# AdaSAM: Boosting Sharpness-Aware Minimization with Adaptive Learning Rate and Momentum for Training Deep Neural Networks

Hao Sun, Li Shen, Qihuang Zhong, Liang Ding, Shixiang Chen
Jingwei Sun, Jing Li, Guangzhong Sun, and Dacheng Tao *Fellow, IEEE*

*Abstract*—Sharpness aware minimization (SAM) optimizer has been extensively explored as it can generalize better for training deep neural networks via introducing extra perturbation steps to flatten the landscape of deep learning models. Integrating SAM with adaptive learning rate and momentum acceleration, dubbed AdaSAM, has already been explored empirically to train large-scale deep neural networks without theoretical guarantee due to the triple difficulties in analyzing the coupled perturbation step, adaptive learning rate and momentum step. In this paper, we try to analyze the convergence rate of AdaSAM in the stochastic non-convex setting. We theoretically show that AdaSAM admits a $\mathcal{O}(1/\sqrt{bT})$ convergence rate, which achieves linear speedup property with respect to mini-batch size $b$. Specifically, to decouple the stochastic gradient steps with the adaptive learning rate and perturbed gradient, we introduce the delayed second-order momentum term to decompose them to make them independent while taking an expectation during the analysis. Then we bound them by showing the adaptive learning rate has a limited range, which makes our analysis feasible. To the best of our knowledge, we are the first to provide the non-trivial convergence rate of SAM with an adaptive learning rate and momentum acceleration. At last, we conduct several experiments on several NLP tasks, which show that AdaSAM could achieve superior performance compared with SGD, AMSGrad, and SAM optimizers.

*Index Terms*—Sharpness-aware minimization, Adaptive learning rate, Non-convex optimization, linear speedup.

## I. INTRODUCTION

S HARPNESS-AWARE minimization (SAM) [1] is a powerful optimizer for training large-scale deep learning models by explicitly minimizing the gap between the training performance and generalization performance. It has achieved remarkable results in training various deep neural networks, including ResNet [1]–[3], vision transformer [4], [5], language models [6]–[8], on extensive benchmarks.

However, SAM-type methods suffer from several issues during training the deep neural networks, especially for huge computation costs and heavily hyper-parameter tuning procedure. In each iteration, SAM needs double gradients computation compared with classic optimizers, like SGD, Adam [9],

Hao Sun, Jingwei Sun, Jing Li and Guangzhong Sun are with School of Computer Science and Technology, University of Science and Technology of China, Hefei, China, 230000. (E-mail: ustcsh@mail.ustc.edu.cn, sunjw@ustc.edu.cn, lj@ustc.edu.cn, gzsun@ustc.edu.cn.)

Qihuang Zhong is with the School of Computer Science, Wuhan University, Hubei, 430000. (E-mail: zhongqihuang@whu.edu.cn)

Li Shen, Liang Ding, Shixiang Chen, and Dacheng Tao are with JD Explore Academy, Beijing, 100000. (E-mail: mathshenli@gmail.com, liangding.liam@gmail.com, chenshxiang@gmail.com, dacheng.tao@gmail.com.)

AMSGrad [10], due to the extra perturbation step. Hence, SAM requires to forward and back propagate twice for one parameter update, resulting in one more computation cost than the classic optimizers. Moreover, as there are two steps during the training process, it needs double hyper-parameters, which makes the learning rate tuning unbearable and costly.

Adaptive learning rate optimization methods [11] scale the gradients based on the history gradient information to accelerate the convergence by tuning the learning rate automatically. These methods, such as Adagrad [12], Adam [9], and AMS-Grad [10], have been proposed for solving the computer vision, natural language process, and generative neural networks tasks [11], [13]–[15]. Recently, several works have tried to ease the learning rate tuning in SAM by inheriting the triplet advantages of SAM, adaptive learning rate, and momentum acceleration. For example, [16] and [17] train ViT models and NLP models with adaptive learning rates and momentum acceleration, respectively. Although remarkable performance has been achieved, their convergences are still unknown since the adaptive learning rate and momentum acceleration are used in SAM. Directly analyzing its convergence is complicated and difficult due to the three coupled steps of optimization, i.e., the adaptive learning rate estimation is coupled with the momentum step and perturbation step of SAM.

In this paper, we analyze the convergence rate of SAM with an adaptive learning rate and momentum acceleration, dubbed AdaSAM, in the non-convex stochastic setting. To circumvent the difficulty in the analysis, we develop a novel technique to decouple the three-step training of SAM from the adaptive learning rate and momentum step. The analysis procedure is mainly divided into three parts. The first part is to analyze the procedure of the SAM. Then we analyze the second step that adopts the adaptive learning rate method. We introduce a second-order momentum term from the previous iteration, which is related to the adaptive learning rate and independent of SAM while taking an expectation. Then we can bound the term composed by the SAM and the previous second-order momentum due to the limited adaptive learning rate. In the last part, we analysis the momentum acceleration that is combined with the SAM and the adaptive learning rate. The momentum acceleration lead to an extra term in convergence analysis. Here, we introduce an auxiliary sequence to absorb it and show that their summation over the all iterations is controllable. We prove that AdaSAM enjoys the property of linear speedup property with respect to the batch size, i.e. $\mathcal{O}(1/\sqrt{bT})$ where

$b$ is the mini-batch size. Empirically, we apply AdaSAM to train RoBERTa model on the GLUE benchmark to evaluate our theoretical findings. We show that AdaSAM achieves the best performance in experiments, where it wins 6 tasks of 8 tasks, and the linear speedup can be clearly observed.

In the end, we summarize our contributions as follows:

- We present the first convergence guarantee of the adaptive SAM method with momentum acceleration under the stochastic non-convex setting. Our results suggest that a large mini-batch can help convergence due to the established linear speedup with respect to batch size.
- We conduct a series of experiments on various tasks. The results show that AdaSAM outperforms most of the state-of-art optimizers and the linear speedup is verified.

## II. PRELIMINARY AND RELATED WORK

In this section, we first describe the basic problem setup and then introduce several related works on the SAM, adaptive learning rate and momentum steps.

### A. Problem Setup

In this work, we focus on stochastic nonconvex optimization

$$\min_{x \in \mathbb{R}^d} f(x) := \mathbb{E}_{\xi \sim D} f_\xi(x), \tag{1}$$

where $d$ is dimension of variable $x$, $D$ is the unknown distribution of the data samples, $f_\xi(x)$ is a smooth and possibly non-convex function, and $f_{\xi_i}(x)$ denotes the objective function at the sampled data point $\xi_i$ according to data distribution $D$. In machine learning, it covers empirical risk minimization as a special case and $f$ is the loss function when the dataset $D$ cover $N$ data points, i.e., $D = \{\xi_i, i = 1, 2, \ldots, N\}$. Problem (1) reduces to the following finite-sum problem:

$$\min_{x \in \mathbb{R}^d} f(x) := \frac{1}{N} \sum_i f_{\xi_i}(x). \tag{2}$$

*a) Notations.:* Without additional declaration, we represent $f_i(x)$ as $f_{\xi_i}(x)$ for simplification, which is the $i$-th loss function while $x \in \mathbb{R}^d$ is the model parameter and $d$ is the parameter dimension. We denote the $l_2$ norm as $\|\cdot\|_2$. A Hadamard product is denoted as $a \odot b$ where $a,b$ are two vectors. For a vector $a \in \mathbb{R}^d$, $\sqrt{a}$ is denoted as a vector that the $j$-th value, $(\sqrt{a})_{(j)}$, is equal to the square root of $a_j$.

### B. Related Work

*a) Sharpness-aware minimization:* Many works try to improve the generalization ability during training the deep learning model. Some methods such as dropout [18], weight decay [19], and regularization methods [20], [21] provide an explicit way to improve generalization. Previous work shows that sharp minima may lead to poor generalization whereas flat minima perform better [22]–[24]. Therefore, it is popular to consider sharpness to be closely related to the generalization. Sharpness-aware minimization (SAM) [1] targets to find flat minimizers explicitly by minimizing the training loss

uniformly in the entire neighborhood. Specifically, SAM aims to solve the following minimax saddle point problem:

$$\min_x \max_{\|\delta\| \leq \rho} f(x + \delta) + \lambda \|x\|_2^2, \tag{3}$$

where $\rho \geq 0$ and $\lambda \geq 0$ are two hyperparameters. That is, the perturbed loss function of $f(x)$ in a neighborhood is minimized instead of the original loss function $f(x)$. By using Taylor expansion of $f(x + \delta)$ with respect to $\delta$, the inner max problem is approximately solved via

$$
\begin{aligned}
\delta^*(x) &= \arg\max_{\|\delta\| \leq \rho} f(x + \delta) \\
&\approx \arg\max_{\|\delta\| \leq \rho} f(x) + \delta^\top \nabla f(x) \\
&= \arg\max_{\|\delta\| \leq \rho} \delta^\top \nabla f(x) = \rho \frac{\nabla f(x)}{\|\nabla f(x)\|}.
\end{aligned}
$$

By dropping the quadratic term, (3) is simplified as the following minimization problem

$$\min_x f\left(x + \rho \frac{\nabla f(x)}{\|\nabla f(x)\|}\right). \tag{4}$$

The stochastic gradient of $f\left(x + \rho \frac{\nabla f(x)}{\|\nabla f(x)\|}\right)$ on a batch data $b$ includes the Hessian-vector product, SAM further approximates the gradient by

$$\nabla_x f_b\left(x + \rho \frac{\nabla f_b(x)}{\|\nabla f_b(x)\|}\right) \approx \nabla_x f_b(x)\big|_{x + \rho \frac{\nabla f_b(x)}{\|\nabla f_b(x)\|}}.$$

Then, along the negative direction $-\nabla_x f_b(x)\big|_{x + \rho \frac{\nabla f_b(x)}{\|\nabla f_b(x)\|}}$, SGD is applied to solve the surrogate minimization problem (4). It is easy to see that SAM requires twice gradient back-propagation, i.e., $\nabla f_b(x)$ and $\nabla_x f_b(x)\big|_{x + \rho \frac{\nabla f_b(x)}{\|\nabla f_b(x)\|}}$. Due to the existence of hyperparameter $\rho$, one needs to carefully tune both $\rho$ and learning rate in SAM. In practice, $\rho$ is predefined to control the radius of the neighborhood.

Recently, Several variants of SAM are proposed to improve its performance. For example, [8], [16], [17] have empirically incorporated adaptive learning rate with SAM and shown impressive generalization accuracy, while their convergence analysis has never been studied. ESAM [25] proposes an efficient method by sparsifying the gradients to alleviate the double computation cost of backpropagation. ASAM [17] modifies SAM by adaptively scaling the neighborhood so that the sharpness is invariant to parameters re-scaling. GSAM [16] simultaneously minimizes the perturbed function and a new defined surrogate gap function to further improve the flatness of minimizers. Liu et al. [26] also study SAM in large-batch training scenario and periodically update the perturbed gradient. Recently, [3], [8] improve the efficiency of SAM by adopting the sparse gradient perturbation technique. [27], [28] extend SAM to the federated learning setting setting with a significant performance gain. On the other hand, there are some works analyzing the convergence of the SAM such as [29] without considering the normalization step, i.e., the normalization in $\frac{\nabla f_b(x)}{\|\nabla f_b(x)\|}$.

*b) Adaptive optimizer:* The adaptive optimizer can automatically adjust the learning rate based on the history gradients methods. The first adaptive method, Adagrad [12], can achieve a better result than other first-order methods under the convex setting. While training the deep neural network, Adagrad will decrease the learning rate rapidly with a degraded performance. Adadelta [30] is proposed to change this situation and introduces a learning rate based on the exponential average history gradients. Adam [9] additionally adds momentum step to stabilize the training process, and it shows great performance in many tasks. However, Reddi et al [10] give a counterexample that it cannot converge even when the objective function is convex and propose an alternative method called AMSGrad with convergence guarantee. Then, many works [31]–[44] have been proposed to study the convergence of adaptive methods and their variants in the nonconvex setting. However, their analysis techniques can not directly extend to establish the convergence of SAM with adaptive learning rate due to the coupled perturbation step and adaptive learning rate.

*c) Momentum acceleration:* Momentum methods such as Polyak's heavy ball method [45], Nestrov's accelerated gradient descent method [46] and accelerated projected method [47] are used to optimize the parameters of the deep neural network. In practice, they have been used to accelerated for federated learning tasks [48], non-negative latent factor model [49] and recommender systems [50]. There are many theoretical works [51]–[53] that focus on analyzing the momentum acceleration for optimizing non-convex problem. [54] shows that it is important for tuning momentum while training deep neural network. [55] first points out linear convergence results for stochastic momentum method. [56] proposes a class of accelerated zeroth-order and first-order momentum method to solve mini-optimization and minimax-optimization problem. [57] extend the momentum method by introducing an RNA scheme and a constrained formulation RNA which has nonlinear updates. [58] propose a heuristic adaptive restart method and [59] propose a scheduled restart momentum accelerated SGD method named SRSGD which helps reduce the training time. [60] adds one momentum term on to the distributed gradient algorithm.

## III. METHODOLOGY

In this section, we introduce SAM with adaptive learning rate and momentum acceleration, dubbed AdaSAM, to stabilize the training process of SAM and ease the learning rate tuning. Then, we present the convergence results of AdaSAM. At last, we give the proof sketch for the main theorem.

### A. AdaSAM Algorithm

AdaSAM for solving Problem (1) is described in Algorithm 1. In each iteration, a mini-batch gradient estimation $g_t$ at point $x + \epsilon(x)$ with batchsize $b$ is computed, i.e.,

$$g_t = \nabla_x f_b(x)|_{x_t + \epsilon(x_t)} = \frac{1}{b} \sum_{i \in B} \nabla f_{\xi_i}(x_t + \delta(x_t)).$$

---

**Algorithm 1:** AdaSAM: SAM with adaptive learning rate and momentum acceleration

**Input:** Initial parameters $x_0$, $m_{-1} = 0$, $\hat{v}_{-1} = \epsilon^2$ (a small positive scalar to avoid the denominator diminishing), base learning rate $\gamma$, neighborhood size $\rho$ and momentum parameters $\beta_1$, $\beta_2$.

**Output:** Optimized parameter $x_{T+1}$

1 **for** *iteration* $t \in \{0, 1, 2, ..., T-1\}$ **do**
2 $\quad$ Sample mini-batch $B = \{\xi_{t_1}, \xi_{t_2}, ..., \xi_{t_{|B|}}\}$;
3 $\quad$ Compute gradient
$\quad\quad s_t = \nabla_x f_B(x)|_{x_t} = \frac{1}{b} \sum_{i \in B} \nabla f_{t_i}(x_t)$;
4 $\quad$ Compute $\delta(x_t) = \rho_t \frac{s_t}{\|s_t\|}$;
5 $\quad$ Compute SAM gradient $g_t = \nabla_x f_B(x)|_{x_t + \delta(x_t)}$;
6 $\quad m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$;
7 $\quad v_t = \beta_2 v_{t-1} + (1 - \beta_2)[g_t]^2$;
8 $\quad \hat{v}_t = \max(\hat{v}_{t-1}, v_t)$;
9 $\quad \eta_t = 1/\sqrt{\hat{v}_t}$;
10 $\quad x_{t+1} = x_t - \gamma m_t \odot \eta_t$;
11 **end**

---

Here, $\delta(x_t)$ is the extra perturbed gradient step in SAM that is given as follows

$$\delta(x_t) = \rho \frac{s_t}{\|s_t\|}, \text{ where } s_t = \nabla_x f_b(x)|_{x_t} = \frac{1}{b} \sum_{i \in B} \nabla f_{\xi_i}(x_t).$$

Then, the momentum term of $g_t$ and the second-order moment term $[g_t]^2$ is accumulatively computed as $m_t$ and $v_t$, respectively. AdaSAM then updates iterate along $-m_t$ with the adaptive learning rate $\gamma \eta_t$.

