SOFTSIGNSGD(S3): AN ENHANCED OPTIMIZER FOR PRACTICAL DNN TRAINING AND LOSS SPIKES MINI MIZATION BEYOND ADAM

Anonymous authors

Paper under double-blind review

ABSTRACT

Adam has been widely successful in training deep neural networks (DNNs), yet the factors contributing to both its practical effectiveness and ineffectiveness remain largely underexplored. In this study, we reveal that the effectiveness of Adam in training complicated DNNs stems primarily from its similarity to SignSGD in managing significant gradient variations, while we also theoretically and empirically uncover that Adam is susceptible to loss spikes due to potential excessively large updates. Building on these insights, we propose a novel optimizer, SignSoft-SGD (S3), which incorporates a generalized sign-like formulation with a flexible p-th order $(p \ge 1)$ momentum in the denominator of the update, replacing the fixed 2-order momentum. We also integrate the memory-efficient Nesterov's accelerated gradient technique into S3, enhancing convergence speed without additional memory overhead. To minimize the risk of loss spikes, we utilize the same coefficient for the momentums in both the numerator and the denominator of the update, which also practically streamlines the tuning overhead. We conduct a theoretical analysis of S3 on a general nonconvex stochastic problem, demonstrating that S3 achieves the optimal convergence rate under weak assumptions. Extensive experimentation across various vision and language tasks demonstrates that S3 not only achieves rapid convergence and improved performance but also rarely encounters loss spikes even at a $10 \times$ larger learning rate. Specifically, S3 delivers performance comparable to or better than AdamW with $2 \times$ the training steps.

031 032

033

006

007

008 009 010

011

013

014

015

016

017

018

019

021

024

025

026

027

028

029

1 INTRODUCTION

Optimizers play a pivotal role in training DNNs.Currently, Adam (Kingma & Ba (2015)) stands out as the dominant optimizer for training Transformers (Vaswani et al. (2017)), especially for the recent phenomenal large language models (LLMs) (Brown et al. (2020); Chowdhery et al. (2023); Touvron et al. (2023)), and large vision models (Radford et al. (2021); Kirillov et al. (2023)). Even in the realm of training the modern convolutional neural networks (CNNs), such as ConvNeXt (Liu et al. (2022); Woo et al. (2023)), Adam also has become the de facto optimizer, although stochastic gradient descent (SGD) was traditionally deemed more suitable for training CNNs (Krizhevsky et al. (2017); He et al. (2016)).

While the practical success of Adam is indisputable, the underlying reasons for its effectiveness
remain poorly understood. The original paper on Adam attributes its success to the effective
coordinate-wise adaptivity (Kingma & Ba (2015)). However, recent work (Chen et al. (2023b))
challenges this perspective by proposing an optimizer that achieves comparable, and sometimes
superior, performance to Adam in training various architectures without leveraging adaptivity.

We first revisit Adam to discern the reasons behind its practical effectiveness. Each coordinate of the update of Adam (*i.e.*, $\frac{m_t}{\sqrt{v_t}}$) exhibits a sign-like characteristic. Empirical evidence from (Kunstner et al. (2023)) demonstrates that the simple sign descent can substantially narrow the performance gap between SGD and Adam in training complicated DNNs, such as Transformers. *This suggests that the sign-like property of Adam is a key factor in its effectiveness*. However, (Kunstner et al. (2023)) uncovered this phenomenon but did not explore why sign descent is effective in training Transformers. This paper explains that wide difference in inter-layer and intra-layer gradients during training is the behind reason, and the effectiveness of Adam is mainly attributed to its conservative



¹This is the reason why we named the new optimizer SoftSignSGD.

108

110

111

112

113

114

115

116 117

118 119

120

121 122 • We conduct a theoretical analysis for S3 on a general nonconvex stochastic problem, achieving the optimal convergence rate under a weak non-uniform smoothness assumption. (Section 4)

• We conduct extensive experiments to evaluate S3 against Adam and other related optimizers. The experimental results demonstrate the faster training speed and superior inference performance of S3. Specifically, Specifically, S3 achieves improvements comparable to or exceeding those of AdamW with twice the training steps, while rarely experiencing loss spikes even at significantly higher learning rates. (Section 5)

2 RETHINKING THE EFFECTIVENESS AND INEFFECTIVENESS OF ADAM

In a deep learning task, the optimizer aims to minimize the empirical risk loss of a model on a dataset, *i.e.*,

$$\min_{\boldsymbol{x}\in\mathbb{R}^d} F(\boldsymbol{x}) = \mathbb{E}_{\boldsymbol{\zeta}\sim\mathcal{D}}[f(\boldsymbol{x};\boldsymbol{\zeta})] = \frac{1}{n} \sum_{i=1}^n f(\boldsymbol{x};\omega_i),\tag{1}$$

(2)

where x is the *d*-dimensional model parameter, and ζ is independently and identically sampled from the dataset { $\omega_i : \omega_i \in \mathcal{D}, 1 \le i \le n$ }.

