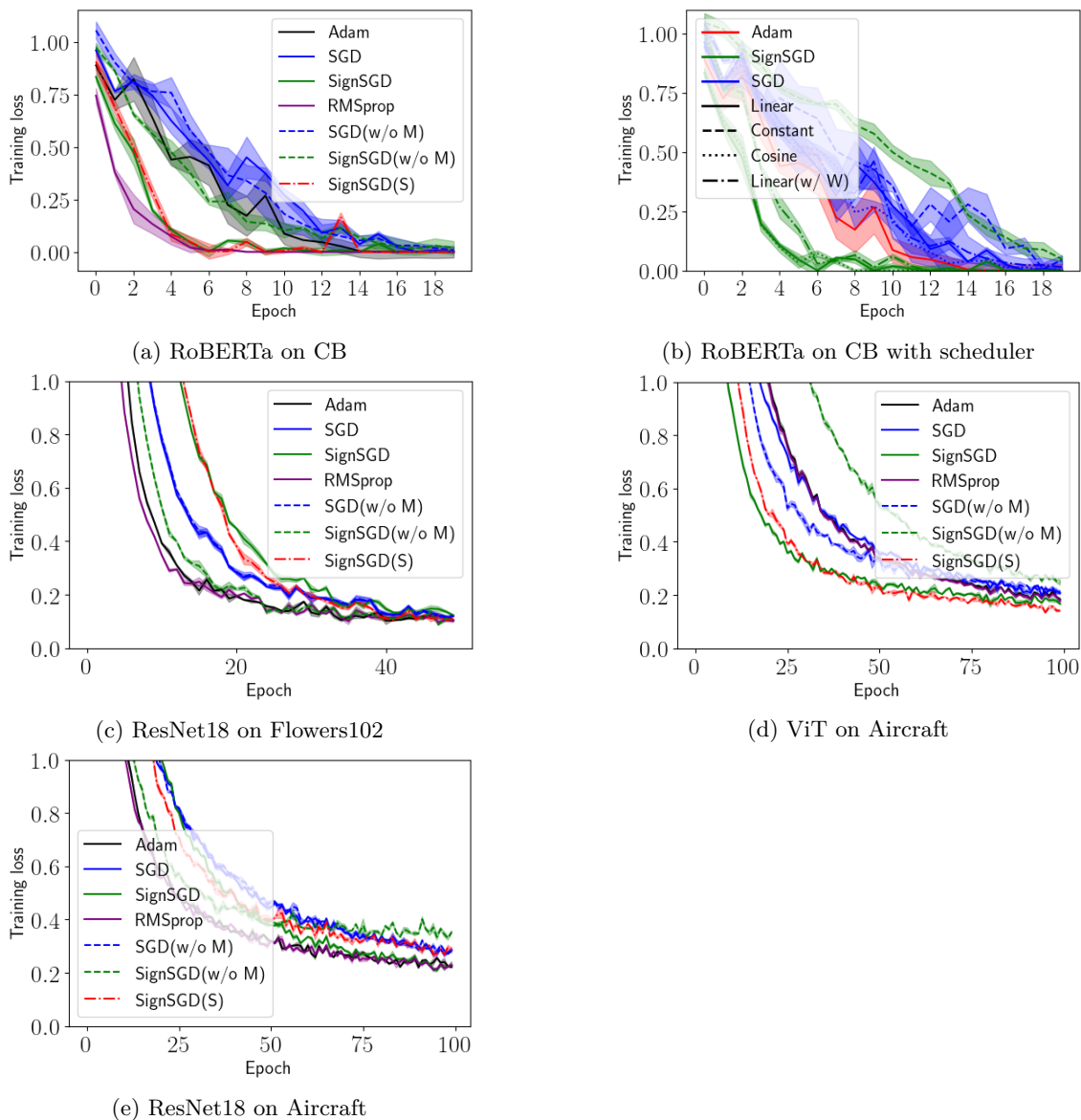


Figure S.7: Training curve with different optimizers. w/ W indicates “with warmup”.

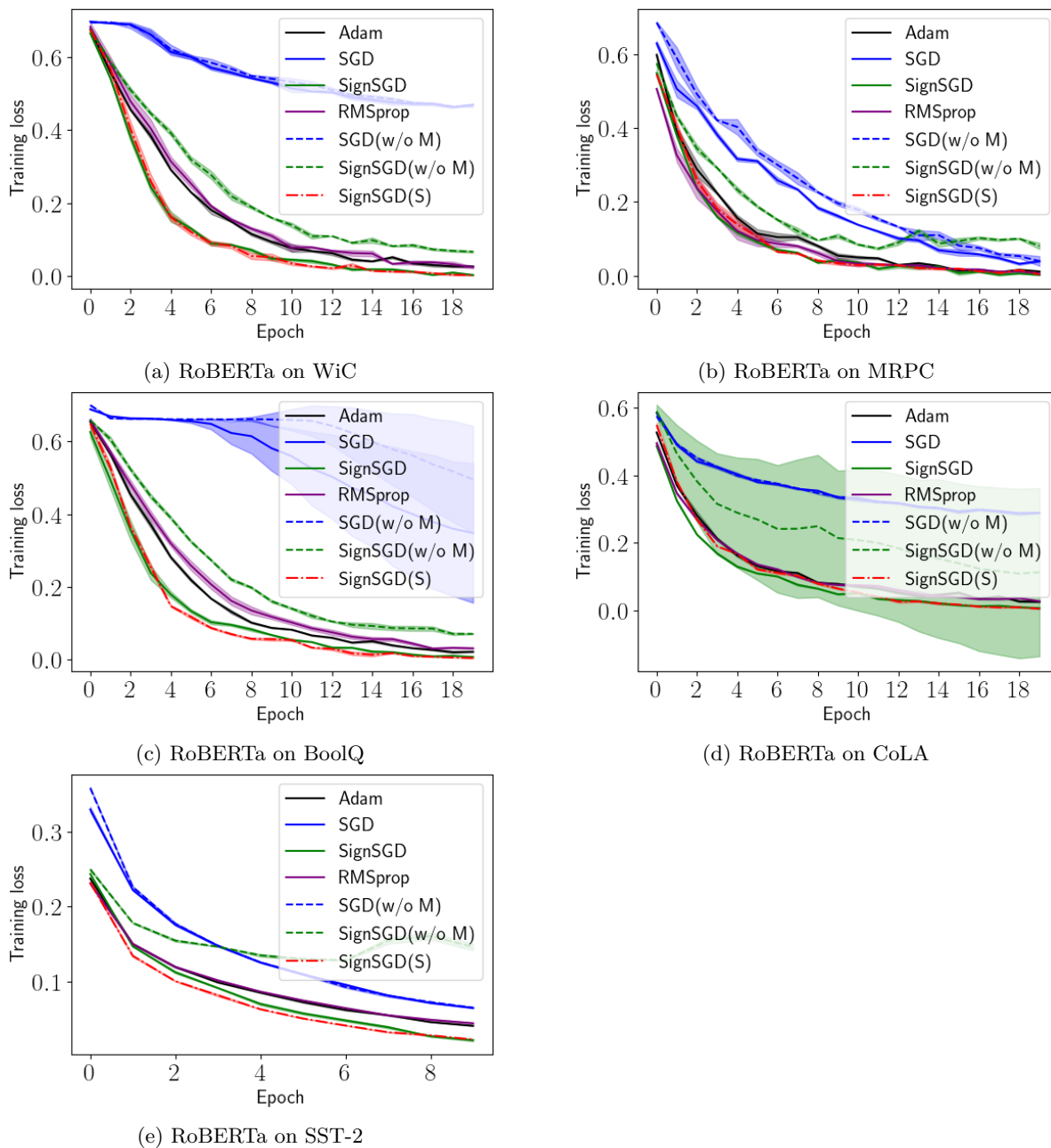


Figure S.8: Training curve with different optimizers.