**Remark 1.** *Below, we give several comments on AdaSAM:*
- *When $\beta_2 = 1$, the adaptive learning rate reduces to the diminishing one as SGD. Then, AdaSAM recovers the classic SAM optimizer.*
- *If we drop out the 8-th line $\hat{v}_t = \max(\hat{v}_{t-1}, v_t)$, then our algorithm becomes the variant of Adam. The counterexample that Adam does not converge in the [10] also holds for the SAM variant, while AdaSAM can converge.*

### B. Convergence Analysis

Before presenting the convergence results of the AdaSAM algorithm, we first introduce some necessary assumptions.

**Assumption 1** (*L*-smooth). *$f_i$ and $f$ is differentiable with gradient Lipschitz property: $\|\nabla f_i(x) - \nabla f_i(y)\| \le L\|x - y\|, \|\nabla f(x) - \nabla f(y)\| \le L\|x - y\|, \forall x, y \in \mathbb{R}^d, i = 1, 2, ..., N$, which also implies the descent inequality, i.e., $f_i(y) \le f_i(x) + \langle \nabla f_i(x), y - x \rangle + \frac{L}{2}\|y - x\|^2$.*

**Assumption 2** (**Bounded variance**). *The estimator of the gradient is unbiased and the variance of the stochastic gradient is bounded. i.e.,*

$$\mathbb{E}\nabla f_i(x) = \nabla f(x), \quad \mathbb{E}\|\nabla f_i(x) - \nabla f(x)\|^2 \le \sigma^2.$$

*When the mini-batch size $b$ is used, we have $\mathbb{E}\|\nabla f_b(x) - \nabla f(x)\|^2 \le \frac{\sigma^2}{b}$.*

**Assumption 3** (**Bounded stochastic gradients**). *The stochastic gradient is uniformly bounded,* i.e.,

$$\|\nabla f_i(x)\|_\infty \le G, for\ any\ i = 1, \dots, N.$$

**Remark 2.** *The above assumptions are commonly used in the proof of convergence for adaptive stochastic gradient methods such as [31], [32], [61], [62].*

Below, we briefly explain the main idea of analyzing the convergence of the AdaSAM algorithm. First, we discuss the difficulty of applying the adaptive learning rate on SAM. We notice that the main step which contains adaptive learning rate in convergence analysis is to estimate the expectation $\mathbb{E}[x_{t+1} - x_t] = -\mathbb{E}m_t \odot \eta_t = -\mathbb{E}(1-\beta_1)g_t \odot \eta_t - \mathbb{E}\beta_1 m_{t-1} \odot \eta_t$, which is conditioned on the filtration $\sigma(x_t)$. In this part, we consider the situation that $\beta_1 = 0$ which does not include the momentum. Then, we apply delay technology to disentangle the dependence between $g_t$ and $\eta_t$, that is

$$\mathbb{E}g_t \odot \eta_t = \mathbb{E}[g_t \odot \eta_{t-1}] + \mathbb{E}[g_t \odot (\eta_t - \eta_{t-1})]$$
$$= \nabla f(x_t) \odot \eta_{t-1} + \mathbb{E}[g_t \odot (\eta_t - \eta_{t-1})].$$

The second term $\mathbb{E}[g_t \odot (\eta_t - \eta_{t-1})]$ is dominated by the first term $\nabla f(x_t) \odot \eta_{t-1}$. Then, it is not difficult to get the convergence result of the stochastic gradient descend with the adaptive learning rate such as AMSGrad. However, when we apply the same strategy to AdaSAM, we find that $\mathbb{E}g_t \odot \eta_{t-1}$ cannot be handled similarly because $\mathbb{E}g_t = \mathbb{E}\nabla_x f_b\left(x + \rho \frac{\nabla f_b(x)}{\|\nabla f_b(x)\|}\right) \ne \nabla f(x_t)$. Inspired by [29, Lemma 16], our key observation is that

$$\mathbb{E}\nabla_x f_b\left(x + \rho \frac{\nabla f_b(x)}{\|\nabla f_b(x)\|}\right) \approx \mathbb{E}\nabla_x f_b\left(x + \rho \frac{\nabla f(x)}{\|\nabla f(x)\|}\right)$$
$$= \nabla_x f\left(x + \rho \frac{\nabla f(x)}{\|\nabla f(x)\|}\right)$$

and we prove the other terms such as $\mathbb{E}\left(\nabla_x f_b\left(x + \rho \frac{\nabla f_b(x)}{\|\nabla f_b(x)\|}\right) - \nabla_x f_b\left(x + \rho \frac{\nabla f(x)}{\|\nabla f(x)\|}\right)\right) \odot \eta_{t-1}$ have small values that do not dominate the convergence rate.

On the other hand, when we apply the momentum steps, we find that the term $\mathbb{E}m_{t-1} \odot \eta_t$ cannot be ignored. By introducing an auxiliary sequence $z_t = x_t + \frac{\beta_1}{1-\beta_1}(x_t - x_{t-1})$, we have $\mathbb{E}[z_{t+1} - z_t] = -\mathbb{E}[\frac{\beta_1}{1-\beta_1}\gamma m_{t-1} \odot (\eta_{t-1} - \eta_t) - \gamma g_t \odot \eta_t]$. The first term contains the momentum term which has a small value due to the difference of the adaptive learning rate $\eta_t$. Thus, it is diminishing without hurting the convergence rate.

**Theorem 1.** *Under the assumptions 1,2,3, and $\gamma$ is a fixed number satisfying $\gamma \le \frac{\epsilon}{16L}$, for the sequence $\{x_t\}$ generated by Algorithm 1, we have the following convergence rate*

$$\frac{1}{T}\sum_{t=0}^{T-1}\mathbb{E}\|\nabla f(x_t)\|_2^2 \le \frac{2G(f(x_0)-f^*)}{\gamma T} + \frac{8G\gamma L}{\epsilon}\frac{\sigma^2}{b\epsilon} + \Phi \quad (5)$$

*where*

$$\Phi = \frac{45GL^2\rho_t^2}{\epsilon} + \frac{2G^3}{(1-\beta_1)T}d(\frac{1}{\epsilon} - \frac{1}{G}) + \frac{6\gamma^2 L^2\beta_1^2}{(1-\beta_1)^2}\frac{dG^3}{\epsilon^3}$$
$$+ \frac{2(4 + (\frac{\beta_1}{1-\beta_1})^2)\gamma LG^3}{T}d(\epsilon^{-2} - G^{-2}) + \frac{8G\gamma L}{\epsilon}\frac{L\rho_t^2}{\epsilon}, \quad (6)$$

*in which $T$ is the number of iteration, $f^*$ is the minimal value of the function $f$, $\gamma$ is the base learning rate, $b$ is the minibatch size, $d$ is the dimension of paramter $x$. $\beta_1$, $G$, $L$, $\epsilon$, $\sigma^2$, $d$ are fixed constants.*

Theorem 1 characterizes the convergence rate of the sequence $\{x_t\}$ generated by AdaSAM with respect to the stochastic gradient residual. The first two terms of the right hand side of Inequality (5) are the terms that dominate the convergence rate. Compared with the first two terms, $\Phi$ is a small value while we set neighborhood size $\rho$ and learning rate $\gamma$ as small values which are related to large iteration number $T$. Then, we obtain the following corollary directly.

**Corollary 1** (**Mini-batch linear speedup**). *Under the same conditions of Theorem 1. Furthermore, when we choose the base learning rate $\gamma = O(\sqrt{\frac{b}{T}})$ and neighborhood size $\rho = O(\sqrt{\frac{1}{bT}})$ , the following result holds:*

$$\frac{1}{T}\sum_{t=0}^{T-1}\mathbb{E}\|\nabla f(x_t)\|_2^2 = O\left(\frac{1}{\sqrt{bT}}\right) + O\left(\frac{1}{bT}\right) + O\left(\frac{1}{T}\right)$$
$$+ O\left(\frac{1}{b^{\frac{1}{2}}T^{\frac{3}{2}}}\right) + O\left(\frac{b^{\frac{1}{2}}}{T^{\frac{3}{2}}}\right) + O\left(\frac{b}{T}\right).$$

*When $T$ is sufficiently large, we achieve the linear speedup convergence rate with respect to mini-batch size $b$, i.e.,*

$$\frac{1}{T}\sum_{t=0}^{T-1}\mathbb{E}\|\nabla f(x_t)\|_2^2 = O\left(\frac{1}{\sqrt{bT}}\right). \quad (7)$$

**Remark 3.** *Two comments are given about the above results:*

- *To reach a $O(\delta)$ stationary point, when the batch size is 1, it needs $T = O(\frac{1}{\delta^2})$ iterations. When the batch size is $b$, we need to run $T = O(\frac{1}{b\delta^2})$ steps. The method with batch size $b$ is $b$ times faster than batch size of 1, which means that it has the mini-batch linear speedup property.*
- *According to [37], [63], [64], AdaSAM can be extended to distributed version and achieves linear speedup property with respect to the number of works in the Parameter-Sever setting.*

*C. Proof Sketch*

In this part, we give the proof sketch of the Theorem 1. For the complete proof, please see Appendix. Below, we first introduce an auxiliary sequence $z_t = x_t + \frac{\beta_1}{1-\beta_1}(x_t - x_{t-1})$. By applying $L$-smooth condition, we have

$$f(z_{t+1}) \le f(z_t) + \langle \nabla f(z_t), z_{t+1} - z_t \rangle + \frac{L}{2}\|z_{t+1} - z_t\|^2. \quad (8)$$

Applying it to the sequence $\{z_t\}$ and using the delay strategy yield

$$f(z_{t+1}) - f(z_t)$$
$$\le \langle \nabla f(z_t), \frac{\gamma\beta_1}{1-\beta_1}m_{t-1} \odot (\eta_{t-1} - \eta_t) \rangle + \frac{L}{2}\|z_{t+1} - z_t\|^2$$
$$+ \langle \nabla f(z_t), \frac{\gamma}{b}\sum_{i\in B}\nabla f_i(x_t + \rho_t \frac{s_t}{\|s_t\|}) \odot (\eta_{t-1} - \eta_t) \rangle$$

TABLE I: Evaluating SGD, SAM, AMSGrad and AdaSAM on the GLUE benchmark with $\beta_1 = 0.9$

| Model | CoLA mcc. | SST-2 Acc. | MRPC Acc./F1 | STS-B Pcor./Scor. | RTE Acc. | MNLI m./mm. | QNLI Acc. | QQP F1/ Acc. | Avg. |
|---|---|---|---|---|---|---|---|---|---|
| SGD | 9.25 | 50.92 | 68.38/ 81.22 | 3.22/ 1.9 | 55.6 | 84.94/ 84.87 | 63.61 | 85.6/ 80.14 | 55.8 |
| SAM($\rho$ =0.01) | 4.64 | 95.87 | 70.58/ 81.98 | 84.74/ 85.57 | 52.71 | 90.5/ 90.19 | 94.44 | 84.7/ 87.88 | 76.98 |
| SAM($\rho$ =0.005) | 66.76 | 95.76 | 68.38/ 81.22 | 2/ 2 | 52.71 | 90.42/ 89.74 | 94.6 | 86.72/ 89.94 | 68.35 |
| SAM(best) | 66.76 | 95.87 | 70.58/ 81.98 | 84.74/ 85.57 | 52.71 | 90.5/ 90.19 | 94.6 | 86.72/ 89.94 | 82.51 |
| AMSGrad | 68.0 | 96.33 | 90.2/ 92.72 | 91.72/ 91.48 | 87.73 | 90.67/ 90.41 | 94.82 | 88.7/ 91.41 | 89.52 |
| AdaSAM($\rho$ =0.01) | 65.29 | 96.33 | 91.18/ 93.64 | 90.13/ 90.36 | 84.84 | 90.97/ 90.42 | 94.65 | 88.55/ 91.23 | 88.97 |
| AdaSAM($\rho$ =0.005) | 68.74 | 96.67 | 90.93/ 93.36 | 91.64/ 91.38 | 87.73 | 90.88/ 90.4 | 94.56 | 88.69/ 91.33 | 89.69 |
| AdaSAM($\rho$ =0.001) | 67.3 | 96.1 | 90.2/ 92.96 | 91.9/ 91.62 | 85.92 | 90.45/ 90.4 | 94.56 | 88.64/ 91.27 | 89.28 |
| AdaSAM(best) | 68.74 | 96.67 | 91.18/ 93.64 | 91.9/ 91.62 | 87.73 | 90.97/ 90.42 | 94.65 | 88.69/ 91.33 | 89.8 |

$$+ \langle \nabla f(z_t) - \nabla f(x_t), -\frac{\gamma}{b} \sum_{i \in B} \nabla f_i(x_t + \rho_t \frac{s_t}{\|s_t\|}) \odot \eta_{t-1} \rangle$$

$$+ \langle \nabla f(x_t), -\frac{\gamma}{b} \sum_{i \in B} \nabla f_i(x_t + \rho_t \frac{\nabla f(x_t)}{\|\nabla f(x_t)\|}) \odot \eta_{t-1} \rangle$$

$$+ \langle \nabla f(x_t), \frac{\gamma}{b} \sum_{i \in B} \nabla f_i(x_t + \rho_t \frac{\nabla f(x_t)}{\|\nabla f(x_t)\|}) \odot \eta_{t-1}$$

$$- \frac{\gamma}{b} \sum_{i \in B} \nabla f_i(x_t + \rho_t \frac{s_t}{\|s_t\|}) \odot \eta_{t-1} \rangle. \tag{9}$$

From the Lemma 5, Lemma 6, Lemma 7 in appendix, we can bound the above terms in (9) as follows

$$\langle \nabla f(z_t), \frac{\gamma}{b} \sum_{i \in B} \nabla f_i(x_t + \rho_t \frac{s_t}{\|s_t\|}) \odot (\eta_{t-1} - \eta_t) \rangle$$

$$\leq \gamma G^2 \|\eta_{t-1} - \eta_t\|_1, \tag{10}$$

$$\langle \nabla f(z_t), \frac{\gamma \beta_1}{1 - \beta_1} m_{t-1} \odot (\eta_{t-1} - \eta_t) \rangle$$

$$\leq \frac{\gamma \beta_1}{1 - \beta_1} G^2 \|\eta_{t-1} - \eta_t\|_1, \tag{11}$$

$$\langle \nabla f(x_t), \frac{\gamma}{b} \sum_{i \in B} \nabla f_i(x_t + \rho_t \frac{\nabla f(x_t)}{\|\nabla f(x_t)\|}) \odot \eta_{t-1}$$

$$- \frac{\gamma}{b} \sum_{i \in B} \nabla f_i(x_t + \rho_t \frac{s_t}{\|s_t\|}) \odot \eta_{t-1} \rangle$$

$$\leq \frac{\gamma}{2\mu^2} \|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2 + \frac{2\mu^2 \gamma L^2 \rho_t^2}{\epsilon}. \tag{12}$$