Nowadays, Adam has emerged as the dominant optimizer for training DNNs. It significantly outperforms SGD in training the increasingly popular Transformers, demonstrating remarkable efficacy. Even for CNN-based models like ConvNeXT, Adam is preferred for achieving superior performance, despite the historical consideration that SGD is more suitable for training CNNs. While the practical success of Adam is indisputable, the factors contributing to its practical effectiveness remain largely underexplored. There is a pressing need to delve into the effectiveness of Adam to facilitate significant advancements in DNN training.

Recalling the updating rule of Adam, we have

$$\tilde{m}_t = \beta_1 \tilde{m}_{t-1} + (1 - \beta_1) g_t,$$

 135
 $m_t = \frac{\tilde{m}_t}{1 - \beta_1^t},$

 136
 $\tilde{v}_t = \beta_2 \tilde{v}_{t-1} + (1 - \beta_2) g_t^2,$

 137
 $\tilde{v}_t = \beta_2 \tilde{v}_{t-1} + (1 - \beta_2) g_t^2,$

 138
 $v_t = \frac{\tilde{v}_t}{1 - \beta_2^t},$

 139
 $x_{t+1} = x_t - \gamma_t \frac{m_t}{\sqrt{v_t}},$

where x_t denotes the model parameter, $g_t = \nabla f(x_t; \zeta_t)$ is the stochastic gradient, γ_t is the learning rate, and β_1 and β_2 represents the exponential moving average coefficients.

144 145 146 147 146 147 In essence, $|\boldsymbol{m}_t|$ and $\sqrt{\boldsymbol{v}_t}$ are of the same order of magnitude. Specifically, if g_t ideally stays stable over a period, Adam in Eq. (2) reduces to SignSGD, *i.e.*. $\boldsymbol{x}_{t+1} = \boldsymbol{x}_t - \gamma_t \frac{\boldsymbol{m}_t}{\sqrt{\boldsymbol{v}_t}} = \boldsymbol{x}_t - \gamma_t \text{Sign}(\boldsymbol{g}_t)$. Therefore, Adam can be viewed as a sign-like optimizer.

The primary reason for Adam's practical effectiveness over SGD in training complex DNNs remains 148 fragmented across prior studies and lacks consolidation. Kunstner et al. (2023) empirically shows 149 that sign descent with momentum yields comparable performance to Adam when training Trans-150 formers, albeit lacking comprehensive analytical justification. More recently, Chen et al. (2023b) 151 employs an AutoML method to discover an effective optimizer, Lion, resembling SignSGD with 152 momentum, and showcases superior performance to Adam across diverse DNN models. Indeed, 153 the effectiveness of both Adam and Lion primarily stems from their shared sign-like property. For 154 deep networks, the gradients of the initial and final layers can differ significantly, as theoretically 155 verified in (Yang et al. (2019); Liu et al. (2020); Xiong et al. (2020); Kim et al. (2021); Qi et al. 156 (2023)) through the mean-field theory and from the perspective of Lipschitz continuity. Further-157 more, even within the same layer of a Transformer, gradients can vary greatly due to the presence 158 of the attention component (Noci et al. (2022)). An illustrative example can be found in Section B.4 159 of the appendix. This substantial gradient discrepancy poses challenges for SGD, which, by directly employing gradients as updates, necessitates a relatively small learning rate to prevent divergence, 160 resulting in noticeable training slowdown. Moreover, another drawback of using SGD is that pa-161 rameters with large gradients undergo substantial changes, while those with small gradients tend to remain close to their initial values. This discrepancy weakens the overall representation ability of networks, thereby degrading final performance. In contrast, Adam's updates remain close to ± 1 despite significant gradient gaps, thanks to its inherent sign-like property. This characteristic renders Adam a conservative yet effective choice for training complex DNNs. In summary, Adam's efficacy in training complex DNNs stems from its conservative sign-like descent, which effectively addresses significant gradient discrepancies.

While Adam effectively trains complex DNNs, it also escalates the risk of training instability and loss spikes with non-trivial probability. This can be inferred from Theorem 1.