## D.7 Test results

Table S.7: Test results corresponding to the training curves shown in Figures 4 and S.7. We report the accuracy and its standard deviation.

| Model        | Dataset    | Adam             | RMSprop          | SGD              | SignSGD          | SGD(w/o M)        | SignSGD(w/o M)   |
|--------------|------------|------------------|------------------|------------------|------------------|-------------------|------------------|
| ViT-Base     | Flowers102 | 95.06 $\pm$ 0.34 | 95.15 $\pm$ 0.41 | 94.22 $\pm$ 0.54 | 94.01 $\pm$ 0.98 | 94.49 $\pm$ 0.62  | 92.45 $\pm$ 1.35 |
|              | Aircraft   | 74.28 $\pm$ 0.59 | 74.86 $\pm$ 0.87 | 71.33 $\pm$ 0.27 | 73.96 $\pm$ 0.73 | 55.25 $\pm$ 0.67  | 75.21 $\pm$ 0.88 |
| ResNet18     | Flowers102 | 93.33 $\pm$ 0.62 | 93.27 $\pm$ 0.71 | 93.40 $\pm$ 0.47 | 94.43 $\pm$ 0.54 | 93.03 $\pm$ 0.62  | 93.10 $\pm$ 0.37 |
|              | Aircraft   | 71.95 $\pm$ 0.69 | 70.53 $\pm$ 0.42 | 72.66 $\pm$ 0.71 | 72.01 $\pm$ 0.40 | 72.16 $\pm$ 0.41  | 70.87 $\pm$ 0.35 |
| RoBERTa-Base | CB         | 76.43 $\pm$ 7.41 | 84.29 $\pm$ 4.96 | 78.21 $\pm$ 6.36 | 83.21 $\pm$ 2.71 | 71.79 $\pm$ 12.46 | 77.86 $\pm$ 2.99 |
|              | RTE        | 75.88 $\pm$ 1.56 | 74.66 $\pm$ 2.89 | 75.31 $\pm$ 3.12 | 75.02 $\pm$ 2.30 | 73.21 $\pm$ 1.83  | 75.74 $\pm$ 2.74 |

## D.8 Effect of layer normalization

Table S.8: Gini coefficients of gradient norms for different normalization. A higher Gini coefficient indicates greater heterogeneity. “No-LN” refers to the architecture without layer normalization.

| Norm Type | Init        | Dataset | Gini Coefficient  |
|-----------|-------------|---------|-------------------|
| No-LN     | Scratch     | RTE     | 0.867 $\pm$ 0.006 |
| Pre-LN    | Scratch     | RTE     | 0.880 $\pm$ 0.004 |
| Post-LN   | Scratch     | RTE     | 0.941 $\pm$ 0.012 |
| Post-LN   | Pre-trained | RTE     | 0.944 $\pm$ 0.005 |
| No-LN     | Scratch     | CB      | 0.850 $\pm$ 0.049 |
| Pre-LN    | Scratch     | CB      | 0.873 $\pm$ 0.017 |
| Post-LN   | Scratch     | CB      | 0.899 $\pm$ 0.018 |
| Post-LN   | Pre-trained | CB      | 0.932 $\pm$ 0.006 |

## D.9 Case study: Quadratic model

Following Zhang et al. (2024a), we consider a synthetic quadratic minimization problem of the form

$$L(\boldsymbol{\theta}) = \frac{1}{2} \boldsymbol{\theta}^\top \mathbf{H} \boldsymbol{\theta},$$

where  $\mathbf{H} = \text{blockdiag}(\{\mathbf{H}_i\}_{i=1}^3)$  is a block-diagonal matrix. Each block  $\mathbf{H}_i$  is constructed as  $\mathbf{H}_i = \mathbf{Q}_i \boldsymbol{\Lambda}_i \mathbf{Q}_i^\top$ , where  $\mathbf{Q}_i$  is an orthogonal matrix and  $\boldsymbol{\Lambda}_i$  is a diagonal matrix containing the eigenvalues of the block.

We consider two settings for the eigenvalue configurations of  $\boldsymbol{\Lambda}_i$ , following Zhang et al. (2024a):

- **Homogeneous (Homo):**  $\{1, 99, 4998\}$ ,  $\{2, 100, 4999\}$ ,  $\{3, 101, 5000\}$ ,
- **Heterogeneous (Hetero):**  $\{1, 2, 3\}$ ,  $\{99, 100, 101\}$ ,  $\{4998, 4999, 5000\}$ ,

for  $i = 1, 2, 3$ .

We use fixed learning rates for both SGD and sign-based optimization methods. For SGD, the learning rate is set to (Nesterov, 2013)

$$\eta = \frac{2}{\lambda_{\min} + \lambda_{\max}},$$

where  $\lambda_{\min} = 1$  and  $\lambda_{\max} = 5000$  denote the smallest and largest eigenvalues of  $\mathbf{H}$ , respectively. For SignSGD, we adopt the theoretically optimal learning rate derived from our analysis, as detailed in Appendix D.9.1.

Figure S.9 shows the evolution of the  $\ell_2$  norm of the gradient during optimization, and Table S.9 reports the weighted Hessian complexities  $\Lambda_G$  and  $\Lambda_P$  defined in Definition 4.4. In this quadratic setting,  $\Lambda_P$  can be computed exactly, whereas  $\Lambda_G$  involves a supremum over the fine-tuning region  $\mathcal{R}_{\text{FT}}$  and cannot be evaluated in closed form. We therefore approximate  $\Lambda_G$  by  $\sup_{t \in \{0, \dots, T\}} \sum_{b=1}^B \frac{\|\nabla L(\boldsymbol{\theta}_t)\|_2^2}{\|\nabla L(\boldsymbol{\theta}_t)\|_2} \|\nabla^2 L(\boldsymbol{\theta}_t)\|_2^{-1}$ .

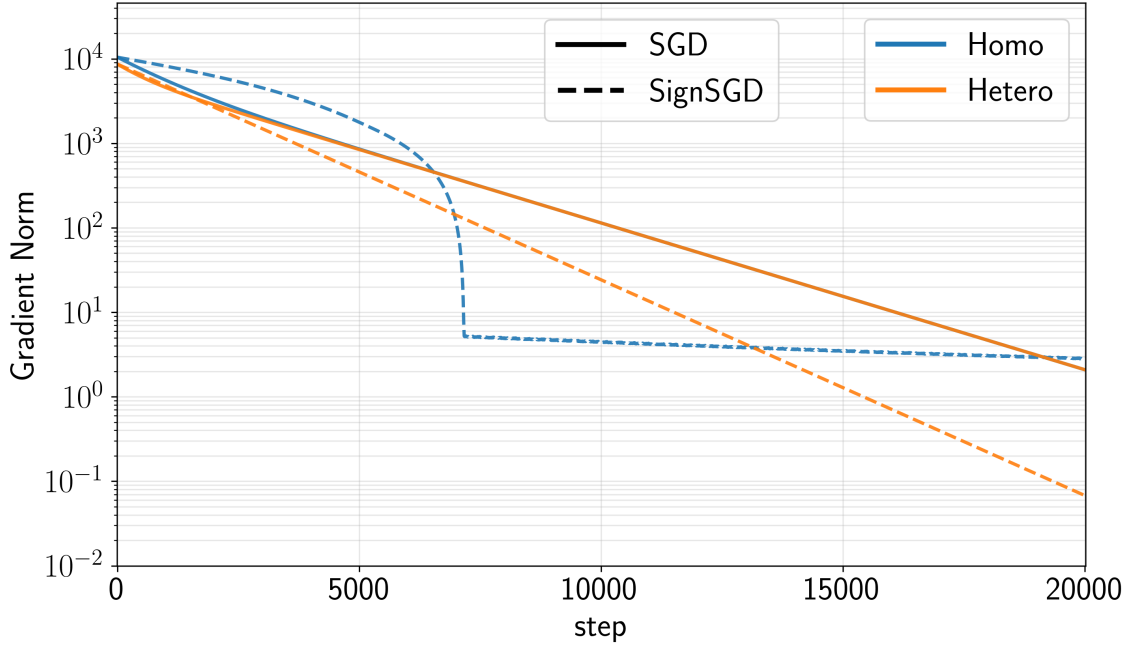
For SGD, the values of  $\Lambda_G$  are nearly identical in the Homo and Hetero settings, and the corresponding gradient norm trajectories exhibit similar behavior. In contrast, for SignSGD,  $\Lambda_P$  is substantially larger in the Homo setting than in the Hetero setting, which is reflected in a slower decay of the gradient norm (except in the early stage) and a larger number of iterations required to reach a small norm. These observations align with our theoretical prediction that the iteration complexity of gradient-based and sign-based methods is characterized by  $\Lambda_G$  and  $\Lambda_P$ , respectively.

Moreover, the pronounced gap between  $\Lambda_G$  and  $\Lambda_P$  in the Hetero setting reflects strong gradient heterogeneity and gradient–Hessian correlation, explaining why sign-based methods are more effective in the Hetero setting. Furthermore, the optimization advantage of Adam over SGD in the heterogeneous quadratic setting reported by Zhang et al. (2024a) is also observed here when comparing SignSGD with SGD.

Table S.9: Values of the weighted Hessian complexities  $\Lambda_G$  and  $\Lambda_P$  in the quadratic model.

|             | Homo      | Hetero    |
|-------------|-----------|-----------|
| $\Lambda_G$ | 4999.9997 | 4999.9997 |
| $\Lambda_P$ | 4999.0000 | 1701.3333 |

<sup>1</sup>While this computation restricts the supremum to the optimization trajectory and is therefore formally smaller than the supremum over  $\boldsymbol{\theta} \in \mathcal{R}_{\text{FT}}$ , the resulting value is nearly maximal over  $\boldsymbol{\theta} \in \mathbb{R}^P$ , indicating that this approximation is inconsequential.

Figure S.9: Evolution of the  $\ell_2$  norm of the gradient in the quadratic model.