Then we substitute them into the (9), and take the conditional expectation to get

$$\mathbb{E} f(z_{t+1}) - f(z_t)$$

$$\leq \mathbb{E} \langle \nabla f(x_t), -\frac{\gamma}{b} \sum_{i \in B} \nabla f_i(x_t + \rho_t \frac{\nabla f(x_t)}{\|\nabla f(x_t)\|}) \odot \eta_{t-1} \rangle$$

$$+ \frac{\gamma}{2\mu^2} \|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2 + \frac{\gamma}{1 - \beta_1} G^2 \|\eta_{t-1} - \eta_t\|_1$$

$$+ \mathbb{E} \langle \nabla f(z_t) - \nabla f(x_t), -\frac{\gamma}{b} \sum_{i \in B} \nabla f_i(x_t + \rho_t \frac{s_t}{\|s_t\|}) \odot \eta_{t-1} \rangle$$

$$+ \frac{2\mu^2 \gamma L^2 \rho_t^2}{\epsilon} + \frac{L}{2} \mathbb{E} \|z_{t+1} - z_t\|^2, \tag{13}$$

where $\mu > 0$ is a constant to be determined. Next, from the Lemma 8, Lemma 10 and Lemma 9 in Appendix, we have

$$\mathbb{E} \langle \nabla f(x_t), -\frac{\gamma}{b} \sum_{i \in B} \nabla f_i(x_t + \rho_t \frac{\nabla f(x_t)}{\|\nabla f(x_t)\|}) \odot \eta_{t-1} \rangle$$

$$\leq -\gamma \|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2 + \mathbb{E} \frac{\gamma}{2\alpha^2} \|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2$$

$$+ \frac{\gamma \alpha^2 L^2 \rho_t^2}{2\epsilon}, \tag{14}$$

$$\frac{L}{2} \mathbb{E} \|z_{t+1} - z_t\|^2 \leq \frac{LG^2 \gamma^2 \beta_1^2}{(1 - \beta_1)^2} \mathbb{E} \|\eta_t - \eta_{t-1}\|^2$$

$$+ \gamma^2 L (3 \frac{1 + \beta}{\beta \epsilon} (\frac{L\rho_t^2}{\epsilon} + \frac{\sigma^2}{b\epsilon} + \mathbb{E} \|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2)$$

$$+ (1 + \beta) G^2 \mathbb{E} \|\eta_t - \eta_{t-1}\|^2), \tag{15}$$

$$\mathbb{E} \langle \nabla f(z_t) - \nabla f(x_t), -\frac{\gamma}{b} \sum_{i \in B} \nabla f_i(x_t + \rho_t \frac{s_t}{\|s_t\|}) \odot \eta_{t-1} \rangle$$

$$\leq \frac{\gamma^3 L^2 \beta_1^2}{2\epsilon(1 - \beta_1)^2} (\frac{1}{\lambda_1^2} + \frac{1}{\lambda_2^2} + \frac{1}{\lambda_3^2}) \frac{dG_\infty^2}{\epsilon^2} + \frac{\gamma L^2 \rho_t^2}{2\epsilon} (\lambda_2^2 + 4\lambda_3^2)$$

$$+ \frac{\gamma \lambda_1^2}{2} \|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2. \tag{16}$$

Next, we substitute it into the (13). Taking the expectation over all history information yields

$$\mathbb{E} f(x_{t+1}) - \mathbb{E} f(x_t)$$

$$\leq -\gamma (1 - \frac{1}{2\mu^2} - \frac{1}{2\alpha^2} - \frac{3\gamma L(1+\beta)}{\beta \epsilon} - \frac{\lambda_1^2}{2}) \mathbb{E} \|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2$$

$$+ \frac{2\mu^2 \gamma L^2 \rho_t^2}{\epsilon} + \frac{\gamma}{1 - \beta_1} G^2 \mathbb{E} \|\eta_{t-1} - \eta_t\|_1 + \frac{\gamma \alpha^2 L^2 \rho^2}{2\epsilon}$$

$$+ \frac{\gamma^3 L^2 \beta_1^2}{2\epsilon(1 - \beta_1)^2} (\frac{1}{\lambda_1^2} + \frac{1}{\lambda_2^2} + \frac{1}{\lambda_3^2}) \frac{dG_\infty^2}{\epsilon^2} + \frac{\gamma L^2 \rho_t^2}{2\epsilon} (\lambda_2^2 + 4\lambda_3^2)$$

$$+ \gamma^2 LG^2 ((\frac{\beta_1}{1 - \beta_1})^2 + 1 + \beta) \mathbb{E} \|\eta_t - \eta_{t-1}\|^2$$

$$+ \frac{3\gamma^2 L(1+\beta)}{\beta \epsilon} (\frac{L\rho_t^2}{\epsilon} + \frac{\sigma^2}{b\epsilon}). \tag{17}$$

We set $\mu^2 = \alpha^2 = 8$, $\beta = 3$, $\lambda_1^2 = \frac{1}{4}$, $\lambda_2^2 = \lambda_3^2 = 1$ and we choose $\frac{2\gamma L}{\epsilon} \leq \frac{1}{8}$. Note that $\eta_t$ is bounded. We have

$$\frac{\gamma}{2G} \mathbb{E} \|\nabla f(x_t)\|^2 \leq \frac{\gamma}{2} \mathbb{E} \|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2 \tag{18}$$

$$\leq -\mathbb{E} f(x_{t+1}) + \mathbb{E} f(x_t) + \frac{45\gamma L^2 \rho_t^2}{2\epsilon} + \frac{4\gamma^2 L}{\epsilon} (\frac{L\rho_t^2}{\epsilon} + \frac{\sigma^2}{b\epsilon})$$

$$+ \frac{\gamma}{1 - \beta_1} G^2 \mathbb{E} \|\eta_{t-1} - \eta_t\|_1 + \frac{3\gamma^3 L^2 \beta_1^2}{(1 - \beta_1)^2} \frac{dG_\infty^2}{\epsilon^3}$$

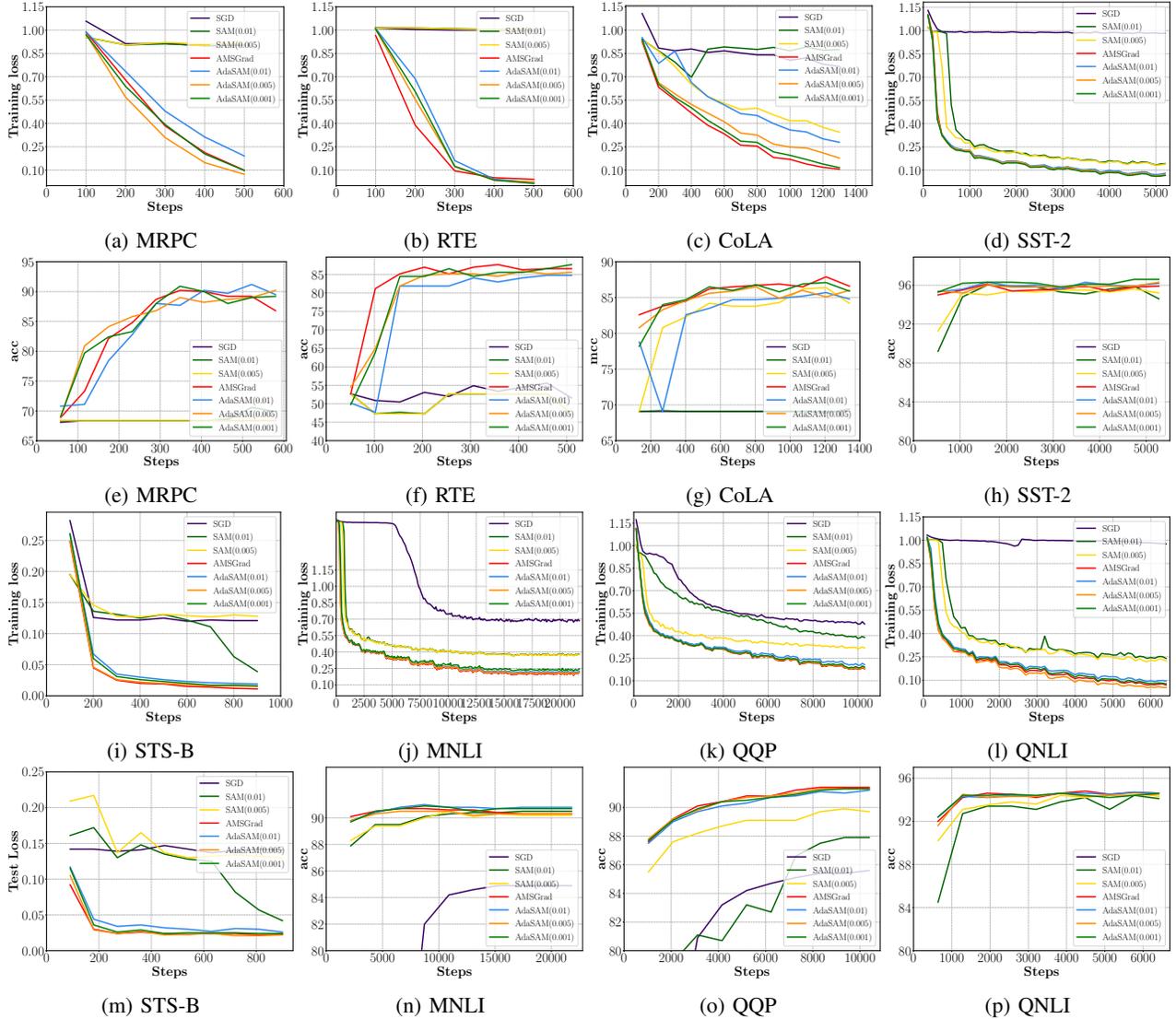$$+ (4 + (\frac{\beta_1}{1 - \beta_1})^2) \gamma^2 LG^2 \mathbb{E} \|\eta_t - \eta_{t-1}\|^2. \tag{19}$$

Fig. 1: The loss and evaluation metric v.s. steps on MRPC, RTE, CoLA, SST-2, STS-B, MNLI, QQP, and QNLI.($\beta_1 = 0.9$)

Then, telescoping it from $t = 0$ to $t = T - 1$, and assuming $\gamma$ is a constant, it follows that

$$
\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(x_t)\|^2 \leq \frac{2G(f(x_0) - f^*)}{\gamma T} + \frac{8G\gamma L}{\epsilon} \frac{\sigma^2}{b\epsilon}
$$
$$
+ \frac{45GL^2\rho_t^2}{\epsilon} + \frac{2G^3}{(1-\beta_1)T} d(\frac{1}{\epsilon} - \frac{1}{G}) + \frac{6\gamma^2 L^2 \beta_1^2}{(1-\beta_1)^2} \frac{dG^3}{\epsilon^3}
$$
$$
+ \frac{8G\gamma L}{\epsilon} \frac{L\rho_t^2}{\epsilon} + \frac{2(4 + (\frac{\beta_1}{1-\beta_1})^2)\gamma LG^3}{T} d(\epsilon^{-2} - G^{-2}), \quad (20)
$$

which completes the proof.

## IV. EXPERIMENTS

In this section, we apply AdaSAM to train language models and compare it with SGD, AMSGrad, and SAM to show its effectiveness. Due to space limitations, more experiments, including visualization, task description, implementation details and results description, are placed in the Appendix.

### A. Experimental Setup

**Tasks and Datasets.** We evaluate AdaSAM on a popular benchmark, *i.e.* General Language Understanding Evaluation (GLUE) [65], which consists of several language understanding tasks including sentiment analysis, question answering and textual entailment. For a fair comparison, we report the results based on single-task, without multi-task or ensemble training. We evaluate the performance with Accuracy ("*Acc*") metric for most tasks, except the F1 scores for QQP and MRPC, the Pearson-Spearman correlations ("*Pcor/Scor*") for STS-B and the Matthew correlations ("*Mcc*") for CoLA. The performance is better as the metric is higher.

**Implementations.** We conduct our experiments using a widely-used pre-trained language model, RoBERTa-large[1] in the open-source toolkit fairseq[2], with 24 transformer layers, a hidden size of 1024. For fine-tuning on each task, we use different combinations of hyper-parameters, including the

[1] https://dl.fbaipublicfiles.com/fairseq/models/roberta.large.tar.gz
[2] https://github.com/facebookresearch/fairseq
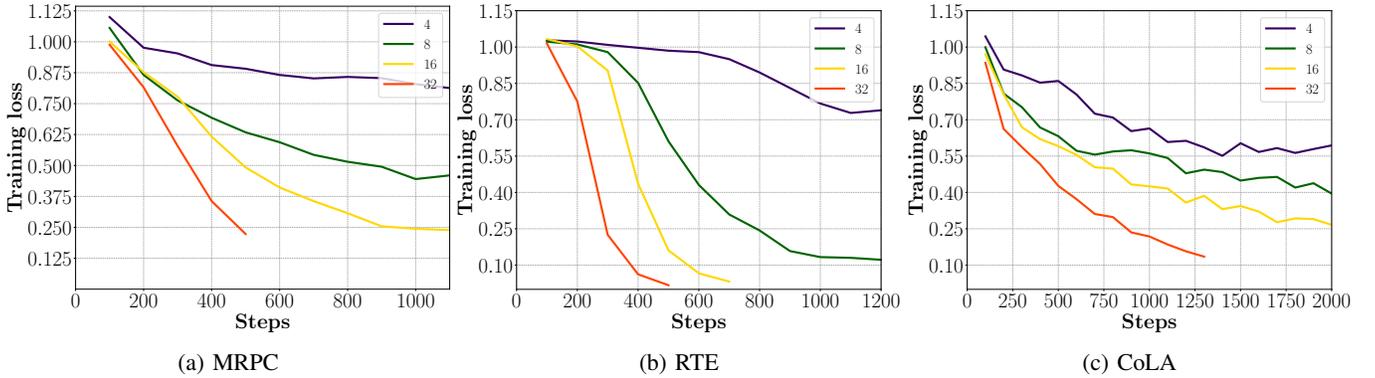
(a) MRPC    (b) RTE    (c) CoLA

Fig. 2: The linear speedup verification of AdaSAM with the number of batch size of 4, 8, 16, 32.