Theorem 1. The sequences $\{m_t\}$ and $\{v_t\}$ are generated by Adam in Eq. (2). If the moving average coefficients satisfy $\beta_1^2 < \beta_2$, then it holds that

173 174 175

176

177 178

$$\frac{|\boldsymbol{m}_{t}^{(j)}|}{\sqrt{\boldsymbol{v}_{t}^{(j)}}} \leq \frac{(1-\beta_{1})\sqrt{1-\beta_{2}^{t}}\sqrt{1-(\frac{\beta_{1}^{t}}{\beta_{2}})^{t}}}{(1-\beta_{1}^{t})\sqrt{1-\beta_{2}}\sqrt{1-\frac{\beta_{1}^{2}}{\beta_{2}}}} \simeq \frac{1-\beta_{1}}{\sqrt{1-\beta_{2}}\sqrt{1-\frac{\beta_{1}^{2}}{\beta_{2}}}},$$
(3)

where $\frac{|\boldsymbol{m}_{t}^{(j)}|}{\sqrt{\boldsymbol{v}_{t}^{(j)}}}$ reach to the largest value if the signs of $\{\boldsymbol{g}_{t}^{(j)}, \boldsymbol{g}_{t-1}^{(j)}, ..., \boldsymbol{g}_{t-k}^{(j)} ...\}$ are the same and $|\boldsymbol{g}_{t}^{(j)}| = \frac{\beta_{2}|\boldsymbol{g}_{t-1}^{(j)}|}{\beta_{1}} = \frac{\beta_{2}^{2}|\boldsymbol{g}_{t-2}^{(j)}|}{\beta_{1}^{2}} = ... = \frac{\beta_{2}^{2}|\boldsymbol{g}_{t-k}^{(j)}|}{\beta_{t}^{k}} ...^{2}.$

179 180 181

182

183

185

191

192



Figure 2: Visualization of the mean update (*i.e.*, $\operatorname{Avg}(|\mathbf{m}_t^{(j)}|/\sqrt{\mathbf{v}_t^{(j)}})$), the maximum update (*i.e.*, $\max_{j \in [d]}(|\mathbf{m}_t^{(j)}|/\sqrt{\mathbf{v}_t^{(j)}})$), and the training loss over 50,000 iterations during GPT-2 (345M) training on Open-WebText using AdamW ($\beta_1 = 0.9, \beta_2 = 0.999$) with a cosine learning rate schedule. The figure illustrates that all loss spikes are preceded by abrupt increases in the mean update for any coordinate can lead to a significant rise in the mean update, which then triggers loss spikes. Moreover, these spikes primarily occur during the early training phase when the learning rate is relatively high. In later stages, as the learning rate decreases, loss spikes become infrequent, even with large maximum updates, like around Iteration 40,000.

200 Theorem 1 indicates that when Adam is employed, there exists a probability that the update of 201 each element $|\mathbf{m}_{t}^{(j)}|/\sqrt{\mathbf{v}_{t}^{(j)}}$ can reach an excessively large value. For instance, with typical settings 202 of $\beta_1 = 0.9$ and $\beta_2 = 0.999$, the update $|\boldsymbol{m}_t^{(j)}|/\sqrt{\boldsymbol{v}_t^{(j)}}$ could approach its theoretical maximum of 203 $1 - \beta_1 / \sqrt{1 - \beta_2} \sqrt{1 - \frac{\beta_1^2}{\beta_2}} \simeq 7.27$, while the normal absolute value of the update is less than 1. While the 204 probability of any specific parameter's update reaching this maximal value is low, the probability that 205 at least one parameter's update reaches this value is high due to the large number of parameters in 206 large models (e.g., LLM). When a parameter's update is excessively large and the learning rate is also 207 high, it is likely to deviate substantially from its intended trajectory. Such deviations may propagate 208 to neighboring parameters through interconnections, triggering a chain reaction that culminates in 209 loss spikes. Supporting examples are clearly illustrated in Figure 1 and Figure 2. This mechanism 210 explains the frequent occurrence of loss spikes during LLM training, particularly in the initial stages 211 when learning rates are higher. Specifically, the probability of loss spikes increases with the size of 212 the LLM. In conclusion, vanilla Adam poses a significant risk of loss spikes during large-scale 213

²¹⁶ ²The update $m_t^{(j)}/\sqrt{v_t^{(j)}}$ with respect to $\{g_k^{(j)}\}_{k=1}^t$ is a continuous function. Thus, when most of the signs 215 of $\{g_k^{(j)}\}_{k=1}^t$ are consistent and the secondary condition is nearly satisfied, $m_t^{(j)}/\sqrt{v_t^{(j)}}$ will be close to the theoretical maximum.

model training. Mitigating this problem requires strategies to constrain the maximum update magnitude for each parameter coordinate.