### D.9.1 Optimal learning rate and upper bounds

**SignSGD.** We derive the optimal learning rate for SignSGD in the quadratic setting. For the quadratic model, we have  $\delta_D = \rho_H = 0$ . Therefore, from Eq. (9), we obtain

$$\begin{aligned} L(\boldsymbol{\theta}_{t+1}^{\text{Sign}}) - L(\boldsymbol{\theta}_t^{\text{Sign}}) &\leq -\eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 + \frac{\eta_t^2}{2} \Lambda_P P + \frac{\eta_t^2}{2} \delta_D P + \eta_t^3 \frac{\rho_H}{6} P^{3/2} \\ &= -\eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 + \frac{\eta_t^2}{2} \Lambda_P P. \end{aligned}$$

The right-hand side is minimized when

$$\eta_t = \frac{\|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1}{P \Lambda_P},$$

which is the learning rate used in our experiments. the steepest descent with respect to the  $\ell_\infty$ -norm as noted by Kelner et al. (2014); Carlson et al. (2015); Balles et al. (2020).

Using this learning rate, we can derive an upper bound on the iteration complexity for the quadratic model by following the same argument as in the general setting:

$$\mathcal{T}_\varepsilon(\{\boldsymbol{\theta}_t^{\text{Sign}}\}_{t=0}^\infty, L, \|\cdot\|_1) \leq \frac{2(L(\boldsymbol{\theta}_0) - L_*)}{P\varepsilon^2} \Lambda_P. \quad (17)$$

**SGD.** Analogously to the SignSGD case, starting from Eq. (8), we obtain the following inequality:

$$\begin{aligned} L(\boldsymbol{\theta}_{t+1}^{\text{Grad}}) - L(\boldsymbol{\theta}_t^{\text{Grad}}) &\leq -\eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 + \frac{\eta_t^2}{2} \Lambda_G \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 + \frac{\eta_t^2}{2} \delta_D \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 + \eta_t^3 \frac{\rho_H}{6} \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^3 \\ &= -\eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2 + \frac{\eta_t^2}{2} \Lambda_G \|\nabla L(\boldsymbol{\theta}_t^{\text{Grad}})\|_2^2, \end{aligned}$$

where we used  $\delta_D = \rho_H = 0$  for the quadratic model. The right-hand side is minimized when

$$\eta_t = \frac{1}{\Lambda_G}. \quad (18)$$

Using this learning rate, we derive an upper bound on the iteration complexity for the quadratic model by following the same argument as in the general setting:

$$\mathcal{T}_\varepsilon(\{\boldsymbol{\theta}_t^{\text{Grad}}\}_{t=0}^\infty, L, \|\cdot\|_2) \leq \frac{2(L(\boldsymbol{\theta}_0) - L_*)}{P\varepsilon^2} \Lambda_G.$$

This bound has the same form as Eq. (17), with  $\Lambda_G$  replacing  $\Lambda_P$ .

From Table S.9, we observe that  $\Lambda_G$  is approximately equal to  $\lambda_{\max}$ . Therefore, comparing Eq. (18) with the classical optimal learning rate for quadratic objectives,

$$\eta = \frac{2}{\lambda_{\min} + \lambda_{\max}}, \tag{19}$$

we obtain

$$\frac{1}{\Lambda_G} \leq \frac{2}{\lambda_{\min} + \lambda_{\max}}.$$

Empirically, we observed faster convergence when using the larger learning rate given by Eq. (19). Consequently, we adopt this learning rate in our experiments.

### D.10 Applicability beyond fine-tuning settings

To test generalization beyond fine-tuning, we trained nanoGPT from scratch on the Shakespeare dataset. Adam outperformed SGD, and SignSGD remained competitive. We also found that gradient heterogeneity in nanoGPT lies between that of ViT/ResNet and RoBERTa. Despite the different setup, the results align with our analysis.

Table S.10: Training loss for nanoGPT trained from scratch on the Shakespeare dataset. “Min” denotes the lowest observed loss during training, and “Last” denotes the final loss at the end of training.

| Optimizer | Min           | Last          |
|-----------|---------------|---------------|
| Adam      | 0.658 ± 0.009 | 0.687 ± 0.019 |
| SGD       | 0.928 ± 0.120 | 0.964 ± 0.122 |
| SignSGD   | 0.791 ± 0.011 | 0.820 ± 0.017 |

Table S.11: Gini coefficient of gradient norms for nanoGPT on the Shakespeare dataset. A higher Gini coefficient indicates greater gradient heterogeneity.

| Model (Dataset)       | Gini Coefficient |
|-----------------------|------------------|
| nanoGPT (Shakespeare) | 0.609 ± 0.004    |

## E Discussion on momentum in SignSGD

The impact of the momentum term used in Adam has not been considered in the analysis so far. However, in sample-wise training, the presence of a momentum term significantly affects the updates of the linear head, particularly for the bias term.

**Model.** The model  $\mathbf{f}$  comprises a pre-trained feature extractor  $\phi(\cdot) : \mathcal{X} \rightarrow \mathbb{R}^h$  and a linear head with weight  $\mathbf{V} \in \mathbb{R}^{C \times h}$  and bias  $\mathbf{b} \in \mathbb{R}^C$ . The output is given by  $\mathbf{f}(\mathbf{x}) = \mathbf{V}\phi(\mathbf{x}) + \mathbf{b}$ .

**Proposition E.1** (SignSGD without momentum). *Let  $\Delta^S\theta$  and  $\Delta^F\theta$  denote the one-epoch updates of a parameter  $\theta$  during sample-wise and full-batch training, respectively. For a linear head trained using the cross-entropy loss and SignSGD with a learning rate  $\eta$ , the updates are as follows:*

For the bias term  $b_k$ :

$$\Delta^S b_k = -\frac{\eta}{N} \sum_{i=1}^N (1 - 2 \cdot \mathbb{1}[y^{(i)} = k]), \quad \Delta^F b_k = -\eta \operatorname{sign} \left( \sum_{i=1}^N \delta_{p_k}^{(i)} \right),$$

and for the weight matrix  $V_{k,l}$ :

$$\Delta^S V_{k,l} = -\frac{\eta}{N} \left( \sum_{y^{(i)} \neq k} s_l^{(i)} - \sum_{y^{(i)} = k} s_l^{(i)} \right), \quad \Delta^F V_{k,l} = -\eta \operatorname{sign} \left( \sum_{i=1}^N \phi(\mathbf{x}^{(i)})_l \delta_{p_k}^{(i)} \right),$$

where  $\delta_{p_k}^{(i)} := \sigma_{\text{SM}}(\mathbf{f}(\mathbf{x}^{(i)}))_k - \mathbb{1}[k = y^{(i)}]$  represents the prediction error for the  $i$ -th sample and class  $k$  and  $s_l^{(i)} := \operatorname{sign}(\phi(\mathbf{x}^{(i)})_l)$  is the sign of the  $l$ -th element of the feature embedding  $\phi(\mathbf{x}^{(i)})_l$ .

**Sign-alignment causes large updates.** In full-batch training, the updates  $\Delta^F b_k$  and  $\Delta^F V_{k,l}$  depend on the model predictions. Because the signs of these updates vary across epochs, these updates remain small. In contrast, in sample-wise training, update signs can align across epochs, resulting in disproportionately large updates. This effect is particularly pronounced for the bias term  $\Delta^S b_k$ , which is independent of model predictions and grows with the number of classes. Similarly, the sign of  $\Delta^S V_{k,l}$ , which depends on the feature extractor output  $\phi(\mathbf{x}^{(i)})$ , may align across epochs.

**Momentum resolves the issue.** Excessively large updates can cause training instability and incorrect predictions. Although the proposition specifically addresses sample-wise updates, similar challenges can arise in batch training. Momentum, which estimates the full-batch gradient using exponential moving averages, effectively mitigates this problem.

## E.1 Experimental results

We show the norm of the linear head for different datasets, models, and optimizers. The results indicate that when the number of classes is large, the bias term of the linear head exhibits a larger norm with SignSGD without momentum compared to other optimizers. In contrast, the weight norm does not necessarily increase under the same conditions, even with SignSGD without momentum. This observation aligns with the theoretical analysis in Proposition E.1, which suggests that a large number of classes leads to an increase in the bias term norm, while the weight norm is influenced by the sign of the feature extractor outputs.