TABLE II: Results of SGD, SAM, AMSGrad and AdaSAM on the GLUE benchmark without momentum, i.e., $\beta_1 = 0$

| Model | CoLA mcc. | SST-2 Acc. | MRPC Acc./F1 | STS-B Pcor./Scor. | RTE Acc. | MNLI m./mm. | QNLI Acc. | QQP F1/ Acc. | Avg. |
|---|---|---|---|---|---|---|---|---|---|
| SGD | 0 | 51.722 | 68.38/ 81.22 | 5.55/ 7.2 | 51.27 | 32.51/ 32.42 | 53.32 | 0/ 63.18 | 37.23 |
| SAM($\rho$ =0.01) | 41.91 | 95.3 | 68.38/ 81.22 | 9.21/ 10.38 | 53.07 | 87.99/ 87.8 | 51.24 | 83.44/ 87.27 | 63.1 |
| SAM($\rho$ =0.005) | 58.79 | 81.54 | 68.38/ 81.22 | 13.52/ 16.6 | 53.79 | 88.42/ 88.15 | 92.95 | 83.84/ 87.7 | 67.91 |
| SAM(best) | 58.79 | 95.3 | 68.38/ 81.22 | 13.52/ 16.6 | 53.79 | 88.42/ 88.15 | 92.95 | 83.84/ 87.7 | 69.06 |
| AMSGrad | 63.78 | 96.44 | 89.71/ 92.44 | 89.98/ 90.35 | 87.36 | 90.65/ 90.35 | 94.53 | 88.59/ 91.27 | 88.79 |
| AdaSAM($\rho$ =0.01) | 69.23 | 96.22 | 89.96/ 92.84 | 88.83/ 89.07 | 87 | 90.83/ 90.41 | 94.8 | 88.67/ 91.38 | 89.1 |
| AdaSAM($\rho$ =0.005) | 68.47 | 96.22 | 89.96/ 92.82 | 91.59/ 91.22 | 73.65 | 90.75/ 90.42 | 94.73 | 88.72/ 91.46 | 88.33 |
| AdaSAM(best) | 69.23 | 96.22 | 89.96/ 92.84 | 91.59/ 91.22 | 87 | 90.83/ 90.42 | 94.8 | 88.72/ 91.46 | 89.52 |

learning rate, the number of epochs, the batch size, *etc* [3]. In particular, for RTE, STS-B and MRPC of GLUE benchmark, we first fine-tune the pre-trained RoBERTa-large model on the MNLI dataset and continue fine-tuning the RoBERTa-large-MNLI model on the corresponding single-task corpus for better performance, as many prior works did [7], [66]. All models are trained on NVIDIA DGX SuperPOD cluster, in which each machine contains 8×40GB A100 GPUs.

### B. Results on GLUE Benchmark

Table I shows the performance of SGD, SAM, AMSGrad, and AdaSAM. For the AdaSAM, we tune the neighborhood size of the perturbation parameter from 0.01, 0.005, and 0.001. The result shows that AdaSAM outperforms AMSGrad on 6 tasks of 8 tasks except for QNLI and QQP. Overall, it improves the 0.28 average score than AMSGrad. On the other hand, Table I indicates that SAM is better than SGD on 7 tasks of 8 tasks except for RTE. And SAM can significantly improve performance. Comparing the results of Table I, we can find that the adaptive learning rate method is better than SGD tuned with handicraft learning rate. AdaSAM achieves the best metric on 6 tasks which is CoLA, SST-2, MRPC, STS-B, RTE, QNLI, and MNLI. In general, AdaSAM is better than the other methods.

In addition, Figure 3 shows the convergence speed of the detailed loss and evaluation metrics vs. the number of steps during training, respectively. The loss curve of AdaSAM decreases faster than SAM and SGD in all tasks, and it has a

similar decreasing speed as the AMSGrad. The evaluation metric curve of AdaSAM and AMSGrad show that the AdaSAM is better than SGD and SAM and decreases the loss value as faster as the AMSGrad in all tasks.

### C. Mini-batch Speedup

In this part, we test the performance with different batch sizes to validate the linear speedup property. The experiments are conducted on the MRPC, RTE, and CoLA tasks. The batch size is set as 4, 8, 16, 32, respectively. We scale the learning rate as $\sqrt{N}$, which is similar as [67], where $N$ is the batch size. The results show that the training loss decreases faster as the batchsize increases, and the loss curve with the batch size of 32 achieves nearly half iterations as the curve with batch size of 16.

### D. Ablation Study

In this subsection, we conduct the experiments the momentum hyper-parameter $\beta_1$ is set to 0 to evaluate the influence of the momentum acceleration and the adaptive learning rate. Table II shows that AdaSAM outperforms AMSGrad on 6 tasks of 8 tasks except for SST-2 and RTE. In Table II, we also compare SGD and SAM, and without the momentum, SAM outperforms SGD on all tasks. Under this situation, AdaSAM without the momentum acceleration method is better than the other methods.

When comparing the result of Table I and Table II, we find that both the adaptive learning rate method and the momentum acceleration are helpful for the model's generalization ability. When there is no momentum term, SAM with an adaptive

---

[3]Due to the space limitation, we show the details of the dataset and training setting in Appendix A.

learning rate improves the 0.74 average score to AMSGrad. With a momentum term, AdaSAM improves the 0.28 average score to AMSGrad. It shows that the adaptive method can improve the performance with or without momentum acceleration and it achieves the best performance with momentum acceleration. And we can find that momentum acceleration improves the performance of SAM, AMSGrad and AdaSAM.

## V. CONCLUSION

In this work, we study the convergence rate of Sharpness aware minimization optimizer with an adaptive learning rate and momentum acceleration, dubbed AdaSAM in the stochastic non-convex setting. To the best of our knowledge, we are the first to provide the non-trivial $\mathcal{O}(1/\sqrt{bT})$ convergence rate of AdaSAM, which achieves a linear speedup property with respect to mini-batch size $b$. We have conducted extensive experiments on several NLP tasks, which verifies that AdaSAM could achieve superior performance compared with AMSGrad and SAM optimizers. Future works include extending AdaSAM to the distributed setting and reducing the twice gradient back-propagation cost.

## REFERENCES

[1] P. Foret, A. Kleiner, H. Mobahi, and B. Neyshabur, "Sharpness-aware minimization for efficiently improving generalization," in *International Conference on Learning Representations*, 2021. 1, 2

[2] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778. 1

[3] P. Mi, L. Shen, T. Ren, Y. Zhou, X. Sun, R. Ji, and D. Tao, "Make sharpness-aware minimization stronger: A sparsified perturbation approach," in *Advances in Neural Information Processing Systems*. 1, 2

[4] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale," in *International Conference on Learning Representations*, 2021. 1

[5] X. Chen, C.-J. Hsieh, and B. Gong, "When vision transformers outperform resnets without pre-training or strong data augmentations," in *International Conference on Learning Representation*, 2022. 1

[6] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018. 1

[7] P. He, X. Liu, J. Gao, and W. Chen, "Deberta: Decoding-enhanced bert with disentangled attention," in *ICLR*, 2020. 1, 7

[8] Q. Zhong, L. Ding, L. Shen, P. Mi, J. Liu, B. Du, and D. Tao, "Improving sharpness-aware minimization with fisher mask for better generalization on language models," *arXiv preprint arXiv:2210.05497*, 2022. 1, 2

[9] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *ICLR (Poster)*, 2015. [Online]. Available: http://arxiv.org/abs/1412.6980 1, 3

[10] S. J. Reddi, S. Kale, and S. Kumar, "On the convergence of adam and beyond," in *International Conference on Learning Representations*, 2018. 1, 3

[11] H. Iiduka, "Appropriate learning rates of adaptive learning rate optimization algorithms for training deep neural networks," *IEEE Trans. Cybern.*, vol. 52, no. 12, pp. 13 250–13 261, 2022. 1

[12] J. Duchi, E. Hazan, and Y. Singer, "Adaptive subgradient methods for online learning and stochastic optimization." *Journal of machine learning research*, vol. 12, no. 7, 2011. 1, 3

[13] S. Ruder, "An overview of gradient descent optimization algorithms," *arXiv preprint arXiv:1609.04747*, 2016. 1

[14] L. Liao, L. Shen, J. Duan, M. Kolar, and D. Tao, "Local adagrad-type algorithm for stochastic convex-concave optimization," *Machine Learning*, pp. 1–20, 2022. 1

[15] J. Zhang, S. P. Karimireddy, A. Veit, S. Kim, S. Reddi, S. Kumar, and S. Sra, "Why are adaptive methods good for attention models?" *Advances in Neural Information Processing Systems*, vol. 33, pp. 15 383–15 393, 2020. 1

[16] J. Zhuang, B. Gong, L. Yuan, Y. Cui, H. Adam, N. C. Dvornek, sekhar tatikonda, J. s Duncan, and T. Liu, "Surrogate gap minimization improves sharpness-aware training," in *International Conference on Learning Representations*, 2022. 1, 2

[17] J. Kwon, J. Kim, H. Park, and I. K. Choi, "Asam: Adaptive sharpness-aware minimization for scale-invariant learning of deep neural networks," in *International Conference on Machine Learning*. PMLR, 2021, pp. 5905–5914. 1, 2

[18] N. Srivastava, G. E. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: a simple way to prevent neural networks from overfitting," *J. Mach. Learn. Res.*, vol. 15, no. 1, pp. 1929–1958, 2014. 2

[19] I. Loshchilov and F. Hutter, "Decoupled weight decay regularization," *arXiv preprint arXiv:1711.05101*, 2017. 2

[20] Z. Li, H. Zhao, Y. Guo, Z. Yang, and S. Xie, "Accelerated log-regularized convolutional transform learning and its convergence guarantee," *IEEE Trans. Cybern.*, vol. 52, no. 10, pp. 10 785–10 799, 2022. 2

[21] Y. Lu, Z. Zhang, G. Lu, Y. Zhou, J. Li, and D. Zhang, "Addi-reg: A better generalization-optimization tradeoff regularization method for convolutional neural networks," *IEEE Trans. Cybern.*, vol. 52, no. 10, pp. 10 827–10 842, 2022. 2

[22] Y. Jiang, B. Neyshabur, H. Mobahi, D. Krishnan, and S. Bengio, "Fantastic generalization measures and where to find them," in *International Conference on Learning Representations*, 2020. 2

[23] N. S. Keskar, D. Mudigere, J. Nocedal, M. Smelyanskiy, and P. T. P. Tang, "On large-batch training for deep learning: Generalization gap and sharp minima," in *5th International Conference on Learning Representations, ICLR 2017*. OpenReview.net, 2017. 2

[24] H. He, G. Huang, and Y. Yuan, "Asymmetric valleys: Beyond sharp and flat local minima," *Advances in neural information processing systems*, vol. 32, 2019. 2

[25] J. Du, H. Yan, J. Feng, J. T. Zhou, L. Zhen, R. S. M. Goh, and V. Tan, "Efficient sharpness-aware minimization for improved training of neural networks," in *International Conference on Learning Representations*, 2022. 2

[26] Y. Liu, S. Mai, X. Chen, C.-J. Hsieh, and Y. You, "Towards efficient and scalable sharpness-aware minimization," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2022, pp. 12 360–12 370. 2

[27] Z. Qu, X. Li, R. Duan, Y. Liu, B. Tang, and Z. Lu, "Generalized federated learning via sharpness aware minimization," in *International Conference on Machine Learning*. PMLR, 2022, pp. 18 250–18 280. 2

[28] Y. Sun, L. Shen, T. Huang, L. Ding, and D. Tao, "Fedspeed: Larger local interval, less communication round, and higher generalization accuracy," in *International Conference on Learning Representations*. 2

[29] M. Andriushchenko and N. Flammarion, "Towards understanding sharpness-aware minimization," in *International Conference on Machine Learning*. PMLR, 2022, pp. 639–668. 2, 4

[30] M. D. Zeiler, "Adadelta: an adaptive learning rate method," *arXiv preprint arXiv:1212.5701*, 2012. 3

[31] D. Zhou, J. Chen, Y. Cao, Y. Tang, Z. Yang, and Q. Gu, "On the convergence of adaptive gradient methods for nonconvex optimization," *arXiv preprint arXiv:1808.05671*, 2018. 3, 4

[32] X. Chen, S. Liu, R. Sun, and M. Hong, "On the convergence of a class of adam-type algorithms for non-convex optimization," in *International Conference on Learning Representations*, 2019. 3, 4

[33] M. Zaheer, S. Reddi, D. Sachan, S. Kale, and S. Kumar, "Adaptive methods for nonconvex optimization," *Advances in neural information processing systems*, vol. 31, 2018. 3

[34] R. Ward, X. Wu, and L. Bottou, "Adagrad stepsizes: Sharp convergence over nonconvex landscapes," in *International Conference on Machine Learning*. PMLR, 2019, pp. 6677–6686. 3

[35] A. Défossez, L. Bottou, F. Bach, and N. Usunier, "On the convergence of adam and adagrad," *CoRR*, vol. abs/2003.02395, 2020. [Online]. Available: https://arxiv.org/abs/2003.02395 3

[36] F. Zou, L. Shen, Z. Jie, W. Zhang, and W. Liu, "A sufficient condition for convergences of adam and rmsprop," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 11 127–11 135. 3

[37] C. Chen, L. Shen, F. Zou, and W. Liu, "Towards practical adam: Non-convexity, convergence theory, and mini-batch acceleration," *arXiv preprint arXiv:2101.05471*, 2021. 3, 4

[38] C. Chen, L. Shen, H. Huang, and W. Liu, "Quantized adam with error feedback," *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 12, no. 5, pp. 1–26, 2021. 3

[39] C. Chen, L. Shen, F. Zou, and W. Liu, "Towards practical adam: Non-convexity, convergence theory, and mini-batch acceleration," *Journal of Machine Learning Research*, vol. 23, pp. 1–47, 2022. 3

[40] C. Chen, L. Shen, W. Liu, and Z.-Q. Luo, "Efficient-adam: Communication-efficient distributed adam with complexity analysis," *arXiv preprint arXiv:2205.14473*, 2022. 3

[41] F. Zou, L. Shen, Z. Jie, J. Sun, and W. Liu, "Weighted adagrad with unified momentum," *arXiv preprint arXiv:1808.03408*, 2018. 3

[42] H. Iiduka, "Appropriate learning rates of adaptive learning rate optimization algorithms for training deep neural networks," *IEEE Transactions on Cybernetics*, vol. 52, no. 12, pp. 13 250–13 261, 2021. 3

[43] S. Sun, Z. Cao, H. Zhu, and J. Zhao, "A survey of optimization methods from a machine learning perspective," *IEEE transactions on cybernetics*, vol. 50, no. 8, pp. 3668–3681, 2019. 3

[44] H. Sakai and H. Iiduka, "Riemannian adaptive optimization algorithm and its application to natural language processing," *IEEE Transactions on Cybernetics*, vol. 52, no. 8, pp. 7328–7339, 2021. 3

[45] B. T. Polyak, "Some methods of speeding up the convergence of iteration methods," *Ussr computational mathematics and mathematical physics*, vol. 4, no. 5, pp. 1–17, 1964. 3

[46] Y. Nesterov, *Introductory lectures on convex optimization: A basic course*. Springer Science & Business Media, 2003, vol. 87. 3