Remark 1 [Loss spikes during LLM Training]. Encountering loss spikes is a common phe-219 nomenon during LLM training (Zeng et al. (2022); Chowdhery et al. (2023); Touvron et al. (2023); 220 Yang et al. (2023)). However, the underlying reasons for this problem were not well explored prior to 221 this. Practitioners had to resort to ad hoc engineering strategies such as skipping some data batches 222 before the spike occurs and restarting training from a nearby checkpoint (Chowdhery et al. (2023); Molybog et al. (2023)), resulting in resource wastage due to frequent rollbacks and checkpointing 224 savings. Some previous works investigated the phenomenon of train instability (Liu et al. (2019)) 225 and loss spikes (Zhu et al. (2023); Zhang & Xu (2023)). (Liu et al. (2019)) demonstrated that the 226 variance of the update $1/\sqrt{v_t^{(j)}}$ is significantly larger, often causing the update to become dispropor-227 tionately large, but the analysis only works in the early stage. The analyses in (Zhu et al. (2023); 228 Zhang & Xu (2023)) were restricted to either linear models or shallow networks with mean squared error (MSE) loss, using (S)GD as the optimizer. Consequently, it is questionable whether these find-229 ings can be directly applied to the context of LLM training. More recently, (Molybog et al. (2023)) 230 suggested that time-domain correlation between gradient estimates of earlier layers contributes to 231 loss spikes during LLM training. The suggested mitigation strategies include lowering the ϵ value in 232 Adam and reducing the batch size. However, the study itself that these methods are not silver bullets 233 for a fundamental solution. To the best of our knowledge, the analyses presented above provide 234 the first formal explanation for the frequent occurrence of loss spikes during LLM training. 235

²³⁶ 3 S3 Algorithm

Analyzing Adam yields insights guiding the construction of a more effective optimizer for training DNNs: 1) The update of the optimizer need only resemble the sign of the gradient, without strictly adhering to the formulation involving the ratio of first-order gradient momentum to second-order gradient momentum. 2) Minimizing the largest value of the update in the optimizer is crucial to mitigate potential loss spikes.

Recently, several studies (Dozat (2016); Xie et al. (2024); Zhou et al. (2023)) introduced the NAG
technique to DNN optimizers, consistently demonstrating faster training convergence and improved
inference performance across a broad spectrum of DNNs compared to the standard Adam. Therefore,
integrating the NAG technique into optimizers is highly worthwhile.

- Given the observations above, we propose a new optimizer, referred to as SoftSignSGD (S3). The detailed implementation of S3 is illustrated in Algorithm 1.
- 249 Key characteristics of S3 are summarized below:

250 First, S3 features a more general sign-like formulation with a flexible p-order momentum, ex-251 tending beyond the fixed 2-order momentum suggested by the original motivation of Adam. 252 According to Theorem 2, a large *p*-order momentum enables the use of a larger learning rate, pro-253 moting faster convergence and better performance (Kong & Tao (2020)). Moreover, during training, 254 abrupt changes occasionally occur in some coordinates of the gradients, potentially leading to train-255 ing instability and even divergence without a proper remedy. In such cases, the gradients become 256 more heterogeneous over time. As indicated in Theorem 2, the update $\frac{|n_t|}{b_t(p)}$ of S3 becomes smaller 257 with a large *p*-order momentum, providing a counteractive effect to stabilize the training process. 258 Figure 1 illustrates this behavior. Additionally, the computational cost of optimization becomes non-259 trivial when training LLMs. Setting p = 1 for S3 involves only a computationally light element-wise 260 absolute operation, reducing overall computational overhead.

Second, S3 shares the same exponential moving average coefficient β for both m_t and r_t , offering the advantages of minimizing the risk of loss spikes and reducing tuning work. In theory, as demonstrated in Theorem 2, the same β guarantees that the largest value of each coordinate of the update $\frac{|n_t|}{b_t(p)}$ is minimized to 1. As discussed in Section 2, this design helps mitigate the lossspike problem. In practice, this design reduces tuning effort by removing one hyperparameter and lowers computational costs by eliminating the bias correction required by Adam.

Third, S3 introduces the NAG technique to further accelerate training without incurring extra
 memory costs. While previous works such as NAdam and Adan have also utilized the NAG technique in their adaptive optimizers, there are significant differences. In S3, NAG is applied to both

270	Algorithm 1. SoftSianSGD (S3)
271	1: Input : the momentum $m_0 = 0$, $s_0 = 0$ the exponential moving average coefficient
272	β within [0, 1], the power factor p within $[1, +\infty)$, and the learning rate sequence $\{\gamma_t\}$.
273	2: for $t = 1,, T$ do
274	3: Randomly sample data and compute the gradient: $g_t \leftarrow \nabla F(x_t; \zeta_t)$
275	4: Update the momentum \boldsymbol{m}_t : $\boldsymbol{m}_t \leftarrow \beta \boldsymbol{m}_{t-1} + (1-\beta)\boldsymbol{g}_t$
276	5: Update the momentum $s_t(p)$: $s_t(p) \leftarrow \beta s_{t-1} + (1-\beta) g_t ^p$
277	6: Compute the Nesterov momentum $n_t: n_t \leftarrow \beta m_t + (1 - \beta)g_t$
278	7: Compute the Nesterov momentum $\boldsymbol{b}_t(p)$: $\boldsymbol{b}_t(p) \leftarrow (\beta \boldsymbol{s}_t(p) + (1-\beta) \boldsymbol{g}_t ^p)^{1/p}$
279	8: Update the model parameter: $x_{t+1} \leftarrow x_t - \gamma_t \frac{n_t}{b_t(p)}$
280	9: end for
281	