## E.2 Proof of Proposition E.1

*Proof.* The partial derivative of the bias and the weight matrix with the cross-entropy loss is given by:

$$\begin{aligned} \frac{\partial \ell(\mathbf{f}(\mathbf{x}^{(i)}), y^{(i)})}{\partial b_k} &= \frac{\partial \ell(\mathbf{f}(\mathbf{x}^{(i)}), y^{(i)})}{\partial \mathbf{f}(\mathbf{x}^{(i)})} \frac{\partial \mathbf{f}(\mathbf{x}^{(i)})}{\partial b_k} \\ &= \frac{\partial \ell(\mathbf{f}(\mathbf{x}^{(i)}), y^{(i)})}{\partial \mathbf{f}(\mathbf{x}^{(i)})} \frac{\partial \mathbf{V} \phi(\mathbf{x}^{(i)}) + \mathbf{b}}{\partial b_k} \\ &= (\sigma_{\text{SM}}(\mathbf{f}(\mathbf{x}^{(i)})) - \mathbf{e}^{(y^{(i)})})^\top \mathbf{e}^{(k)} \\ &= \sigma_{\text{SM}}(\mathbf{f}(\mathbf{x}^{(i)}))_k - \mathbb{1}[k = y^{(i)}] \\ \frac{\partial \ell(\mathbf{f}(\mathbf{x}^{(i)}), y^{(i)})}{\partial V_{k,l}} &= \frac{\partial \ell(\mathbf{f}(\mathbf{x}^{(i)}), y^{(i)})}{\partial \mathbf{f}(\mathbf{x}^{(i)})} \frac{\partial \mathbf{V} \phi(\mathbf{x}^{(i)}) + \mathbf{b}}{\partial V_{k,l}} \\ &= (\sigma_{\text{SM}}(\mathbf{f}(\mathbf{x}^{(i)})) - \mathbf{e}^{(y^{(i)})})^\top \phi(\mathbf{x}^{(i)})_l \mathbf{e}^{(k)} \\ &= \phi(\mathbf{x}^{(i)})_l (\sigma_{\text{SM}}(\mathbf{f}(\mathbf{x}^{(i)}))_k - \mathbb{1}[k = y^{(i)}]) \end{aligned}$$

The one-epoch updates of the bias and the weight matrix with the sample-wise training are given by:

$$\Delta^S b_k = -\frac{\eta}{N} \sum_{i=1}^N \operatorname{sign} \left( \frac{\partial \ell(\mathbf{f}(\mathbf{x}^{(i)}), y^{(i)})}{\partial b_k} \right)$$

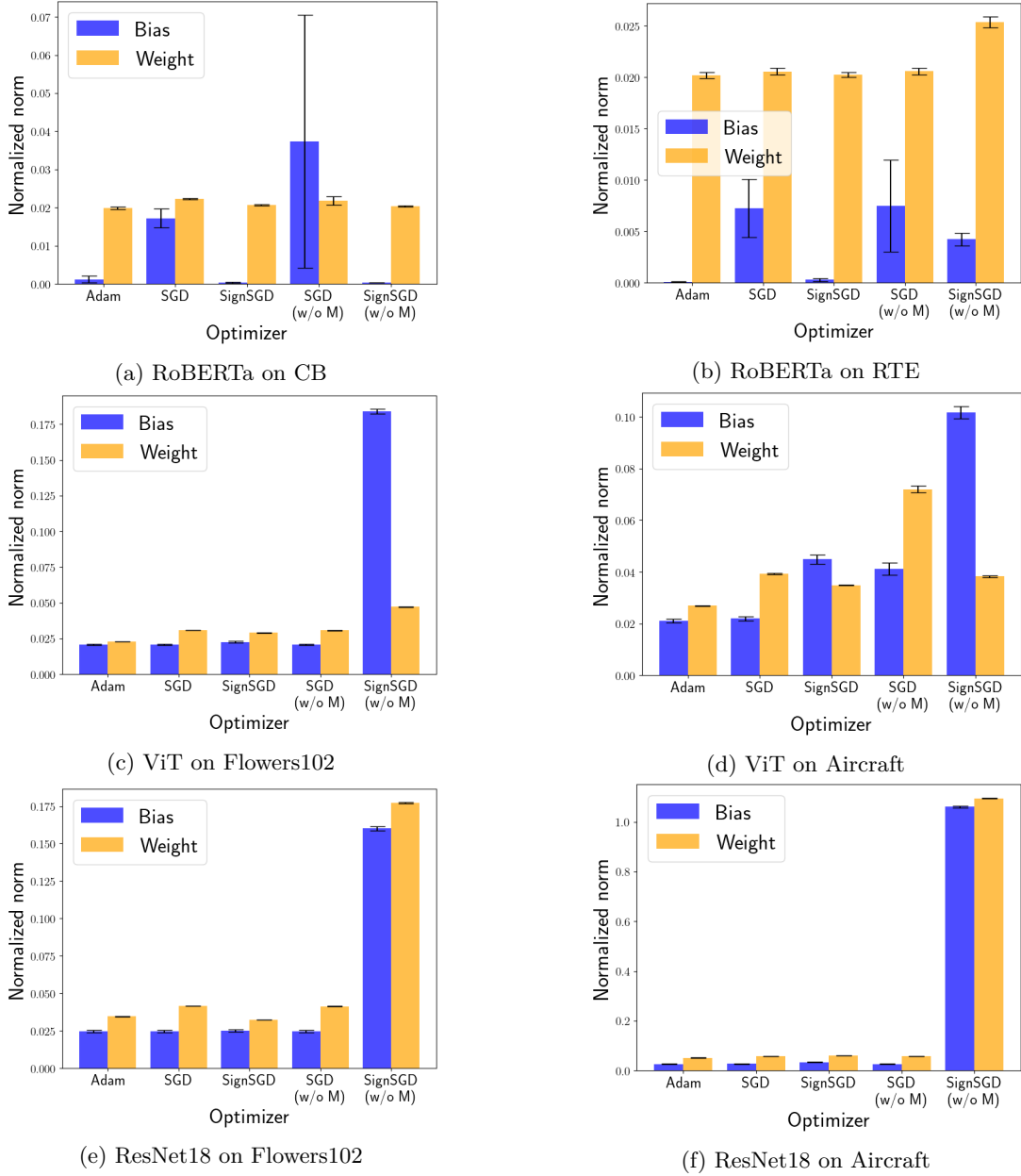


Figure S.10: Norm of the linear head.

$$\begin{aligned}
&= -\frac{\eta}{N} \sum_{i=1}^N \text{sign} \left( \sigma_{\text{SM}}(\mathbf{f}(\mathbf{x}^{(i)}))_k - \mathbb{1}[k = y^{(i)}] \right) \\
&= -\frac{\eta}{N} \sum_{i=1}^N (1 - 2 \cdot \mathbb{1}[y^{(i)} = k])
\end{aligned}$$

and

$$\Delta^S V_{k,l} = -\frac{\eta}{N} \sum_{i=1}^N \text{sign} \left( \frac{\partial \ell(\mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}))}{\partial V_{k,l}} \right)$$

$$\begin{aligned}
&= -\frac{\eta}{N} \sum_{i=1}^N \text{sign} \left( \phi(\mathbf{x}^{(i)})_l (\boldsymbol{\sigma}_{\text{SM}}(\mathbf{f}(\mathbf{x}^{(i)}))_k - \mathbb{1}[k = y^{(i)}]) \right) \\
&= -\frac{\eta}{N} \sum_{i=1}^N \text{sign} \left( \phi(\mathbf{x}^{(i)})_l \right) \text{sign} \left( \boldsymbol{\sigma}_{\text{SM}}(\mathbf{f}(\mathbf{x}^{(i)}))_k - \mathbb{1}[k = y^{(i)}] \right) \\
&= -\frac{\eta}{N} \left( \sum_{y^{(i)} \neq k} \text{sign} \left( \phi(\mathbf{x}^{(i)})_l \right) - \sum_{y^{(i)} = k} \text{sign} \left( \phi(\mathbf{x}^{(i)})_l \right) \right)
\end{aligned}$$

The one-epoch updates of the bias and the weight matrix with the full-batch training are given by:

$$\begin{aligned}
\Delta^{\text{F}} b_k &= -\eta \text{sign} \left( \frac{1}{N} \sum_{i=1}^N \frac{\partial \ell(\mathbf{f}(\mathbf{x}^{(i)}), y^{(i)})}{\partial b_k} \right) \\
&= -\eta \text{sign} \left( \frac{1}{N} \sum_{i=1}^N \left( \boldsymbol{\sigma}_{\text{SM}}(\mathbf{f}(\mathbf{x}^{(i)}))_k - \mathbb{1}[k = y^{(i)}] \right) \right) \\
&= -\eta \text{sign} \left( \sum_{i=1}^N \delta_{p_k}^{(i)} \right)
\end{aligned}$$

and

$$\begin{aligned}
\Delta^{\text{F}} V_{k,l} &= -\eta \text{sign} \left( \frac{1}{N} \sum_{i=1}^N \frac{\partial \ell(\mathbf{f}(\mathbf{x}^{(i)}), y^{(i)})}{\partial V_{k,l}} \right) \\
&= -\eta \text{sign} \left( \frac{1}{N} \sum_{i=1}^N \phi(\mathbf{x}^{(i)})_l (\boldsymbol{\sigma}_{\text{SM}}(\mathbf{f}(\mathbf{x}^{(i)}))_k - \mathbb{1}[k = y^{(i)}]) \right) \\
&= -\eta \text{sign} \left( \sum_{i=1}^N \phi(\mathbf{x}^{(i)})_l \delta_{p_k}^{(i)} \right).
\end{aligned}$$

□

## F More discussion on Transformers

In this section, we provide additional discussion on the gradient heterogeneity in Transformers, focusing on the self-attention mechanism.