[47] B. O'Donoghue and E. J. Candès, "Adaptive restart for accelerated gradient schemes," *Found. Comput. Math.*, vol. 15, no. 3, pp. 715–732, 2015. 3

[48] W. Liu, L. Chen, Y. Chen, and W. Zhang, "Accelerating federated learning via momentum gradient descent," *IEEE Transactions on Parallel and Distributed Systems*, vol. 31, no. 8, pp. 1754–1766, 2020. 3

[49] X. Luo, Z. Liu, S. Li, M. Shang, and Z. Wang, "A fast non-negative latent factor model based on generalized momentum method," *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, vol. 51, no. 1, pp. 610–620, 2018. 3

[50] M. Shang, Y. Yuan, X. Luo, and M. Zhou, "An $\alpha-\beta$-divergence-generalized recommender for highly accurate predictions of missing user preferences," *IEEE transactions on cybernetics*, vol. 52, no. 8, pp. 8006–8018, 2021. 3

[51] T. Yang, Q. Lin, and Z. Li, "Unified convergence analysis of stochastic momentum methods for convex and non-convex optimization," *arXiv preprint arXiv:1604.03257*, 2016. 3

[52] S. S. Mannelli and P. Urbani, "Analytical study of momentum-based acceleration methods in paradigmatic high-dimensional non-convex problems," in *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, M. Ranzato, A. Beygelzimer, Y. N. Dauphin, P. Liang, and J. W. Vaughan, Eds., 2021, pp. 187–199. 3

[53] X. Gao, M. Gürbüzbalaban, and L. Zhu, "Global convergence of stochastic gradient hamiltonian monte carlo for nonconvex stochastic optimization: Nonasymptotic performance bounds and momentum-based acceleration," *Operations Research*, vol. 70, no. 5, pp. 2931–2947, 2022. 3

[54] I. Sutskever, J. Martens, G. Dahl, and G. Hinton, "On the importance of initialization and momentum in deep learning," in *International conference on machine learning*. PMLR, 2013, pp. 1139–1147. 3

[55] B. Can, M. Gürbüzbalaban, and L. Zhu, "Accelerated linear convergence of stochastic momentum methods in wasserstein distances," in *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, ser. Proceedings of Machine Learning Research, K. Chaudhuri and R. Salakhutdinov, Eds., vol. 97. PMLR, 2019, pp. 891–901. 3

[56] F. Huang, S. Gao, J. Pei, and H. Huang, "Accelerated zeroth-order and first-order momentum methods from mini to minimax optimization," *J. Mach. Learn. Res.*, vol. 23, pp. 36:1–36:70, 2022. 3

[57] R. Bollapragada, D. Scieur, and A. d'Aspremont, "Nonlinear acceleration of momentum and primal-dual algorithms," *Mathematical Programming*, pp. 1–38, 2022. 3

[58] B. O'donoghue and E. Candes, "Adaptive restart for accelerated gradient schemes," *Foundations of computational mathematics*, vol. 15, pp. 715–732, 2015. 3

[59] B. Wang, T. M. Nguyen, T. Sun, A. L. Bertozzi, R. G. Baraniuk, and S. J. Osher, "Scheduled restart momentum for accelerated stochastic gradient descent," *SIAM J. Imaging Sci.*, vol. 15, no. 2, pp. 738–761, 2022. 3

[60] B. Liu, L. Chai, and J. Yi, "Convergence analysis of distributed gradient descent algorithms with one and two momentum terms," *IEEE Transactions on Cybernetics*, 2022. 3

[61] A. Cutkosky and F. Orabona, "Momentum-based variance reduction in non-convex sgd," *Advances in neural information processing systems*, vol. 32, 2019. 4

[62] F. Huang, J. Li, and H. Huang, "Super-adam: faster and universal framework of adaptive gradients," *Advances in Neural Information Processing Systems*, vol. 34, 2021. 4

[63] M. Li, D. G. Andersen, A. J. Smola, and K. Yu, "Communication efficient distributed machine learning with the parameter server," *Advances in Neural Information Processing Systems*, vol. 27, 2014. 4

[64] M. Li, T. Zhang, Y. Chen, and A. J. Smola, "Efficient mini-batch training for stochastic optimization," in *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, 2014, pp. 661–670. 4

[65] A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, and S. Bowman, "Glue: A multi-task benchmark and analysis platform for natural language understanding," in *EMNLP*, 2018. 6

[66] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, "Roberta: A robustly optimized bert pretraining approach," *arXiv*, 2019. 7

[67] X. Li, B. Karimi, and P. Li, "On distributed adaptive optimization with gradient compression," in *International Conference on Learning Representations*, 2021. 7

In this supplementary material, we give additional discussion on this paper. In Appendix A, detailed experimental settings such as some hyper-parameters are listed. In Appendix B, we first give the proof, then we give some useful lemmas to help proving the main theorem. In Appendix C, we provide additional experiment illustration.

## APPENDIX A
### EXPERIMENTAL SETTINGS

TABLE III: Experimental settings and data divisions upon different downstream tasks. Notably, for each tasks in GLUE benchmark, we provide the number of classes ("classes"), the learning rate ("lr"), the batch size ("bsz"), the total number of updates ("total"), the number of warmup updates ("warmup") and the number of GPUs ("GPUs") during fine-tuning, respectively.

|  | MNLI | QNLI | QQP | RTE | SST-2 | MRPC | CoLA | STS-B |
|---|---|---|---|---|---|---|---|---|
| *experimental settings upon different downstream tasks* | | | | | | | | |
| –classes | 3 | 2 | 2 | 2 | 2 | 2 | 2 | 1 |
| –lr | 1e-5 | 1e-5 | 1e-5 | 2e-5 | 1e-5 | 1e-5 | 1e-5 | 2e-5 |
| –bsz | 256 | 128 | 256 | 32 | 64 | 32 | 32 | 32 |
| –total | 15,484 | 8,278 | 14,453 | 1,018 | 10,467 | 1,148 | 2,668 | 1,799 |
| –warmup | 929 | 496 | 867 | 61 | 628 | 68 | 160 | 107 |
| –GPUs | 4 | 4 | 8 | 2 | 2 | 2 | 2 | 2 |
| *data divisions for each dataset* | | | | | | | | |
| train | 392,720 | 104,743 | 363,870 | 2,491 | 67,350 | 5,801 | 8,551 | 5,749 |
| dev | 9,815 | 5,463 | 40,431 | 277 | 873 | 4,076 | 1,043 | 1,500 |
| test | 9,796 | 5,461 | 390,956 | 3,000 | 1,821 | 1,725 | 1,063 | 1,379 |

The GLUE benchmark contains 8 tasks, they are RTE, STS-B, CoLA, SST-2, MNLI, MRPC, QNLI and QQP. CoLA is a single sentence task. Each sentence has a label 1 and -1. 1 represents that it is a grammatical sentence, while -1 represents that it is illegal. Matthews correlation coefficient, dubbed **mcc** is used as our evaluation metric. STS-B is a similarity and paraphrase task. Each sample has a pair of a paragraph. People annotated the sample from 1 to 5 based on the similarity between the two paragraphs. The metric is Pearson and Spearman, dubbed **p/s** correlation coefficients. RTE is an inference task. Each sample has two sentences. If two sentences have a relation of entailment, we view them as a positive sample. If not, they compose of a negative sample. In the RTE task, the metric is the accuracy, dubbed **acc**. SST-2 is a single sentence task and its metric is the accuracy. MNLI is a sentence-level task that has 3 classes. They are entailment, contradiction and neutral. MRPC is a task to classify whether the sentences in the pair are equivalent. QNLI is a question-answering task. If the sentence contains the answer to the question, then it is a positive sample. QQP is a social question-answering task that consists of question pairs from Quora. It determines whether the questions are equivalent. The metric of MNLI, MRPC, QNLI, QQP is accuracy.

## APPENDIX B
### PROOF OF THE MAIN RESULTS

We set $z_t = x_t + \frac{\beta_1}{1-\beta_1}(x_t - x_{t-1})$ for $t \geq 0$ and we assume $x_{-1} = 0$ and $m_{-1} = 0$.

We have that

$$z_{t+1} - z_t = x_{t+1} + \frac{\beta_1}{1-\beta_1}(x_{t+1} - x_t) - x_t - \frac{\beta_1}{1-\beta_1}(x_t - x_{t-1}) \tag{21}$$

$$= \frac{1}{1-\beta_1}(x_{t+1} - x_t) - \frac{\beta_1}{1-\beta_1}(x_t - x_{t-1}) \tag{22}$$

$$= -\frac{1}{1-\beta_1}\gamma m_t \odot \eta_t + \frac{\beta_1}{1-\beta_1}(x_t - x_{t-1})\gamma m_{t-1} \odot \eta_{t-1} \tag{23}$$

$$= -\frac{1}{1-\beta_1}\gamma(\beta_1 m_{t-1} + (1-\beta_1)g_t) \odot \eta_t + \frac{\beta_1}{1-\beta_1}(x_t - x_{t-1})\gamma m_{t-1} \odot \eta_{t-1} \tag{24}$$

$$= \frac{\beta_1}{1-\beta_1}\gamma m_{t-1} \odot (\eta_{t-1} - \eta_t) - \gamma g_t \odot \eta_t \tag{25}$$

By applying L-smooth, we have

$$f(z_{t+1}) \leq f(z_t) + \langle \nabla f(z_t), z_{t+1} - z_t \rangle + \frac{L}{2}\|z_{t+1} - z_t\|^2 \tag{26}$$

We re-organize it, and we have

$$f(z_{t+1}) - f(z_t)$$
$$\leq \langle \nabla f(z_t), z_{t+1} - z_t \rangle + \frac{L}{2}\|z_{t+1} - z_t\|^2 \tag{27}$$

$$= \langle \nabla f(z_t), \frac{\gamma \beta_1}{1-\beta_1} m_{t-1} \odot (\eta_{t-1} - \eta_t) \rangle + \langle \nabla f(z_t), -\gamma g_t \odot \eta_t \rangle + \frac{L}{2} \|z_{t+1} - z_t\|^2 \tag{28}$$

$$= \langle \nabla f(z_t), \frac{\gamma \beta_1}{1-\beta_1} m_{t-1} \odot (\eta_{t-1} - \eta_t) \rangle + \frac{L}{2} \|z_{t+1} - z_t\|^2$$
$$+ \langle \nabla f(z_t), \frac{\gamma_t}{b} \sum_{i \in B} \nabla f_i(x_t + \rho_t \frac{s_t}{\|s_t\|}) \odot (\eta_{t-1} - \eta_t) \rangle$$
$$+ \langle \nabla f(z_t), -\frac{\gamma_t}{b} \sum_{i \in B} \nabla f_i(x_t + \rho_t \frac{s_t}{\|s_t\|}) \odot \eta_{t-1} \rangle \tag{29}$$

$$= \langle \nabla f(z_t), \frac{\gamma \beta_1}{1-\beta_1} m_{t-1} \odot (\eta_{t-1} - \eta_t) \rangle + \frac{L}{2} \|z_{t+1} - z_t\|^2$$
$$+ \langle \nabla f(z_t), \frac{\gamma_t}{b} \sum_{i \in B} \nabla f_i(x_t + \rho_t \frac{s_t}{\|s_t\|}) \odot (\eta_{t-1} - \eta_t) \rangle$$
$$+ \langle \nabla f(z_t) - \nabla f(x_t), -\frac{\gamma_t}{b} \sum_{i \in B} \nabla f_i(x_t + \rho_t \frac{s_t}{\|s_t\|}) \odot \eta_{t-1} \rangle$$
$$+ \langle \nabla f(x_t), -\frac{\gamma_t}{b} \sum_{i \in B} \nabla f_i(x_t + \rho_t \frac{s_t}{\|s_t\|}) \odot \eta_{t-1} \rangle \tag{30}$$

$$= \langle \nabla f(z_t), \frac{\gamma \beta_1}{1-\beta_1} m_{t-1} \odot (\eta_{t-1} - \eta_t) \rangle + \frac{L}{2} \|z_{t+1} - z_t\|^2$$
$$+ \langle \nabla f(z_t), \frac{\gamma_t}{b} \sum_{i \in B} \nabla f_i(x_t + \rho_t \frac{s_t}{\|s_t\|}) \odot (\eta_{t-1} - \eta_t) \rangle$$
$$+ \langle \nabla f(z_t) - \nabla f(x_t), -\frac{\gamma_t}{b} \sum_{i \in B} \nabla f_i(x_t + \rho_t \frac{\sum \nabla f_i(x_t)}{\|\sum \nabla f_i(x_t)\|}) \odot \eta_{t-1} \rangle$$
$$+ \langle \nabla f(x_t), \frac{\gamma_t}{b} \sum_{i \in B} \nabla f_i(x_t + \rho_t \frac{\nabla f(x_t)}{\|\nabla f(x_t)\|}) \odot \eta_{t-1} - \frac{\gamma_t}{b} \sum_{i \in B} \nabla f_i(x_t + \rho_t \frac{s_t}{\|s_t\|}) \odot \eta_{t-1} \rangle$$
$$+ \langle \nabla f(x_t), -\frac{\gamma_t}{b} \sum_{i \in B} \nabla f_i(x_t + \rho_t \frac{\nabla f(x_t)}{\|\nabla f(x_t)\|}) \odot \eta_{t-1} \rangle. \tag{31}$$

From the Lemma 5, Lemma 6, Lemma 7, we have

$$\langle \nabla f(z_t), \frac{\gamma_t}{b} \sum_{i \in B} \nabla f_i(x_t + \rho_t \frac{s_t}{\|s_t\|}) \odot (\eta_{t-1} - \eta_t) \rangle \le \gamma_t G^2 \|\eta_{t-1} - \eta_t\|_1, \tag{32}$$

$$\langle \nabla f(z_t), \frac{\gamma \beta_1}{1-\beta_1} m_{t-1} \odot (\eta_{t-1} - \eta_t) \rangle \le \frac{\gamma \beta_1}{1-\beta_1} G^2 \|\eta_{t-1} - \eta_t\|_1, \tag{33}$$

$$\langle \nabla f(x_t), \frac{\eta_t}{b} \sum_{i \in B} \nabla f_i(x_t + \rho_t \frac{\nabla f(x_t)}{\|\nabla f(x_t)\|}) \odot \eta_{t-1} - \frac{\gamma_t}{b} \sum_{i \in B} \nabla f_i(x_t + \rho_t \frac{s_t}{\|s_t\|}) \odot \eta_{t-1} \rangle$$
$$\le \frac{\gamma_t}{2\mu^2} \|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2 + \frac{2\mu^2 \gamma_t L^2 \rho_t^2}{\epsilon}. \tag{34}$$

Taking conditional expectation, we have

$$\mathbb{E} f(z_{t+1}) - f(z_t) \tag{35}$$
$$\le \mathbb{E} \langle \nabla f(x_t), -\frac{\gamma_t}{b} \sum_{i \in B} \nabla f_i(x_t + \rho_t \frac{\nabla f(x_t)}{\|\nabla f(x_t)\|}) \odot \eta_{t-1} \rangle + \frac{L}{2} \mathbb{E} \|z_{t+1} - z_t\|^2$$
$$+ \frac{\gamma_t}{2\mu^2} \|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2 + \frac{2\mu^2 \gamma_t L^2 \rho_t^2}{\epsilon} + \frac{\gamma}{1-\beta_1} G^2 \|\eta_{t-1} - \eta_t\|_1$$
$$+ \mathbb{E} \langle \nabla f(z_t) - \nabla f(x_t), -\frac{\gamma_t}{b} \sum_{i \in B} \nabla f_i(x_t + \rho_t \frac{s_t}{\|s_t\|}) \odot \eta_{t-1} \rangle \tag{36}$$

where $\mu > 0$ is to be determined.