the numerator and denominator of the update. The Nesterov momentum estimators in S3 follow the NAG (II) formulation, shown to be equivalent to vanilla NAG (I) in Theorem 7. The key ad-vantage of this formulation is that S3 avoids additional memory usage. Conversely, NAdam (Dozat (2016)) only incorporates the Nesterov momentum in the numerator of the update and relies on complex bias-correction operations to stabilize training. Adan (Xie et al. (2024)) also employs Nesterov momentum estimators in both the numerator and the denominator of the update. However, its formulations, akin to NAG (III), demand more memory for storing the previous gradient g_{t-1} and the new momentum r_k compared to Adam. Consequently, Adam may not be ideal for training LLMs due to its memory demands. Furthermore, it introduces a new momentum coefficient, increasing the need for tuning.

Theorem 2. The sequences $\{n_t\}$ and $b_t(p)$ are generated S3 in Algorithm 1. If the moving average coefficients for m_t , n_t and s_t , b_t of $ar\beta_1$ and β_2 which satisfy $\beta_1 < \beta_2^{1/p}$ and $p \ge 1$, it holds that

(1). The upper bound of each element of the update $\frac{n_t^{(j)}}{\mathbf{b}_{\cdot}^{(j)}}$ is

$$\frac{|\boldsymbol{n}_t^{(j)}|}{\boldsymbol{b}_t^{(j)}(p)} \le \frac{(1-\beta_1)}{(1-\beta_2)^{1/p} \left(1-\frac{\beta_1^q}{\beta_2^{q/p}}\right)^{1/q}},\tag{4}$$

where $\frac{1}{p} + \frac{1}{q} = 1$.

(2). When $\beta_1 = \beta_2$, the upper bound of each element of the update $\frac{\mathbf{n}_t^{(j)}}{\mathbf{b}_t^{(j)}}$ reaches to the smallest 1, $|\boldsymbol{n}_t^{(j)}| \leq 1$

i.e.,
$$\frac{1}{\boldsymbol{b}_t^{(j)}(p)} \leq 1$$

(3). Let $1 \le p_1 \le p_2$, and then $b_t(p_1) \le b_t(p_2)$.

Theorem 3. The three formulations of NAG are listed in the following. Let $x_t = \tilde{x}_t - \gamma \beta m_{t-1}$, the three formulations are equivalent, i.e.,

NAG (I):
$$\begin{cases} \boldsymbol{g}_t = \nabla f(\tilde{\boldsymbol{x}}_t - \gamma \beta \boldsymbol{m}_{t-1}; \zeta_t) \\ \boldsymbol{m}_t = \beta \boldsymbol{m}_{t-1} + (1-\beta) \boldsymbol{g}_t \\ \tilde{\boldsymbol{x}}_{t+1} = \tilde{\boldsymbol{x}}_t - \gamma \boldsymbol{m}_t \end{cases}$$
(5)

NAG (II):
$$\begin{cases} \boldsymbol{g}_t = \nabla f(\boldsymbol{x}_t; \zeta_t) \\ \boldsymbol{m}_t = \beta \boldsymbol{m}_{t-1} + (1-\beta) \boldsymbol{g}_t \\ \boldsymbol{x}_{t+1} = \boldsymbol{x}_t - \gamma (\beta \boldsymbol{m}_t + (1-\beta) \boldsymbol{g}_t) \end{cases},$$
(6)

(7)

$$\int oldsymbol{g}_t =
abla f(oldsymbol{x}_t; \zeta_t)$$

NAG (III): $\begin{cases} \boldsymbol{m}_{t} = \beta \boldsymbol{m}_{t-1} + (1-\beta)\boldsymbol{g}_{t} \\ \boldsymbol{r}_{t} = \beta \boldsymbol{r}_{t-1} + (1-\beta)(\boldsymbol{g}_{t} - \boldsymbol{g}_{t-1}) \\ \boldsymbol{x}_{t+1} = \boldsymbol{x}_{t} - \gamma(\boldsymbol{m}_{t} + \beta \boldsymbol{r}_{t}) \end{cases}$

Moreover, if $\tilde{x}_{t+1} \to \tilde{x}_t$ as $m_t \to 0$, x_t will converge to \tilde{x}_t .

4 THEORETICAL CONVERGENCE ANALYSIS 325

To present the theoretical convergence guarantee for S3 (Algorithm 1) to optimize the nonconvex problem in Eq. (1), we first introduce some necessary assumptions.

Assumption 1. [Bounded infimum] There exists a constant F^* , the objective function follows $F(\mathbf{x}) \geq F^*$ for any $\mathbf{x} \in \mathbb{R}^d$.