**Additional notation.** The  $k$ -th standard basis vector is denoted by  $\mathbf{e}^{(k)}$  with  $e_l^{(k)} = \delta_{kl}$ , where  $\delta_{kl}$  is the Kronecker delta. Function  $\text{vec}(\cdot)$  denotes row-wise vectorization. Frobenius norm and the Kronecker product is denoted by  $\|\cdot\|_F$  and  $\otimes$ , respectively.

### F.1 Transformer architecture

The Transformer architecture (Vaswani, 2017) relies on the self-attention mechanism, which assigns importance to each token in the input sequence.

For an input sequence of  $n$  tokens, each of dimension  $d$ , represented by  $\mathbf{X} \in \mathbb{R}^{n \times d}$ , single-head self-attention is defined as:

$$\text{SA}(\mathbf{X}) := \sigma_{\text{SM}} \left( \frac{\mathbf{X} \mathbf{W}_Q (\mathbf{X} \mathbf{W}_K)^\top}{\sqrt{d_k}} \right) \mathbf{X} \mathbf{W}_V,$$

where  $\mathbf{W}_Q, \mathbf{W}_K \in \mathbb{R}^{d \times d_k}$  and  $\mathbf{W}_V \in \mathbb{R}^{d \times d_v}$  are learnable projection matrices for queries, keys, and values, respectively. Multi-head attention concatenates the outputs of parallel single-head self-attention mechanisms and applies a linear transformation, followed by a feed-forward network.

### F.2 Gradient of self-attention mechanism

We analyze the gradients in self-attention, focusing on the value and query/key weight matrices. Using Lemma A.2 from Noci et al. (2022), the Frobenius norms of these gradients are:

$$\begin{aligned} \left\| \frac{\partial \text{SA}(\mathbf{X})}{\partial \mathbf{W}_V} \right\|_F &= \|\mathbf{P} \mathbf{X} \otimes \mathbf{I}_{d_v}\|_F \\ &\leq \underbrace{\sqrt{d_v} \|\mathbf{P}\|_F \|\mathbf{X}\|_F}_{=: \mathcal{U}_V}, \end{aligned} \quad (20)$$

$$\begin{aligned} &\left\| \frac{\partial \text{SA}(\mathbf{X})}{\partial \mathbf{W}_Q} \right\|_F \\ &= \|(\mathbf{I}_n \otimes \mathbf{W}_V \mathbf{X}^\top) \frac{\partial \mathbf{P}}{\partial \mathbf{M}} \frac{\mathbf{X} \otimes \mathbf{X} \mathbf{W}_K}{\sqrt{d_k}}\|_F \\ &\leq \underbrace{\sqrt{n} \|\mathbf{W}_V \mathbf{X}^\top\|_F \left\| \frac{\partial \mathbf{P}}{\partial \mathbf{M}} \right\|_F \frac{\|\mathbf{X}\|_F \|\mathbf{X} \mathbf{W}_K\|_F}{\sqrt{d_k}}}_{=: \mathcal{U}_Q}, \end{aligned} \quad (21)$$

where  $\mathbf{M} := \mathbf{X} \mathbf{W}_Q \mathbf{W}_K^\top \mathbf{X}^\top / \sqrt{d_k}$ ,  $\mathbf{P} := \sigma_{\text{SM}}(\mathbf{M})$ , and  $\mathcal{U}_V$  and  $\mathcal{U}_Q$  represent the upper bounds for the gradients of the value and query weight matrices, respectively. The derivation of the gradient for the key weight matrix is omitted, as it is analogous to that of the query weight matrix.

Focusing on the attention matrix  $\mathbf{P}$ , we derive the following result.

**Proposition F.1** (Gradients and attention matrices). *In Transformers, one-hot attention matrices uniquely maximize the upper bound of the Frobenius norm of the gradient with respect to the value weight matrix  $\mathcal{U}_V$  and uniquely minimize that with respect to the query weight matrix  $\mathcal{U}_Q$ , as follows:*

$$\arg \max_{\mathbf{P}} \mathcal{U}_V = \arg \min_{\mathbf{P}} \mathcal{U}_Q = \mathcal{P}_{\text{one-hot}},$$

where

$$\mathcal{P}_{\text{one-hot}} := \{\mathbf{P} \mid \forall i, \exists k_i \text{ s.t. } \mathbf{P}_{i,:} = \mathbf{e}^{(k_i)}\}$$

is the set of one-hot matrices.

The proof of the proposition is provided in Appendix F.4. The statement about the query weight matrix also applies to the key weight matrix due to their analogous gradients. The proposition demonstrates that the gradients of the value and query/key weight matrices exhibit opposing behaviors with respect to one-hot attention matrices: the gradient of the value weight matrix is maximized, while those of the query/key weight matrices are minimized.

Previous studies (Noci et al., 2022; Wang et al., 2021) observed that the gradient of the value weight matrix is typically larger than those of the query/key weight matrices, consistent with our experimental findings in Section 5.2. Together with Proposition F.1, these results suggest that attention matrices close to one-hot amplify gradient heterogeneity in the self-attention mechanism.

### F.3 Uniformity of the attention matrix

In Figure S.11, we compare the attention matrices of pre-trained RoBERTa and ViT. The attention matrix of ViT is more uniform than that of RoBERTa, reflecting the differences between NLP and vision tasks. In NLP, the use of special tokens and stronger interrelations between input tokens lead to less uniform attention, with only a few tokens receiving attention (Clark, 2019). Conversely, vision tasks, which prioritize holistic information (Torralba, 2003; Rabinovich et al., 2007; Shotton et al., 2009), produce more uniform attention matrices, where all tokens are attended to. This observation aligns with Hyeon-Woo et al. (2023), who also reported uniform attention matrices in ViT. Notably, more uniform attention matrices are farther from one-hot matrices, indicating reduced dominance by individual tokens.

Combined with the analysis in Appendix F.2, which shows that attention matrices closer to one-hot matrices amplify gradient heterogeneity, this suggests that gradient heterogeneity in the self-attention mechanism is more pronounced in NLP tasks than in vision tasks.

### F.4 Proof of Proposition F.1

*Proof of  $\mathcal{U}_V$ .* As defined in Eq.(20), the upper bound of the gradient is given by:

$$\mathcal{U}_V = \sqrt{d_v} \|\mathbf{P}\|_F \|\mathbf{X}\|_F.$$

We observe that:

$$\begin{aligned} \arg \max_{\mathbf{P}} \mathcal{U}_V &= \arg \max_{\mathbf{P}} \|\mathbf{P}\|_F \\ &= \arg \max_{\mathbf{P}} \|\mathbf{P}\|_F^2 \\ &= \arg \max_{\mathbf{P}} \sum_{i=1}^n \|\mathbf{P}_{i,:}\|_2^2. \end{aligned}$$

Since the rows of the attention matrix are independent, we focus on the  $i$ -th row. The  $i$ -th row of the attention matrix satisfies the following constraints:

$$1 \leq j \leq n, \quad P_{i,j} \geq 0, \quad \sum_{j=1}^n P_{i,j} = 1.$$

We define the Lagrangian function as:

$$\mathcal{L}_V = -\sum_{j=1}^n P_{i,j}^2 - \sum_{j=1}^n \mu_j P_{i,j} + \lambda \left( \sum_{j=1}^n P_{i,j} - 1 \right),$$

where  $\lambda$  and  $\mu_j$  are the Lagrange multipliers. To minimize the Lagrangian function, the solution must satisfy the following KKT conditions:

$$\frac{\partial \mathcal{L}_V}{\partial P_{i,j}} = -2P_{i,j} - \mu_j + \lambda = 0, \quad 1 \leq j \leq n, \quad (22)$$

$$\sum_{j=1}^n P_{i,j} - 1 = 0, \quad (23)$$

$$P_{i,j} \geq 0, \quad 1 \leq j \leq n, \quad (24)$$

$$\mu_j \geq 0, \quad 1 \leq j \leq n, \quad (25)$$

$$\mu_j P_{i,j} = 0, \quad 1 \leq j \leq n. \quad (26)$$

From Equations (23) and (24), it follows that  $P_{i,j} > 0$  for some  $j$ . Let  $k$  ( $1 \leq k \leq n$ ) denote the number of non-zero elements in  $\mathbf{P}_{i,:}$ , and suppose  $P_{i,j_l} > 0$  for  $1 \leq l \leq k$ . From Equation (26), we have  $\mu_{j_l} = 0$ , and thus, from Equation (22), we deduce that  $P_{i,j_l} = \frac{\lambda}{2}$  for  $1 \leq l \leq k$ . Using Equation (23), we get  $\sum_{l=1}^k \frac{\lambda}{2} = 1$ , which gives  $\lambda = 2/k$ . For  $j \notin \{j_l \mid 1 \leq l \leq k\}$ , we have  $P_{i,j} = 0$  and  $\mu_j = \lambda = 2/k$ , satisfying Eq.(25).