For the term

$$\mathbb{E} \langle \nabla f(x_t), -\frac{\gamma_t}{b} \sum_{i \in B} \nabla f_i(x_t + \rho_t \frac{\nabla f(x_t)}{\|\nabla f(x_t)\|}) \odot \eta_{t-1} \rangle, \tag{37}$$

the term

$$\frac{L}{2} \mathbb{E} \|z_{t+1} - z_t\|^2, \tag{38}$$

and the term

$$\mathbb{E}\langle \nabla f(z_t) - \nabla f(x_t), -\frac{\gamma_t}{b}\sum_{i\in B}\nabla f_i(x_t + \rho_t \frac{s_t}{\|s_t\|}) \odot \eta_{t-1}\rangle, \tag{39}$$

we introduce the Lemma 8, the Lemma 10 and the Lemma 9. We take the expectation over the whole processing and we have

$$\mathbb{E}f(z_{t+1}) - \mathbb{E}f(z_t)$$

$$\leq \frac{\gamma_t}{2\mu^2}\mathbb{E}\|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2 + \frac{2\mu^2\gamma_t L^2\rho_t^2}{\epsilon} + \frac{\gamma}{1-\beta_1}G^2\mathbb{E}\|\eta_{t-1}-\eta_t\|_1$$

$$- \gamma_t\mathbb{E}\|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2 + \mathbb{E}\frac{\gamma_t}{2\alpha^2}\mathbb{E}\|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2 + \frac{\gamma_t\alpha^2 L^2\rho^2}{2\epsilon} + \frac{LG^2\gamma^2\beta_1^2}{(1-\beta_1)^2}\mathbb{E}\|\eta_t-\eta_{t-1}\|^2$$

$$+ \gamma_t^2 L(3\frac{1+\beta}{\beta\epsilon}(\mathbb{E}\|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2 + \frac{L\rho_t^2}{\epsilon} + \frac{\sigma^2}{b\epsilon}) + (1+\beta)G^2\mathbb{E}\|\eta_t-\eta_{t-1}\|^2)$$

$$+ \frac{\gamma^3 L^2\beta_1^2}{2\epsilon(1-\beta_1)^2}(\frac{1}{\lambda_1^2} + \frac{1}{\lambda_2^2} + \frac{1}{\lambda_3^2})\frac{dG_\infty^2}{\epsilon^2} + \frac{\gamma\lambda_1^2}{2}\|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2 + \frac{\gamma L^2\rho_t^2}{2\epsilon}(\lambda_2^2 + 4\lambda_3^2) \tag{40}$$

$$= -\gamma_t(1 - \frac{1}{2\mu^2} - \frac{1}{2\alpha^2} - \frac{3\gamma L(1+\beta)}{\beta\epsilon} - \frac{\lambda_1^2}{2})\mathbb{E}\|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2 + \frac{2\mu^2\gamma_t L^2\rho_t^2}{\epsilon} + \frac{\gamma}{1-\beta_1}G^2\mathbb{E}\|\eta_{t-1}-\eta_t\|_1$$

$$+ \frac{\gamma_t\alpha^2 L^2\rho^2}{2\epsilon} + \frac{3\gamma_t^2 L(1+\beta)}{\beta\epsilon}(\frac{L\rho_t^2}{\epsilon} + \frac{\sigma^2}{b\epsilon}) + \gamma_t^2 LG^2((\frac{\beta_1}{1-\beta_1})^2 + 1 + \beta)\mathbb{E}\|\eta_t-\eta_{t-1}\|^2$$

$$+ \frac{\gamma^3 L^2\beta_1^2}{2\epsilon(1-\beta_1)^2}(\frac{1}{\lambda_1^2} + \frac{1}{\lambda_2^2} + \frac{1}{\lambda_3^2})\frac{dG_\infty^2}{\epsilon^2} + \frac{\gamma L^2\rho_t^2}{2\epsilon}(\lambda_2^2 + 4\lambda_3^2). \tag{41}$$

We set $\mu^2 = \alpha^2 = 8$, $\beta = 3$, $\lambda_1^2 = \frac{1}{4}$, $\lambda_2^2 = \lambda_3^2 = 1$ and we choose $\frac{2\gamma_t L}{\epsilon} \leq \frac{1}{8}$. So we have

$$\mathbb{E}f(x_{t+1}) - \mathbb{E}f(x_t)$$

$$\leq -\frac{\gamma_t}{2}\mathbb{E}\|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2 + \frac{16\gamma_t L^2\rho_t^2}{\epsilon} + \frac{\gamma}{1-\beta_1}G^2\mathbb{E}\|\eta_{t-1}-\eta_t\|_1$$

$$+ \frac{4\gamma_t L^2\rho^2}{\epsilon} + \frac{4\gamma_t^2 L}{\epsilon}(\frac{L\rho_t^2}{\epsilon} + \frac{\sigma^2}{b\epsilon}) + (4 + (\frac{\beta_1}{1-\beta_1})^2)\gamma_t^2 LG^2\mathbb{E}\|\eta_t-\eta_{t-1}\|^2$$

$$+ \frac{3\gamma^3 L^2\beta_1^2}{\epsilon(1-\beta_1)^2}\frac{dG_\infty^2}{\epsilon^2} + \frac{5\gamma L^2\rho_t^2}{2\epsilon} \tag{42}$$

We re-arrange it and $\eta_t$ is bounded. We have

$$\frac{\gamma_t}{2G}\mathbb{E}\|\nabla f(x_t)\|^2 \leq \frac{\gamma_t}{2}\mathbb{E}\|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2 \tag{43}$$

$$\leq -\mathbb{E}f(x_{t+1}) + \mathbb{E}f(x_t) + \frac{45\gamma_t L^2\rho_t^2}{2\epsilon} + \frac{\gamma}{1-\beta_1}G^2\mathbb{E}\|\eta_{t-1}-\eta_t\|_1$$

$$+ \frac{4\gamma_t^2 L}{\epsilon}(\frac{L\rho_t^2}{\epsilon} + \frac{\sigma^2}{b\epsilon}) + (4 + (\frac{\beta_1}{1-\beta_1})^2)\gamma_t^2 LG^2\mathbb{E}\|\eta_t-\eta_{t-1}\|^2 + \frac{3\gamma^3 L^2\beta_1^2}{(1-\beta_1)^2}\frac{dG_\infty^2}{\epsilon^3}. \tag{44}$$

We summary it from $t = 0$ to $t = T - 1$, and we assume $\gamma_t$ is a constant, and we have

$$\frac{1}{T}\sum_{t=0}^{T-1}\mathbb{E}\|\nabla f(x_t)\|^2 \leq 2G\frac{\mathbb{E}f(x_0) - \mathbb{E}f(x_{t+1})}{\gamma_t T} + \frac{45GL^2\rho_t^2}{\epsilon} + \frac{2G^3}{(1-\beta_1)T}\mathbb{E}\sum_{t=0}^{T-1}\|\eta_{t-1}-\eta_t\|_1$$

$$+ \frac{8G\gamma_t L}{\epsilon}(\frac{L\rho_t^2}{\epsilon} + \frac{\sigma^2}{b\epsilon}) + \frac{2(4 + (\frac{\beta_1}{1-\beta_1})^2)\gamma_t LG^3}{T}\mathbb{E}\sum_{t=0}^{T-1}\|\eta_t-\eta_{t-1}\|^2 + \frac{6\gamma^2 L^2\beta_1^2}{(1-\beta_1)^2}\frac{dG^3}{\epsilon^3} \tag{45}$$

$$\leq \frac{2G(f(x_0) - f^*)}{\gamma_t T} + \frac{45GL^2\rho_t^2}{\epsilon} + \frac{2G^3}{(1-\beta_1)T}d(\frac{1}{\epsilon} - \frac{1}{G}) + \frac{8G\gamma_t L}{\epsilon}(\frac{L\rho_t^2}{\epsilon} + \frac{\sigma^2}{b\epsilon})$$

$$+ \frac{2(4 + (\frac{\beta_1}{1-\beta_1})^2)\gamma_t LG^3}{T}d(\epsilon^{-2} - G^{-2}) + \frac{6\gamma^2 L^2\beta_1^2}{(1-\beta_1)^2}\frac{dG^3}{\epsilon^3} \tag{46}$$

$$= \frac{2G(f(x_0) - f^*)}{\gamma_t T} + \frac{8G\gamma_t L}{\epsilon}\frac{\sigma^2}{b\epsilon} + \frac{45GL^2\rho_t^2}{\epsilon} + \frac{2G^3}{(1-\beta_1)T}d(\frac{1}{\epsilon} - \frac{1}{G}) + \frac{8G\gamma_t L}{\epsilon}\frac{L\rho_t^2}{\epsilon}$$

$$+ \frac{2(4 + (\frac{\beta_1}{1-\beta_1})^2)\gamma_t LG^3}{T}d(\epsilon^{-2} - G^{-2}) + \frac{6\gamma^2 L^2\beta_1^2}{(1-\beta_1)^2}\frac{dG^3}{\epsilon^3}. \tag{47}$$

*A. Technical Lemma*

**Lemma 1.** *Given two vectors $a$, $b \in \mathbb{R}^d$, we have $\langle a, b \rangle \leq \frac{\lambda^2}{2}\|a\|^2 + \frac{1}{2\lambda^2}\|b\|^2$ for parameter $\lambda$, $\forall \lambda \in (1, +\infty)$.*

*Proof.*

$$RHS = \frac{\lambda^2}{2}\sum_{j=1}^{d}(a)_j^2 + \frac{1}{2\lambda^2}\sum_{j=1}^{d}(b)_j^2 \geq \sum_{j=1}^{d}2\sqrt{\frac{\lambda^2}{2}(a)_j^2 \times \frac{1}{2\lambda^2}(b)_j^2} = \sum_{j=1}^{d}|(a)_j| \times |(b)_j| \geq LHS. \tag{48}$$

$\square$

**Lemma 2.** *For any vector $x, y \in \mathbb{R}^d$, we have*

$$\|x \odot y\|^2 \leq \|x\|^2 \times \|y\|_\infty^2 \leq \|x\|^2 \times \|y\|^2. \tag{49}$$

*Proof.* The first inequality can be derived from that $\sum_{i=1}^{d}(x_i^2 y_i^2) \leq \sum_{i=1}^{d}(x_i^2\|y\|_\infty^2)$. The second inequality follows from that $\|y\|_\infty^2 \leq \|y\|^2$. $\square$

**Lemma 3.** *$\eta$ is bounded, i.e., $\frac{1}{G_\infty} \leq (\eta_t)_j \leq \frac{1}{\epsilon}$.*

*Proof.* As the gradient is bounded by $G$ and $(\eta_t)_j = \frac{1}{\sqrt{(\hat{v}_t)_j}}$. Follow the update rule, we have $\frac{1}{G_\infty} \leq (\eta_t)_j \leq \frac{1}{\epsilon}$. $\square$

**Lemma 4.** *For the term defined in the algorithm, we have*

$$\frac{1}{T}\mathbb{E}\sum_{t=0}^{T-1}\|\eta_{t-1} - \eta_t\|^1 \leq \frac{d}{T}\left(\frac{1}{\epsilon} - \frac{1}{G}\right) \tag{50}$$

*Proof.* $(\eta_t)_i$, the i-th dimension of $\eta_t$ deceases as t increases. So we have

$$\frac{1}{T}\mathbb{E}\sum_{t=0}^{T-1}\|\eta_{t-1} - \eta_t\|^1 = \mathbb{E}\frac{1}{T}\sum_{i=1}^{d}\sum_{t=0}^{T-1}|(\eta_{t-1})_i - (\eta_t)_i|$$

$$\leq \mathbb{E}\frac{1}{T}\sum_{i=1}^{d}((\eta_{-1})_i - (\eta_{T-1})_i) \leq \mathbb{E}\frac{1}{T}\sum_{i=1}^{d}\left(\frac{1}{\epsilon} - \frac{1}{G}\right) = \frac{d}{T}\left(\frac{1}{\epsilon} - \frac{1}{G}\right) \tag{51}$$

$\square$

**Lemma 5.** *For the term defined in the algorithm, we have*

$$\langle \nabla f(z_t), \frac{\gamma_t}{b}\sum_{i \in B}\nabla f_i(x_t + \rho_t\frac{s_t}{\|s_t\|}) \odot (\eta_{t-1} - \eta_t)\rangle \leq \gamma_t G^2\|\eta_{t-1} - \eta_t\|_1 \tag{52}$$

*Proof.*

$$\langle \nabla f(z_t), \frac{\gamma_t}{b}\sum_{i \in B}\nabla f_i(x_t + \rho_t\frac{s_t}{\|s_t\|}) \odot (\eta_{t-1} - \eta_t)\rangle$$

$$\leq \gamma_t\sum_{j=1}^{d}|(\nabla f(z_t))_{(j)}| \times |(\frac{1}{b}\sum_{i \in B}\nabla f_i(x_t + \rho_t\frac{\sum \nabla f_i(x_t)}{\|\sum \nabla f_i(x_t)\|}) \odot (\eta_{t-1} - \eta_t))_{(j)}| \tag{53}$$

$$\leq \gamma_t G\sum_{j=1}^{d}|((\frac{1}{b}\sum_{i \in B}\nabla f_i(x_t + \rho_t\frac{\sum \nabla f_i(x_t)}{\|\sum \nabla f_i(x_t)\|}) \odot (\eta_{t-1} - \eta_t))_{(j)}| \tag{54}$$

$$\leq \frac{\gamma_t G}{b}\sum_{j=1}^{d}\sum_{i \in B}|((\nabla f_i(x_t + \rho_t\frac{\sum \nabla f_i(x_t)}{\|\sum \nabla f_i(x_t)\|}) \odot (\eta_{t-1} - \eta_t))_{(j)}| \tag{55}$$

$$= \frac{\gamma_t G}{b}\sum_{j=1}^{d}\sum_{i \in B}|(\nabla f_i(x_t + \rho_t\frac{\sum \nabla f_i(x_t)}{\|\sum \nabla f_i(x_t)\|})_{(j)} \times (\eta_{t-1} - \eta_t)_{(j)}| \tag{56}$$

$$\leq \frac{\gamma_t G^2}{b}\sum_{j=1}^{d}\sum_{i \in B}|(\eta_{t-1} - \eta_t)_{(j)}| \tag{57}$$

$$= \gamma_t G^2\|\eta_{t-1} - \eta_t\|_1 \tag{58}$$

$\square$

**Lemma 6.** *For the term defined in the algorithm, we have*

$$\langle \nabla f(z_t), \frac{\gamma\beta_1}{1-\beta_1} m_{t-1} \odot (\eta_{t-1} - \eta_t) \rangle \leq \frac{\gamma\beta_1}{1-\beta_1} G^2 \|\eta_{t-1} - \eta_t\|_1 \tag{59}$$