Assumption 2. [Generalized Smoothness] *There exist constants* $L_0, L_1, R \ge 0$, for any $x, y \in \mathbb{R}^d$ with $||x - y||_2 \le R$, the objective function follows,

$$\|\nabla F(\boldsymbol{y}) - \nabla F(\boldsymbol{x})\|_{2} \le (L_{0} + L_{1} \|\nabla F(\boldsymbol{x})\|_{2}) \|\boldsymbol{x} - \boldsymbol{y}\|_{2}.$$
(8)

Assumption 3. [Unbias noisy gradient and bounded variance] There exists a constant σ . For $x_t \in \mathbb{R}^d$ at any time, the noisy gradient of the objective function obeys follows

$$\mathbb{E}[\boldsymbol{g}_t] = \mathbb{E}[\nabla f(\boldsymbol{x}_t; \zeta_t)] = \nabla F(\boldsymbol{x}_t), \qquad \mathbb{E}[\|\boldsymbol{g}_t - \nabla F(\boldsymbol{x}_t)\|_2^2] \le \sigma^2.$$
(9)

Under the assumptions above, we then present the theoretical convergence for S3 in Theorem 4.

Theorem 4. $\{x_t\}_{t=1}^T$ is generated by Algorithm 1 under Assumption 1-4. Let the hyperparameters be set as $\beta = 1 - \frac{1}{\sqrt{T}}$ and $\gamma = \frac{1}{L_0 T^{3/4}}$. If $u_t = \frac{|\boldsymbol{n}_t^{(j)}|}{\boldsymbol{b}_s^{(j)}} \ge \frac{1}{U_{\text{max}}}$, then

$$\frac{1}{T} \sum_{t=1}^{T} \mathbb{E}[\|\nabla F(\boldsymbol{x}_{t})\|_{1}] \leq \frac{2L_{0}U_{\max}(F(\boldsymbol{x}_{1}) - F(\boldsymbol{x}^{*}))}{T^{1/4}} + \frac{4\beta U_{\max}\sqrt{d}\mathbb{E}\left[\|\nabla F(\boldsymbol{x}_{1})\|_{2}\right]}{T^{1/2}} + \frac{4U_{\max}\sqrt{d}\sigma}{T^{1/4}} + \frac{4\beta^{2}U_{\max}d}{T^{3/4}} + \frac{U_{\max}d}{T^{3/4}}.$$
(10)

347 348 349

328

330

331

332 333 334

335

340 341

342

350 Remark 2 [Adopting Weaker Assumption]. The theoretical convergence analysis for S3 in The-351 orem 4 requires only a general non-uniform smoothness condition (Assumption 2). In contrast, 352 previous works that developed convergence analyses for Adam required stronger assumptions or 353 achieved weaker conclusions. (Chen et al. (2018); Défossez et al. (2020)) proved convergence for 354 non-convex objectives under the assumption that gradients are bounded. (De et al. (2018)) required that the signs of gradients remain the same along the trajectory, despite considering Nesterov accel-355 eration. (Zhang et al. (2022)) assumed uniform L-smoothness but only proved convergence to some 356 neighborhood of stationary points with a constant radius. Very recently, (Li et al. (2023); Hong 357 & Lin (2024)) offered convergence bounds for Adam under the general non-uniform smoothness 358 assumption, but the convergence is in probability. 359

Remark 3 [Achieving Optimal Convergence Rate]. As illustrated in Theorem 4, the convergence order of S3 is $O(\frac{1}{T^{1/4}})$, consistent with the established lower bound for optimal convergence in non-convex stochastic optimization (Arjevani et al. (2023)).

³⁶³ 5 EXPERIMENT

We compare S3 with representative optimizers, including SGDM(Robbins & Monro (1951)), 365 AdamW (Loshchilov & Hutter (2017)), NAdam(Dozat (2016)), Adan (Xie et al. (2024)), and Lion 366 (Chen et al. (2023b)), for both vision and language tasks. For vision tasks, we evaluate the classi-367 fication accuracy by training the CNN-type ResNet-50 (He et al. (2016)) and the Transformer-type 368 ViT-Base (Dosovitskiy et al. (2020)) on ImageNet. In language tasks, we assess the pre-training 369 performance by training GPT-2 (345M) and GPT-2 (7B) (Radford et al. (2019)) on OpenWebText 370 and the refined CommonCrawl. Results on downstream tasks for the pre-trained GPT-2 (345M) and 371 GPT-2 (7B) are provided in the Appendix. 372

373 374

5.1 EXPERIMENTS FOR VISION TASKS

Compared with the baseline AdamW, Figure 3(a) and Figure 3 (c) illustrate that S3 exhibits obvious faster convergence and achieves a significantly smaller final training loss. As shown in Figure 3(b), Figure 3(d), and Table 1, S3 achieves test accuracies that are 1.47% and 1.36% higher for training ResNet-50 and ViT-B, respectively. This represents a significant improvement in training speed and