With  $k$  non-zero elements of  $\mathbf{P}_{i,:}$ , the value of the Lagrangian function becomes  $-\sum_{j=1}^n P_{i,j}^2 = -\sum_{l=1}^k (\frac{\lambda}{2})^2 = -\frac{\lambda^2}{4}k = -\frac{1}{k}$ . The minimum value of the Lagrangian function is achieved if and only if  $k = 1$ , which implies  $\mathbf{P}_{i,:} = \mathbf{e}^{(k_i)}$  for some  $k_i$ . Therefore, we conclude:

$$\arg \max_{\mathbf{P}} \mathcal{U}_V = \{\mathbf{P} \mid \forall i, \exists k_i \text{ s.t. } \mathbf{P}_{i,:} = \mathbf{e}^{(k_i)}\}.$$

□

*Proof of  $\mathcal{U}_Q$ .* As defined in Eq.(21), the upper bound of the gradient is given by:

$$\mathcal{U}_Q = \sqrt{n} \|\mathbf{W}_V \mathbf{X}^\top\|_F \left\| \frac{\partial \mathbf{P}}{\partial \mathbf{M}} \right\|_F \frac{\|\mathbf{X}\|_F \|\mathbf{X} \mathbf{W}_K\|_F}{\sqrt{d_k}}.$$

The partial derivative is expressed as:

$$\begin{aligned} \frac{\partial \mathbf{P}}{\partial \mathbf{M}} &= \frac{\partial \sigma_{\text{SM}}(\mathbf{M})}{\partial \mathbf{M}} \\ &= \text{blockdiag}(\left\{ \frac{\partial \sigma_{\text{SM}}(\mathbf{M}_{i,:})}{\partial \mathbf{M}_{i,:}} \right\}_{i=1}^n) \\ &= \text{blockdiag}(\left\{ \text{diag}(\mathbf{P}_{i,:}) - \mathbf{P}_{i,:} \mathbf{P}_{i,:}^\top \right\}_{i=1}^n). \end{aligned}$$

Considering the attention matrix  $\mathbf{P}$ , we obtain:

$$\begin{aligned} \arg \min_{\mathbf{P}} \mathcal{U}_Q &= \arg \min_{\mathbf{P}} \left\| \frac{\partial \mathbf{P}}{\partial \mathbf{M}} \right\|_F \\ &= \arg \min_{\mathbf{P}} \sum_{i=1}^n \left\| \text{diag}(\mathbf{P}_{i,:}) - \mathbf{P}_{i,:} \mathbf{P}_{i,:}^\top \right\|_F^2. \end{aligned}$$

As in the proof of  $\mathcal{U}_V$ , we focus on the value of the  $i$ -th row:

$$\left\| \text{diag}(\mathbf{P}_{i,:}) - \mathbf{P}_{i,:} \mathbf{P}_{i,:}^\top \right\|_F^2 = \sum_{j=1}^n (P_{i,j} - P_{i,j}^2)^2 + \sum_{j \neq l} P_{i,j}^2 P_{i,l}^2,$$

subject to the constraints  $1 \leq j \leq n$ ,  $P_{i,j} \geq 0$ ,  $\sum_{j=1}^n P_{i,j} = 1$ . Since both the first term and the second term are non-negative, the minimum value is attained if and only if both terms are 0. This condition is satisfied if  $\mathbf{P}_{i,:}$  is a one-hot vector. Conversely, if  $\mathbf{P}_{i,:}$  is not a one-hot vector, the second term becomes positive, and the minimum value cannot be attained. Thus, we have shown that the minimum value of the objective function is achieved if and only if  $\mathbf{P}_{i,:}$  is a one-hot vector. Therefore:

$$\arg \min_{\mathbf{P}} \mathcal{U}_Q = \{\mathbf{P} \mid \forall i, \exists k_i \text{ s.t. } \mathbf{P}_{i,:} = \mathbf{e}^{(k_i)}\}.$$

□

## F.5 Experimental results

**Heatmap of attention matrices.** In Figure S.11, we show the attention matrices computed from pre-trained models. These matrices are calculated for a randomly sampled sequence from the training data and are averaged across all heads.

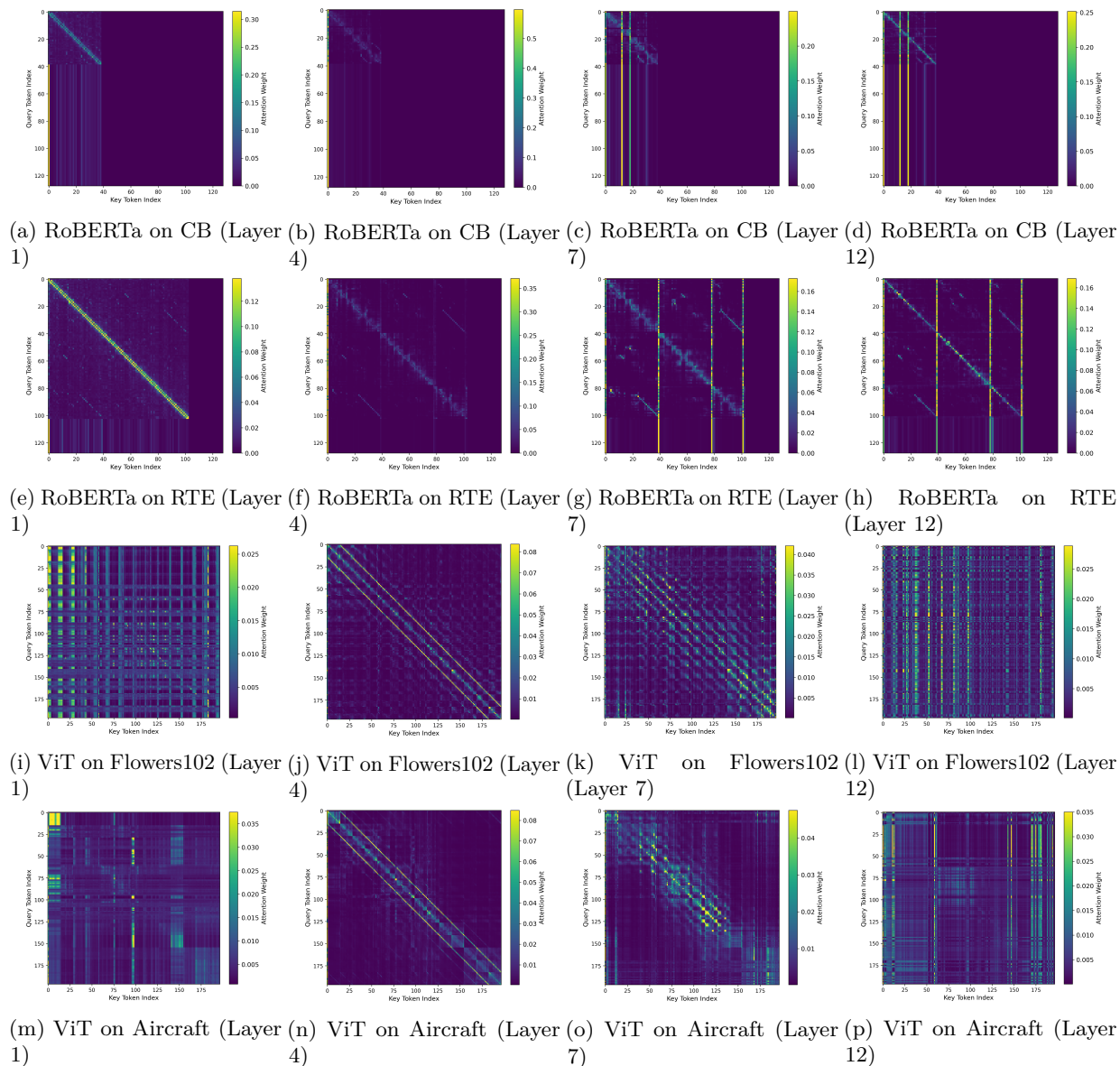


Figure S.11: Attention matrices of the pre-trained RoBERTa and ViT.