*Proof.*

$$\langle \nabla f(z_t), \frac{\gamma\beta_1}{1-\beta_1} m_{t-1} \odot (\eta_{t-1} - \eta_t) \rangle$$

$$\leq \frac{\gamma\beta_1}{1-\beta_1} \sum_{j=1}^{d} |(\nabla f(z_t))_{(j)}| \times |(m_{t-1} \odot (\eta_{t-1} - \eta_t))_{(j)}| \tag{60}$$

$$\leq \frac{\gamma\beta_1}{1-\beta_1} G \sum_{j=1}^{d} |(m_{t-1} \odot (\eta_{t-1} - \eta_t))_{(j)}| \tag{61}$$

$$= \frac{\gamma\beta_1}{1-\beta_1} \sum_{j=1}^{d} |(m_{t-1})_{(j)} \times (\eta_{t-1} - \eta_t)_{(j)}| \tag{62}$$

$$\leq \frac{\gamma\beta_1}{1-\beta_1} G^2 \sum_{j=1}^{d} |(\eta_{t-1} - \eta_t)_{(j)}| \tag{63}$$

$$= \frac{\gamma\beta_1}{1-\beta_1} G^2 \|\eta_{t-1} - \eta_t\|_1 \tag{64}$$

$\square$

**Lemma 7.** *For the term defined in the algorithm, we have*

$$\langle \nabla f(x_t), \frac{\gamma_t}{b} \sum_{i\in B} \nabla f_i(x_t + \rho_t \frac{\nabla f(x_t)}{\|\nabla f(x_t)\|}) \odot \eta_{t-1} - \frac{\gamma_t}{b} \sum_{i\in B} \nabla f_i(x_t + \rho_t \frac{s_t}{\|s_t\|}) \odot \eta_{t-1} \rangle$$

$$\leq \frac{\gamma_t}{2\mu^2} \|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2 + \frac{2\mu^2\gamma_t L^2\rho_t^2}{\epsilon}. \tag{65}$$

*Proof.*

$$\langle \nabla f(x_t), \frac{\gamma_t}{b} \sum_{i\in B} \nabla f_i(x_t + \rho_t \frac{\nabla f(x_t)}{\|\nabla f(x_t)\|}) \odot \eta_{t-1} - \frac{\gamma_t}{b} \sum_{i\in B} \nabla f_i(x_t + \rho_t \frac{s_t}{\|s_t\|}) \odot \eta_{t-1} \rangle$$

$$= \langle \nabla f(x_t) \odot \sqrt{\eta_{t-1}}, \frac{\gamma_t}{b} \sum_{i\in B} (\nabla f_i(x_t + \rho_t \frac{\nabla f(x_t)}{\|\nabla f(x_t)\|}) - \nabla f_i(x_t + \rho_t \frac{\sum_{i\in B} \nabla f_i(x_t)}{\|\sum_{i\in B} \nabla f_i(x_t)\|})) \odot \sqrt{\eta_{t-1}} \rangle \tag{66}$$

$$\leq \frac{\mu^2\gamma_t}{2b^2} \|\sum (\nabla f_i(x_t + \rho_t \frac{\nabla f(x_t)}{\|\nabla f(x_t)\|}) - \nabla f_i(x_t + \rho_t \frac{\sum_{i\in B} \nabla f_i(x_t)}{\|\sum_{i\in B} \nabla f_i(x_t)\|})) \odot \sqrt{\eta_{t-1}}\|^2$$

$$+ \frac{\gamma_t}{2\mu^2} \|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2 \tag{67}$$

$$\leq + \frac{\mu^2\gamma_t}{2b} \sum \|\nabla f_i(x_t + \rho_t \frac{\nabla f(x_t)}{\|\nabla f(x_t)\|}) - \nabla f_i(x_t + \rho_t \frac{\sum_{i\in B} \nabla f_i(x_t)}{\|\sum_{i\in B} \nabla f_i(x_t)\|}) \odot \sqrt{\eta_{t-1}}\|^2$$

$$+ \frac{\gamma_t}{2\mu^2} \|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2 \tag{68}$$

$$\leq + \frac{\mu^2\gamma_t}{2b} \sum \|\nabla f_i(x_t + \rho_t \frac{\nabla f(x_t)}{\|\nabla f(x_t)\|}) - \nabla f_i(x_t + \rho_t \frac{\sum_{i\in B} \nabla f_i(x_t)}{\|\sum_{i\in B} \nabla f_i(x_t)\|})\|^2 \times \|\sqrt{\eta_{t-1}}\|_\infty^2$$

$$+ \frac{\gamma_t}{2\mu^2} \|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2 \tag{69}$$

$$\leq \frac{\gamma_t}{2\mu^2} \|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2 + \frac{\mu^2\gamma_t L^2\rho_t^2}{2b\epsilon} \sum \|\frac{\nabla f(x_t)}{\|\nabla f(x_t)\|} - \frac{\sum_{i\in B} \nabla f_i(x_t)}{\|\sum_{i\in B} \nabla f_i(x_t)\|}\|^2 \tag{70}$$

$$\leq \frac{\gamma_t}{2\mu^2} \|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2 + \frac{2\mu^2\gamma_t L^2\rho_t^2}{\epsilon}. \tag{71}$$

$\square$

**Lemma 8.** *For the term defined in the algorithm, we have*

$$\mathbb{E}\langle \nabla f(x_t), -\frac{\gamma_t}{b}\sum_{i\in B}\nabla f_i(x_t + \rho_t\frac{\nabla f(x_t)}{\|\nabla f(x_t)\|})\odot\eta_{t-1}\rangle$$

$$\leq -\gamma_t\|\nabla f(x_t)\odot\sqrt{\eta_{t-1}}\|^2 + \mathbb{E}\frac{\gamma_t}{2\alpha^2}\|\nabla f(x_t)\odot\sqrt{\eta_{t-1}}\|^2 + \frac{\gamma_t\alpha^2 L^2\rho_t^2}{2\epsilon} \tag{72}$$

*Proof.*

$$\mathbb{E}\langle \nabla f(x_t), -\frac{\gamma_t}{b}\sum_{i\in B}\nabla f_i(x_t + \rho_t\frac{\nabla f(x_t)}{\|\nabla f(x_t)\|})\odot\eta_{t-1}\rangle$$

$$= -\gamma_t\|\nabla f(x_t)\odot\sqrt{\eta_{t-1}}\|^2 + \mathbb{E}\langle \nabla f(x_t), \frac{\gamma_t}{b}\sum_{i\in B}(\nabla f(x_t) - \nabla f_i(x_t + \rho_t\frac{\nabla f(x_t)}{\|\nabla f(x_t)\|}))\odot\eta_{t-1}\rangle \tag{73}$$

$$= -\gamma_t\|\nabla f(x_t)\odot\sqrt{\eta_{t-1}}\|^2 + \mathbb{E}\langle \nabla f(x_t), \frac{\gamma_t}{b}\sum_{i\in B}(\nabla f_i(x_t) - \nabla f_i(x_t + \rho_t\frac{\nabla f(x_t)}{\|\nabla f(x_t)\|}))\odot\eta_{t-1}\rangle \tag{74}$$

$$\leq -\gamma_t\|\nabla f(x_t)\odot\sqrt{\eta_{t-1}}\|^2 + \mathbb{E}\frac{\gamma_t}{2\alpha^2}\|\nabla f(x_t)\odot\sqrt{\eta_{t-1}}\|^2$$

$$+ \frac{\gamma_t\alpha^2}{2}\mathbb{E}\|\frac{1}{b}\sum_{i\in B}(\nabla f_i(x_t) - \nabla f_i(x_t + \rho_t\frac{\nabla f(x_t)}{\|\nabla f(x_t)\|}))\odot\sqrt{\eta_{t-1}}\|^2 \tag{75}$$

$$\leq -\gamma_t\|\nabla f(x_t)\odot\sqrt{\eta_{t-1}}\|^2 + \mathbb{E}\frac{\gamma_t}{2\alpha^2}\|\nabla f(x_t)\odot\sqrt{\eta_{t-1}}\|^2$$

$$+ \frac{\gamma_t\alpha^2}{2\epsilon}\mathbb{E}\|\frac{1}{b}\sum_{i\in B}(\nabla f_i(x_t) - \nabla f_i(x_t + \rho_t\frac{\nabla f(x_t)}{\|\nabla f(x_t)\|}))\|^2 \tag{76}$$

$$\leq -\gamma_t\|\nabla f(x_t)\odot\sqrt{\eta_{t-1}}\|^2 + \mathbb{E}\frac{\gamma_t}{2\alpha^2}\|\nabla f(x_t)\odot\sqrt{\eta_{t-1}}\|^2$$

$$+ \frac{\gamma_t\alpha^2}{2b\epsilon}\mathbb{E}\sum_{i\in B}\|(\nabla f_i(x_t) - \nabla f_i(x_t + \rho_t\frac{\nabla f(x_t)}{\|\nabla f(x_t)\|}))\|^2 \tag{77}$$

$$\leq -\gamma_t\|\nabla f(x_t)\odot\sqrt{\eta_{t-1}}\|^2 + \mathbb{E}\frac{\gamma_t}{2\alpha^2}\|\nabla f(x_t)\odot\sqrt{\eta_{t-1}}\|^2 + \frac{\gamma_t\alpha^2 L^2\rho_t^2}{2b\epsilon}\mathbb{E}\sum_{i\in B}\|\frac{\nabla f(x_t)}{\|\nabla f(x_t)\|}\|^2 \tag{78}$$

$$= -\gamma_t\|\nabla f(x_t)\odot\sqrt{\eta_{t-1}}\|^2 + \mathbb{E}\frac{\gamma_t}{2\alpha^2}\|\nabla f(x_t)\odot\sqrt{\eta_{t-1}}\|^2 + \frac{\gamma_t\alpha^2 L^2\rho_t^2}{2\epsilon} \tag{79}$$

$$\square$$

**Lemma 9.** *For the term defined in the algorithm, we have*

$$\mathbb{E}\langle \nabla f(z_t) - \nabla f(x_t), -\frac{\gamma_t}{b}\sum_{i\in B}\nabla f_i(x_t + \rho_t\frac{s_t}{\|s_t\|})\odot\eta_{t-1}\rangle$$

$$\leq \frac{\gamma^3 L^2\beta_1^2}{2\epsilon(1-\beta_1)^2}(\frac{1}{\lambda_1^2} + \frac{1}{\lambda_2^2} + \frac{1}{\lambda_3^2})\frac{dG_\infty^2}{\epsilon^2} + \frac{\gamma\lambda_1^2}{2}\|\nabla f(x_t)\odot\sqrt{\eta_{t-1}}\|^2 + \frac{\gamma L^2\rho_t^2}{2\epsilon}(\lambda_2^2 + 4\lambda_3^2). \tag{80}$$

*Proof.*

$$\mathbb{E}\langle \nabla f(z_t) - \nabla f(x_t), -\frac{\gamma_t}{b}\sum_{i\in B}\nabla f_i(x_t + \rho_t\frac{s_t}{\|s_t\|})\odot\eta_{t-1}\rangle \tag{81}$$

$$= \gamma\mathbb{E}\langle (\nabla f(x_t) - \nabla f(z_t))\odot\sqrt{\eta_{t-1}}, \frac{1}{b}\sum_{i\in B}\nabla f_i(x_t + \rho_t\frac{\sum_{i\in B}\nabla f_i(x_t)}{\|\sum_{i\in B}\nabla f_i(x_t)\|})\odot\sqrt{\eta_{t-1}}\rangle \tag{82}$$

$$= \gamma\mathbb{E}\langle (\nabla f(x_t) - \nabla f(z_t))\odot\sqrt{\eta_{t-1}}, \nabla f(x_t)\odot\sqrt{\eta_{t-1}}\rangle$$

$$+ \gamma\mathbb{E}\langle (\nabla f(x_t) - \nabla f(z_t))\odot\sqrt{\eta_{t-1}}, \frac{1}{b}\sum_{i\in B}(\nabla f_i(x_t + \rho_t\frac{\nabla f(x_t)}{\|\nabla f(x_t)\|}) - \nabla f_i(x_t))\odot\sqrt{\eta_{t-1}}\rangle$$

$$+ \gamma\mathbb{E}\langle (\nabla f(x_t) - \nabla f(z_t))\odot\sqrt{\eta_{t-1}}, \frac{1}{b}\sum_{i\in B}(\nabla f_i(x_t + \rho_t\frac{\sum_{i\in B}\nabla f_i(x_t)}{\|\sum_{i\in B}\nabla f_i(x_t)\|}) - \nabla f_i(x_t + \rho_t\frac{\nabla f(x_t)}{\|\nabla f(x_t)\|}))\odot\sqrt{\eta_{t-1}}\rangle \tag{83}$$

$$\leq \frac{\gamma}{2}(\frac{1}{\lambda_1^2} + \frac{1}{\lambda_2^2} + \frac{1}{\lambda_3^2})\mathbb{E}\|(\nabla f(x_t) - \nabla f(z_t))\odot\sqrt{\eta_{t-1}}\|^2 + \frac{\gamma\lambda_1^2}{2}\|\nabla f(x_t)\odot\sqrt{\eta_{t-1}}\|^2$$

$$+ \frac{\gamma\lambda_2^2}{2}\mathbb{E}\|\frac{1}{b}\sum_{i\in B}(\nabla f_i(x_t + \rho_t\frac{\nabla f(x_t)}{\|\nabla f(x_t)\|}) - \nabla f_i(x_t))\odot\sqrt{\eta_{t-1}}\|^2$$

$$+ \frac{\gamma\lambda_3^2}{2}\mathbb{E}\|\frac{1}{b}\sum_{i\in B}(\nabla f_i(x_t + \rho_t \frac{\sum_{i\in B}\nabla f_i(x_t)}{\|\sum_{i\in B}\nabla f_i(x_t)\|}) - \nabla f_i(x_t + \rho_t \frac{\nabla f(x_t)}{\|\nabla f(x_t)\|}) \odot \sqrt{\eta_{t-1}}\|^2 \tag{84}$$

$$\leq \frac{\gamma}{2}(\frac{1}{\lambda_1^2} + \frac{1}{\lambda_2^2} + \frac{1}{\lambda_3^2})\mathbb{E}\|(\nabla f(x_t) - \nabla f(z_t)) \odot \sqrt{\eta_{t-1}}\|^2 + \frac{\gamma\lambda_1^2}{2}\|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2$$

$$+ \frac{\gamma\lambda_2^2 L^2 \rho_t^2}{2\epsilon} + \frac{2\lambda_3^2 \gamma L^2 \rho_t^2}{\epsilon} \tag{85}$$