378 4.5 379 AdamW, epoch=300, Ir=3e-3 AdamW, epoch=150, Ir=3e-3 SGDM, epoch=150, Ir=3e-1 4.2 380 NAdam, epoch=150, Ir=6e-3 AdamW, epoch=300, Ir=3e-3 381 Adan, epoch=150, /r=6e-3 Lion,epoch=150, /r=1e-3 AdamW, epoch=150, Ir=3e SGDM, epoch=150, lr=3e-1 382 Train 3.3 S3, epoch=150, Ir=6e-3 NAdam, epoch=150, Ir=6e-3 Adan, epoch=150, *lr*=6e-3 Lion, epoch=150, *lr*=1e-3 383 384 S3, epoch=150, /r=6e-3 100 200 250 385 Epoch 386 (a) ResNet-50, Train Loss (b) ResNet-50, Test Accuracy 387 388 AdamW, epoch=300, /r=3e-3 AdamW, epoch=150, /r=3e-3 80 389 75 SGDM, epoch=150, /r=3e-1 NAdam, epoch=150, *lr*=6e-3 Adan, epoch=150, *lr*=6e-3 390 damW, epoch=300, Ir= 391 Lion,epoch=150, Ir=1e-3 AdamW, epoch=150, Ir= .u 3.5 est \$3, epoch=150, /r=6e-3 SGDM, epoch=150, Ir=3e-1 Nadam, epoch=150, *lr*=6e 392 3.0 Adan, epoch=150, Ir=6e-3 Lion enorb=150 /r=1e-3 393 2.5 S3, epoch=150, /r=6e-3 250 200 250 Epoch Ep (c) ViT-B/16, Train Loss (d) ViT-B/16, Test Accuracy 396

Figure 3: Comparison of train loss and test accuracy on ImageNet for training ReNet-50 and ViT-B/16 with AdamW, SGDM, NAadm, Adan, Lion and S3.

inferencing accuracy. Even when we increase the training epochs by $2 \times$ for AdamW, S3 remains comparable and even slightly better. Moreover, AdamW experiences loss spikes during the training of ViT-B, while S3 maintains training stability even with a large learning rate.

In addition, other competitive optimizers, including SGDM, Adan, and Lion, are also evaluated and presented in Figure 3 and Table 1. While SGDM performs comparably to AdamW on the CNN-type ResNet-50, it exhibits poor performance on the Transformer-type ViT-B, consistent with the analyses in Section 2. Due to the introduction of the NAG technique, Adan and Lion outperform AdamW in terms of training speed and test accuracy, yet they still fall short compared to S3. It is noteworthy that Adan and Lion also encounter issues of instability during training.

Table 1: Test accuracy (%) on ImageNet for training ResNet-50 and ViT-B/16 with AdamW, SGDM, Adan, Lion and S3.

Naturali		300 epochs					
INCLWOIK	AdamW	SGDM	NAdam	Adan	Lion	S3	AdamW
ResNet-50	77.29	77.50	77.36	78.23	77.14	78.76	78.46
ViT-B/16	79.52	60.99	80.31	80.11	80.32	80.93	80.13

415 416 417

418

410

411

412 413 414

397

399

5.2 EXPERIMENTS FOR LANGUAGE TASKS

419 As illustrated in Figure 4 and Table 2, S3 consistently achieves faster train convergence and lower 420 validation perplexity, compared to AdamW. The superiority becomes more obvious as the model size 421 increases. Importantly, the improvement on the 345M model brought by S3 is comparable to that 422 achieved by AdamW with twice the number of steps. This can translate into a significant reduction in 423 the number of steps and total computation needed to reach the same loss level, providing substantial time and cost savings for LLM pre-training. Moreover, while AdamW frequently experiences loss 424 spikes, S3 rarely encounters this issue, even with a learning rate that is $10 \times$ larger. Additionally, a 425 large p-order momentum for S3 allows a large learning rate, leading to further training acceleration 426 and performance improvement. 427

Moreover, Lion and Adan are also investigated on the GPT-2 (345M) model. Although Lion converges faster than AdamW at the beginning, it does not showcase superiority in the final validation perplexity. Adan slightly outperforms AdamW in train speed and validation performance, but as analyzed in Section 3, Adan requires more memory and hyperparameters to tune, which is not appealing for pre-training LLMs. Additionally, Adan is also prone to experiencing loss spikes.



Figure 4: Comparison of train loss and validation loss for pre-training GPT-2 (345M) and GPT-2 (7B) with AdamW, NAdam, Adan, Lion and S3.