**Gradient and entropy of attention matrices.** In Figure S.12 (a) and (c), we show the ratio of the mean entropy relative to the maximum entropy of the attention matrix for each layer of the Transformer model. Error bars indicate the standard deviation. Specifically, we plot:

$$\frac{1}{HNS} \sum_{h=1}^H \sum_{i=1}^N \sum_{s=1}^S \left( \sum_{j=1}^S A_{s,j}^{(i,h,l)} \log(A_{s,j}^{(i,h,l)}) / \log(S) \right),$$

for each layer  $l$ , where  $H$  is the number of heads,  $S$  is the sequence length, and  $\mathbf{A}^{(i,h,l)} \in \mathbb{R}^{S \times S}$  is the attention matrix of the  $h$ -th head in the  $l$ -th layer for sample  $\mathbf{x}^{(i)}$ .

In Figure S.12 (b) and (d), we show the ratio of the mean gradient norm relative to the sum of the gradient norms of the attention matrix for each layer. Specifically, we plot:

$$\frac{G_p^{(l)}}{G_Q^{(l)} + G_K^{(l)} + G_V^{(l)}},$$

for each layer  $l$  and  $p \in \{Q, K, V\}$ , where  $G_Q^{(l)}$ ,  $G_K^{(l)}$ , and  $G_V^{(l)}$  are the full-batch gradient norms of the query, key, and value weight matrices in the  $l$ -th layer of the Transformer model, respectively.

The results show that the entropy of the attention matrix is higher in RoBERTa than in ViT, and the gradient norm of the attention matrix is more heterogeneous in RoBERTa than in ViT. This observation is consistent with the theoretical analysis in Appendix F.3.

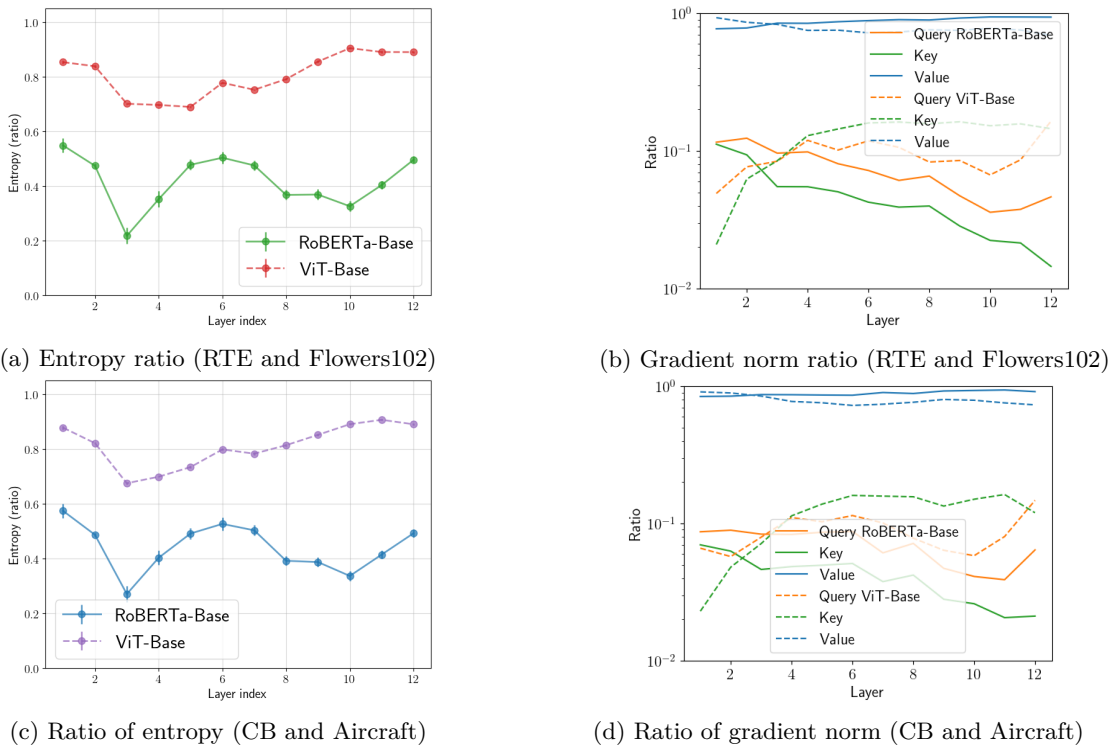


Figure S.12: Comparison of entropy and gradient norms in attention matrices for RoBERTa and ViT. (a) and (c): the ratio of entropy relative to the maximum possible entropy. (b) and (d): the ratio of the gradient norm for self-attention parameters relative to the total gradient norm.

## G More discussion on the sign-based sequence the stochastic setting

In this section, we further examine the iteration complexity of the sign-based sequence under stochastic settings. Specifically, we present iteration complexity results that account for a learning rate adapted to the noise level.

**Theorem G.1.** *Assume that  $\delta_D < \Lambda_P/3$ ,  $\varepsilon < \frac{5\Lambda_P^2}{3(1-2\sigma_2)\rho_H\sqrt{P}}$ , and  $\sigma_2 < \frac{1}{2}$  hold and that the learning rate at time  $t$  satisfies  $\eta_t = \zeta_t \min(\frac{3(1-2\sigma_2)\|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1}{5\Lambda_P P}, \sqrt{\frac{3(1-2\sigma_2)\|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1}{5\rho_H P^{3/2}}})$ , where  $\zeta_t \in [\zeta_0, 1]$ . Then, the iteration complexity for the sign-based sequence the stochastic setting are bounded as follows.*

$$\mathcal{T}_\varepsilon(\{\boldsymbol{\theta}_t^{\text{Sign}}\}_{t=0}^\infty, L, \|\cdot\|_1) \leq \frac{20(L(\boldsymbol{\theta}_0) - L_*)}{3(1-2\sigma_2)^2 P \varepsilon^2 \zeta_0} \Lambda_P.$$

*Proof.* We start with Eq. (16) in Appendix A.4. Let  $\varepsilon < \frac{\alpha\Lambda_P^2}{\rho_H\sqrt{P}}$  and set the learning rate as  $\eta_t = \zeta_t \min(\frac{\|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1}{\alpha\Lambda_P P}, \sqrt{\frac{\|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1}{\alpha\rho_H P^{3/2}}})$ , where  $\zeta_t \in [\zeta_0, 1]$  and  $\alpha > \frac{5}{6(1-2\sigma_2)}$ . Then, we have:

$$\begin{aligned} & \mathbb{E} \left[ L(\boldsymbol{\theta}_{t+1}^{\text{Sign}}) - L(\boldsymbol{\theta}_t^{\text{Sign}}) \mid \boldsymbol{\theta}_t^{\text{Sign}} \right] \\ & \leq -\eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 + \frac{\eta_t^2}{2} \Lambda_P P + \frac{\eta_t^2}{2} \delta_D P + \eta_t^3 \frac{\rho_H}{6} P^{3/2} + 2\sigma_2 \eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 \\ & \leq -\eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 + \frac{\eta_t}{2\alpha} \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 + \frac{\eta_t}{6\alpha} \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 + \frac{\eta_t}{6\alpha} \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 + 2\sigma_2 \eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 \\ & \quad (\text{From } \eta_t \leq \min(\frac{\|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1}{\alpha\Lambda_P P}, \sqrt{\frac{\|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1}{\alpha\rho_H P^{3/2}}}) \text{ and } \delta_D < \Lambda_P/3) \\ & = -\frac{(6\alpha(1-2\sigma_2) - 5)\eta_t}{6\alpha} \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 \end{aligned}$$

Assume that the probability of the event  $\mathcal{E}(T) = \{\forall s \leq T, \|\nabla L(\boldsymbol{\theta}_s^{\text{Sign}})\|_1 \geq P\varepsilon\}$  satisfies  $\mathbb{P}(\mathcal{E}(T)) \geq \frac{1}{2}$ . By applying the telescoping sum and taking expectations, and noting that  $\boldsymbol{\theta}_0 = \boldsymbol{\theta}_0^{\text{Sign}}$ , we have:

$$\begin{aligned} & \mathbb{E} \left[ L(\boldsymbol{\theta}_T^{\text{Sign}}) \right] - L(\boldsymbol{\theta}_0) \\ & \leq -\frac{(6\alpha(1-2\sigma_2) - 5)\eta_t}{6\alpha} \sum_{t=0}^{T-1} \mathbb{E} \left[ \eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 \right] \\ & = -\frac{(6\alpha(1-2\sigma_2) - 5)\eta_t}{6\alpha} \sum_{t=0}^{T-1} \left( \mathbb{E} \left[ \eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 \mid \mathcal{E}(T) \right] \mathbb{P}(\mathcal{E}(T)) + \mathbb{E} \left[ \bar{\eta}_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 \mid \overline{\mathcal{E}(T)} \right] \mathbb{P}(\overline{\mathcal{E}(T)}) \right) \\ & \leq -\frac{(6\alpha(1-2\sigma_2) - 5)\eta_t}{6\alpha} \sum_{t=0}^{T-1} \mathbb{E} \left[ \eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 \mid \mathcal{E}(T) \right] \mathbb{P}(\mathcal{E}(T)) \\ & \leq -\frac{(6\alpha(1-2\sigma_2) - 5)\eta_t}{12\alpha} \sum_{t=0}^{T-1} \mathbb{E} \left[ \eta_t \|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1 \mid \mathcal{E}(T) \right] \\ & \leq -\frac{(6\alpha(1-2\sigma_2) - 5)\eta_t \zeta_0}{12\alpha} \sum_{t=0}^{T-1} \mathbb{E} \left[ \min\left(\frac{\|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1^2}{\alpha\Lambda_P P}, \frac{\|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1^{3/2}}{\sqrt{\alpha\rho_H P^{3/2}}}\right) \mid \mathcal{E}(T) \right] \\ & \quad (\text{From } \eta_t \geq \zeta_0 \min(\frac{\|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1}{\alpha\Lambda_P P}, \sqrt{\frac{\|\nabla L(\boldsymbol{\theta}_t^{\text{Sign}})\|_1}{\alpha\rho_H P^{3/2}}})) \\ & \leq -\frac{(6\alpha(1-2\sigma_2) - 5)\eta_t \zeta_0}{12\alpha} \sum_{t=0}^{T-1} \min\left(\frac{P\varepsilon^2}{\alpha\Lambda_P}, P\varepsilon \sqrt{\frac{\varepsilon}{\alpha\rho_H P^{1/2}}}\right) \end{aligned}$$

$$= -\frac{(6\alpha(1-2\sigma_2)-5)TP\varepsilon^2\zeta_0}{12\alpha^2\Lambda_P} \quad (\text{From } \varepsilon < \frac{\alpha\Lambda_P^2}{\rho_H\sqrt{P}}).$$

Therefore, we have:

$$\begin{aligned} T &\leq \frac{12\alpha^2(L(\boldsymbol{\theta}_0) - \mathbb{E}[L(\boldsymbol{\theta}_T^{\text{Sign}})])}{(6\alpha(1-2\sigma_2)-5)P\varepsilon^2\zeta_0} \Lambda_P \\ &\leq \frac{12\alpha^2(L(\boldsymbol{\theta}_0) - L_*)}{(6\alpha(1-2\sigma_2)-5)P\varepsilon^2\zeta_0} \Lambda_P. \end{aligned}$$

This means that when we take  $T > \frac{12\alpha^2(L(\boldsymbol{\theta}_0)-L_*)}{(6\alpha(1-2\sigma_2)-5)P\varepsilon^2\zeta_0} \Lambda_P$ , we have  $\mathbb{P}(\mathcal{E}(T)) < \frac{1}{2}$ . Therefore, we have

$$\mathcal{T}_\varepsilon(\{\boldsymbol{\theta}_t^{\text{Sign}}\}_{t=0}^\infty, L, \|\cdot\|_1) \leq \frac{12\alpha^2(L(\boldsymbol{\theta}_0) - L_*)}{(6\alpha(1-2\sigma_2)-5)P\varepsilon^2\zeta_0} \Lambda_P,$$

for any  $\alpha > \frac{5}{6(1-2\sigma_2)}$ . Setting  $\alpha = \frac{5}{3(1-2\sigma_2)}$  to minimize the right-hand side completes the proof.  $\square$

## H Additional background and related work

### H.1 Steepest descent

Beyond the descent direction, steepest descent can be formulated as a full update obtained by minimizing a local smoothness-based upper bound of the objective (Bernstein & Newhouse, 2024).

We say that  $L$  is  $L_p$ -smooth if its gradient is Lipschitz continuous with respect to the  $\ell_p$  norm, that is,

$$\|\nabla L(\boldsymbol{\theta}) - \nabla L(\boldsymbol{\theta}')\|_q \leq L_p \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_p,$$

where  $\|\cdot\|_q$  is the dual norm of  $\|\cdot\|_p$ . For an  $L_p$ -smooth function in a deterministic setting, the steepest descent update is given by

$$\boldsymbol{\theta}_{t+1} \in \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^P} \left( \langle \nabla L(\boldsymbol{\theta}_t), \boldsymbol{\theta} - \boldsymbol{\theta}_t \rangle + \frac{L_p}{2} \|\boldsymbol{\theta} - \boldsymbol{\theta}_t\|_p^2 \right).$$

For the  $\ell_2$  norm, this reduces to the standard gradient descent update

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \frac{1}{L_2} \nabla L(\boldsymbol{\theta}_t),$$

whereas for the  $\ell_\infty$  norm, a closed-form solution is given by

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \frac{\|\nabla L(\boldsymbol{\theta}_t)\|_1}{L_\infty} \text{sign}(\nabla L(\boldsymbol{\theta}_t)),$$

which corresponds to a scaled sign-based update.

### H.2 Extended related work

**Transformer architecture and layer normalization.** The original Transformer architecture (Vaswani, 2017), referred to as Post-LN, applies layer normalization after the residual connection. In contrast, the Pre-LN architecture places layer normalization before the residual connection. Wang et al. (2019b) demonstrated that Post-LN Transformers are difficult to train when the number of layers is large, a finding later theoretically confirmed by Xiong et al. (2020) using mean field theory. Other architectures such as Reformer (He et al., 2021) were also introduced. Shi et al. (2022) showed that a large standard deviation in layer normalization leads to rank collapse in Post-LN Transformers. Furthermore, Wu et al. (2024) observed that sparse masked attention mitigates rank collapse in the absence of layer normalization and that layer normalization induces equilibria ranging from rank one to full rank.

**Attention sparsity.** Sparse attention mechanisms have been proposed to reduce the computational costs of Transformers. For example, ETC (Ainslie et al., 2020) introduces efficient sparse attention, and Zaheer et al. (2020) proposed BigBird, which they theoretically demonstrated to be as expressive as full attention. These sparse attention mechanisms are widely used in language models with large context windows, such as Longformer (Beltagy et al., 2020) and Mistral 7B (Jiang et al., 2023). In NLP, Clark (2019) found that attention of pre-trained BERT focuses on specific tokens. In vision, Hyeon-Woo et al. (2023) showed that while uniform attention is challenging to learn with the softmax function, ViT successfully learns uniform attention, which is key to its success. Additionally, Zhai et al. (2023) suggested that low attention entropy contributes to training instability in Transformers, a phenomenon they termed *entropy collapse*. Furthermore, Bao et al. (2024) demonstrated that a small eigenspectrum variance of query and key matrices leads to localized attention and mitigates both rank and entropy collapse.

## I Notation

Table S.12 shows our notations.

Table S.12: Table of notations.

| Variable  | Definition   |
|---|--|
| $a_k$   | $k$ -th element of vector $\mathbf{a}$                                       |
| $\mathbf{A}_{k,:}, \mathbf{A}_{:,j}, A_{k,j}$   | $k$ -th row, $j$ -th column, and $(k, j)$ -th element of matrix $\mathbf{A}$ |
| $[\mathbf{A}]_b, [\mathbf{a}]_b$                | $b$ -th block of matrix $\mathbf{A}$ and vector $\mathbf{a}$                 |
| $B$   | number of blocks in parameters   |
| $\mathbf{1}_a$                                  | all-ones vector of size $a$  |
| $\mathbf{I}_a$                                  | identity matrix of size $a \times a$   |
| $\text{vec}(\cdot), \text{blockdiag}(\cdot)$    | row-wise vectorization, block diagonal matrix                                |
| $\otimes$                                       | Kronecker product  |
| $C, N$  | number of classes and training samples                                       |
| $P, P_b$  | dimensions of model parameters, and $b$ -th block of parameters              |
| $\mathcal{X}$                                   | sample space   |
| $\theta$  | model parameter  |
| $\mathbf{f}(\cdot), \phi(\cdot)$                | model and feature extractor  |
| $\mathbf{V}, \mathbf{b}$                        | weight matrix and bias of the linear head                                    |
| $h, d$  | dimensions of features and tokens  |
| $\mathbf{x}^{(i)}, y^{(i)}$                     | $i$ -th training sample and label  |
| $L(\cdot)$                                      | training loss  |
| $\hat{L}(\cdot)$                                | mini-batch loss  |
| $\eta_t$  | learning rate at iteration $t$   |
| $\ell(\cdot, \cdot)$                            | cross entropy loss function  |
| $\sigma_{\text{SM}}(\cdot), \text{sign}(\cdot)$ | softmax and sign function  |
| $\mathcal{R}_{\text{FT}}$                       | parameter region of fine-tuning  |
| $L_* = L(\theta_*)$                             | local minimum of training loss   |
| $\rho_H$  | Lipschitz constant of the Hessian matrix                                     |
| $L_D$   | block-diagonal approximation of the Hessian matrix                           |
| $\delta_D$                                      | upper bound of the approximation of $L_D$                                    |
| $\sigma_2, \sigma_3$                            | constants in the upper bound of the gradient error                           |
| $\text{SA}(\cdot)$                              | single-head self-attention   |
| $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V$      | query, key, and value weight matrix  |
| $d_k, d_v$                                      | dimensions of key/query and value  |