$$\leq \frac{\gamma L^2}{2\epsilon}(\frac{1}{\lambda_1^2} + \frac{1}{\lambda_2^2} + \frac{1}{\lambda_3^2})\mathbb{E}\|z_t - x_t\|^2 + \frac{\gamma\lambda_1^2}{2}\|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2$$

$$+ \frac{\gamma\lambda_2^2 L^2 \rho_t^2}{2\epsilon} + \frac{2\lambda_3^2 \gamma L^2 \rho_t^2}{\epsilon} \tag{86}$$

$$= \frac{\gamma^3 L^2 \beta_1^2}{2\epsilon(1-\beta_1)^2}(\frac{1}{\lambda_1^2} + \frac{1}{\lambda_2^2} + \frac{1}{\lambda_3^2})\|m_{t-1} \odot \eta_t - 1\|^2 + \frac{\gamma\lambda_1^2}{2}\|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2$$

$$+ \frac{\gamma\lambda_2^2 L^2 \rho_t^2}{2\epsilon} + \frac{2\lambda_3^2 \gamma L^2 \rho_t^2}{\epsilon} \tag{87}$$

$$\leq \frac{\gamma^3 L^2 \beta_1^2}{2\epsilon(1-\beta_1)^2}(\frac{1}{\lambda_1^2} + \frac{1}{\lambda_2^2} + \frac{1}{\lambda_3^2})\frac{dG_\infty^2}{\epsilon^2} + \frac{\gamma\lambda_1^2}{2}\|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2 + \frac{\gamma L^2 \rho_t^2}{2\epsilon}(\lambda_2^2 + 4\lambda_3^2). \tag{88}$$

$$\square$$

**Lemma 10.** *For the term defined in the algorithm, we have*

$$\frac{L}{2}\mathbb{E}\|z_{t+1} - z_t\|^2 \leq \frac{LG^2\gamma^2\beta_1^2}{(1-\beta_1)^2}\mathbb{E}\|\eta_t - \eta_{t-1}\|^2$$

$$+ \gamma_t^2 L(3\frac{1+\beta}{\beta\epsilon}(\mathbb{E}\|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2 + \frac{L\rho_t^2}{\epsilon} + \frac{\sigma^2}{b\epsilon}) + (1+\beta)G^2\mathbb{E}\|\eta_t - \eta_{t-1}\|^2) \tag{89}$$

*Proof.*

$$\frac{L}{2}\mathbb{E}\|z_{t+1} - z_t\|^2$$

$$= \frac{L}{2}\mathbb{E}\|\frac{\gamma\beta_1}{1-\beta_1}m_{t-1} \odot (\eta_t - \eta_{t-1}) - \gamma g_t \odot \eta_t\|^2 \tag{90}$$

$$\leq \frac{L\gamma^2\beta_1^2}{(1-\beta_1)^2}\mathbb{E}\|m_{t-1} \odot (\eta_t - \eta_{t-1})\|^2 + L\mathbb{E}\|\frac{\gamma_t}{b}\sum(\nabla f_i(x_t + \rho_t \frac{s_t}{\|s_t\|})) \odot \eta_t\|^2 \tag{91}$$

$$\leq \frac{LG^2\gamma^2\beta_1^2}{(1-\beta_1)^2}\mathbb{E}\|\eta_t - \eta_{t-1}\|^2 + L\mathbb{E}\|\frac{\gamma_t}{b}\sum(\nabla f_i(x_t + \rho_t \frac{s_t}{\|s_t\|})) \odot \eta_t\|^2 \tag{92}$$

$$= \gamma_t^2 L\mathbb{E}\|\frac{1}{b}\sum(\nabla f_i(x_t + \rho_t \frac{s_t}{\|s_t\|})) \odot \eta_{t-1} + \frac{1}{b}\sum(\nabla f_i(x_t + \rho_t \frac{s_t}{\|s_t\|})) \odot (\eta_t - \eta_{t-1})\|^2$$

$$+ \frac{LG^2\gamma^2\beta_1^2}{(1-\beta_1)^2}\mathbb{E}\|\eta_t - \eta_{t-1}\|^2 \tag{93}$$

$$\leq \frac{LG^2\gamma^2\beta_1^2}{(1-\beta_1)^2}\mathbb{E}\|\eta_t - \eta_{t-1}\|^2 + \gamma_t^2 L((1+\frac{1}{\beta})\mathbb{E}\|\frac{1}{b}\sum(\nabla f_i(x_t + \rho_t \frac{s_t}{\|s_t\|})) \odot \eta_{t-1}\|^2$$

$$+ (1+\beta)\mathbb{E}\|\frac{1}{b}\sum(\nabla f_i(x_t + \rho_t \frac{s_t}{\|s_t\|})) \odot (\eta_t - \eta_{t-1})\|^2) \tag{94}$$

$$\leq \gamma_t^2 L((1+\frac{1}{\beta})\mathbb{E}\|\frac{1}{b}\sum(\nabla f_i(x_t + \rho_t \frac{s_t}{\|s_t\|})) \odot \eta_{t-1}\|^2 + (1+\beta)G^2\mathbb{E}\|\eta_t - \eta_{t-1}\|^2)$$

$$+ \frac{LG^2\gamma^2\beta_1^2}{(1-\beta_1)^2}\mathbb{E}\|\eta_t - \eta_{t-1}\|^2 \tag{95}$$

$$\leq \gamma_t^2 L((1+\frac{1}{\beta})\mathbb{E}\|\frac{1}{b}\sum(\nabla f_i(x_t + \rho_t \frac{s_t}{\|s_t\|})) \odot \sqrt{\eta_{t-1}}\|^2 \times \|\sqrt{\eta_{t-1}}\|_\infty^2$$

$$+ (1+\beta)G^2\mathbb{E}\|\eta_t - \eta_{t-1}\|^2) + \frac{LG^2\gamma^2\beta_1^2}{(1-\beta_1)^2}\mathbb{E}\|\eta_t - \eta_{t-1}\|^2 \tag{96}$$

$$\leq \gamma_t^2 L(\frac{1+\beta}{\beta\epsilon}\mathbb{E}\|\frac{1}{b}\sum(\nabla f_i(x_t + \rho_t \frac{s_t}{\|s_t\|})) \odot \sqrt{\eta_{t-1}}\|^2 + (1+\beta)G^2\mathbb{E}\|\eta_t - \eta_{t-1}\|^2)$$

$$+ \frac{LG^2\gamma^2\beta_1^2}{(1-\beta_1)^2}\mathbb{E}\|\eta_t - \eta_{t-1}\|^2 \tag{97}$$

$$\leq \gamma_t^2 L(3\frac{1+\beta}{\beta\epsilon}\mathbb{E}(\|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2 + \|(\frac{1}{b}\sum \nabla f_i(x_t) - \nabla f(x_t)) \odot \sqrt{\eta_{t-1}}\|^2$$

$$+ \|\frac{1}{b}\sum(\nabla f_i(x_t + \rho_t\frac{\sum_{i\in B}\nabla f_i(x_t)}{\|\sum_{i\in B}\nabla f_i(x_t)\|}) - \nabla f_i(x_t)) \odot \sqrt{\eta_{t-1}}\|^2) + (1+\beta)G^2\mathbb{E}\|\eta_t - \eta_{t-1}\|^2)$$

$$+ \frac{LG^2\gamma^2\beta_1^2}{(1-\beta_1)^2}\mathbb{E}\|\eta_t - \eta_{t-1}\|^2 \tag{98}$$

$$\leq \gamma_t^2 L(3\frac{1+\beta}{\beta\epsilon}(\mathbb{E}\|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2 + \mathbb{E}\|\frac{1}{b}\sum(\nabla f_i(x_t + \rho_t\frac{\sum_{i\in B}\nabla f_i(x_t)}{\|\sum_{i\in B}\nabla f_i(x_t)\|}) - \nabla f_i(x_t)) \odot \sqrt{\eta_{t-1}}\|^2$$

$$+ \frac{\sigma^2}{b\epsilon}) + (1+\beta)G^2\mathbb{E}\|\eta_t - \eta_{t-1}\|^2) + \frac{LG^2\gamma^2\beta_1^2}{(1-\beta_1)^2}\mathbb{E}\|\eta_t - \eta_{t-1}\|^2 \tag{99}$$

$$\leq \gamma_t^2 L(3\frac{1+\beta}{\beta\epsilon}(\mathbb{E}\|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2 + \frac{1}{\epsilon}\mathbb{E}\|\frac{1}{b}\sum(\nabla f_i(x_t + \rho_t\frac{\sum_{i\in B}\nabla f_i(x_t)}{\|\sum_{i\in B}\nabla f_i(x_t)\|}) - \nabla f_i(x_t))\|^2$$

$$+ \frac{\sigma^2}{b\epsilon}) + (1+\beta)G^2\mathbb{E}\|\eta_t - \eta_{t-1}\|^2) + \frac{LG^2\gamma^2\beta_1^2}{(1-\beta_1)^2}\mathbb{E}\|\eta_t - \eta_{t-1}\|^2 \tag{100}$$

$$\leq \gamma_t^2 L(3\frac{1+\beta}{\beta\epsilon}(\mathbb{E}\|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2 + \frac{1}{\epsilon b}\mathbb{E}\sum\|\nabla f_i(x_t + \rho_t\frac{\sum_{i\in B}\nabla f_i(x_t)}{\|\sum_{i\in B}\nabla f_i(x_t)\|}) - \nabla f_i(x_t)\|^2$$

$$+ \frac{\sigma^2}{b\epsilon}) + (1+\beta)G^2\mathbb{E}\|\eta_t - \eta_{t-1}\|^2) + \frac{LG^2\gamma^2\beta_1^2}{(1-\beta_1)^2}\mathbb{E}\|\eta_t - \eta_{t-1}\|^2 \tag{101}$$

$$\leq \gamma_t^2 L(3\frac{1+\beta}{\beta\epsilon}(\mathbb{E}\|\nabla f(x_t) \odot \sqrt{\eta_{t-1}}\|^2 + \frac{L\rho_t^2}{\epsilon} + \frac{\sigma^2}{b\epsilon}) + (1+\beta)G^2\mathbb{E}\|\eta_t - \eta_{t-1}\|^2)$$

$$+ \frac{LG^2\gamma^2\beta_1^2}{(1-\beta_1)^2}\mathbb{E}\|\eta_t - \eta_{t-1}\|^2. \tag{102}$$

□

# APPENDIX C
## ADDITIONAL EXPERIMENT ILLUSTRATIONS

### A. Experiment Illustrations

In the ablation study, we conduct the experiments on the GLUE benchmark with AdaSAM, AMSGrad, SAM and SGD, respectively. The optimizers do not have the momentum part ($\beta_1 = 0$). As a supplement to Table II, Figure 3 show the detailed loss and evaluation metrics versus number of steps curves during training. The loss curve of AdaSAM decreases faster than SAM and SGD in all tasks, and it has a similar decreasing speed as the AMSGrad. The metric curve of AdaSAM and AMSGrad show that the adaptive learning rate method is better than SGD and SAM. And AdaSAM decrease as faster as the AMSGrad in all tasks.
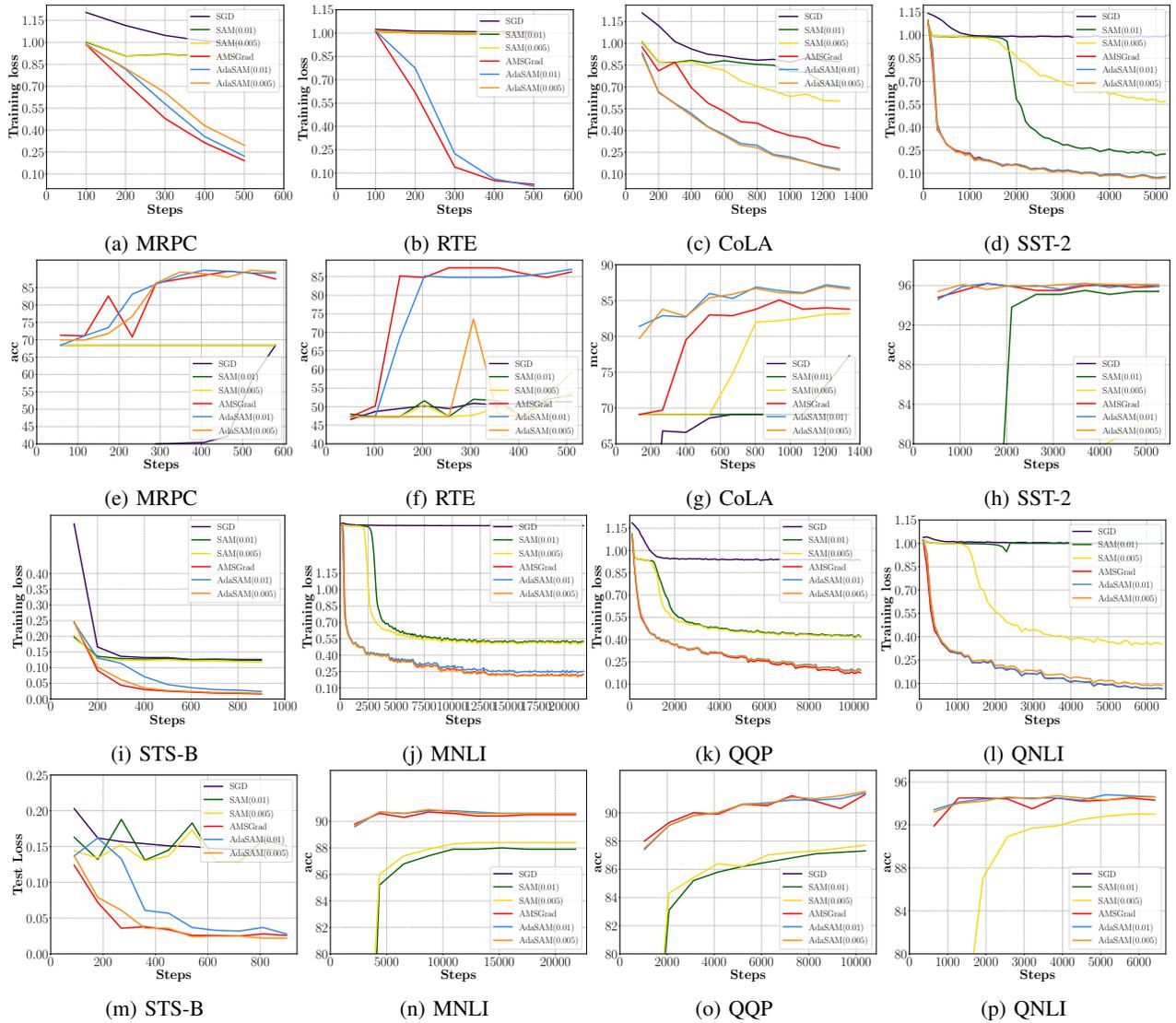
Fig. 3: The loss and evaluation metric v.s. steps on MRPC, RTE, CoLA, SST-2, STS-B, MNLI, QQP and QNLI.($\beta_1 = 0$)