Table 2. Validation perpressity (the lower, the better) for training Of 1-2 (34500) and Of 1-2 (7B).								
		50k steps					100k steps	
Network	Dataset	AdamW	NAdam	Adan	Lion	S3	AdamW	
		(<i>lr</i> =3e-4)	(<i>lr</i> =3e-4)	(<i>lr</i> =1e-3)	(<i>lr</i> =6e-5)	(p=3,lr=3e-3)	(<i>lr</i> =3e-4)	
GPT-2 (345M)	OpenWebText	4.78	4.71	4.69	4.76	4.59	4.57	
GPT-2 (7B)	CommonCrawl	21.13	-	-	-	19.69	-	

5.3 ABLATION STUDY

We implement ablation experiments for training ViT-B/16 to clarify the contributions of each modification of S3 over Adam. Figure 5 and Table 3 showcase that employing a large learning rate and sharing the same β alone have little or even a negative impact on performance (e.g., Exp. ① vs. Exp. ②, Exp. ① vs. Exp. ③), while their combination results in a notable improvement (e.g., Exp. ②, Exp. ③ vs. Exp. ⑥, Exp. ⑧ vs. Exp. ⑨). In contrast, harnessing a larger p can have an individually beneficial effect on performance (e.g., Exp. ① vs. Exp. ⑤, Exp. ③ vs. Exp. ⑧), and the performance gain from the benefits of NAG is more pronounced than other modification (e.g., Exp. ① vs. Exp. ④, Exp. ⑥ vs. Exp. ⑦, Exp. ⑨ vs. Exp. ⑩).

Table 3:	Ablation stud	y on test accuracie	s (%) of S3	for training	ViT-B/16.
----------	---------------	---------------------	-------------	--------------	-----------

	Exp.	large <i>lr</i>	same β	NAG	flexible p	Test Accuracy
[1	-	-	-	-	79.52 (AdamW,lr=3e-3)
	2	v	-	-	-	79.45 (AdamW, <i>lr</i> =6e-3)
	3	-	 Image: A set of the set of the	-	-	79.48 (S3, lr =3e-3, same β , w/o NAG, $p = 2$)
	4	-	-	×	-	80.17 (S3, lr =3e-3, diff. β , w/ NAG, $p = 2$)
	5	-	-	-	~	79.74 (S3, lr =3e-3, diff. β , w/o NAG, $p = 3$)
	6	v	 Image: A set of the set of the	-	-	80.25 (S3, lr =6e-3, same β , w/o NAG, $p = 2$)
	7	 Image: A second s	 Image: A second s		-	80.82 (S3, lr =6e-3, same β , w/ NAG, $p = 2$)
	8	 Image: A set of the set of the	-	-	~	79.98 (S3, lr =6e-3, diff. β , w/o NAG, $p = 3$)
	9	 Image: A set of the set of the	 Image: A set of the set of the	-	~	80.31 (S3, lr =6e-3, same β , w/o NAG, $p = 3$)
	10	 Image: A start of the start of	 Image: A set of the set of the	 Image: A start of the start of	 ✓ 	80.93 (S3, lr =6e-3, same β , w/ NAG, $p = 3$)

5.4 SENSITIVITY ANALYSIS FOR HYPERPARAMETERS

We perform a grid search to verify the sensitivity to the momentum order p and the momentum coefficient β of S3 on ViT-B/16 with 150 training epoches. As shown in Figure 6, all combinations achieve an accuracy of 80.20%+, surpassing the 80.13% accuracy achieved by Adam with 300





Figure 5: Ablation study on train loss of S3 for training ViT-B/16 on ImageNet.

494

495

496 497

498

499

500

501 502

516 517

518

519

526

527 528

529

Figure 6: Impact of the momentum order (p) and the momentum coefficient (β) on the Accuracy of S3 training ViT-B/16 on ImageNet.

training epochs. The performance of S3 is not sensitive when p > 1, and p = 1 achieves a slightly lower accuracy. However, as pointed out in Section 3, the computation cost becomes lower when p = 1. Another interesting phenomenon is that setting β to 0.95 obtains the highest accuracies across different p, and p = 3 performs well in most cases.

5.5 VERIFYING THE REASON FOR LOSS SPIKES FROM ADAM

In this subsection, we further experimentally verify that the claim that the potential overlarge update 504 of Adam with relative large learning rate is underlying reason for loss spikes, as discussed in Section 505 2. Figure 7 visually illustrates that convergence of vanilla Adam is attained at the baseline learning 506 rate of 3×10^{-4} despite sporadic spikes, and more frequent spikes and higher loss are observed 507 at a learning rate of 1×10^{-3} , with the same iteration count. Moreover, employing a $10 \times$ higher 508 learning rate of 3×10^{-3} results in premature divergence with pronounced spikes. Noted that all of 509 the phenomenons are aligned with analysis in Section 2. As showcased in Figure 7, naively clipping 510 the Adam update to [-1,1] the range reduces the frequency of loss spikes, but they still occur. This 511 indicates that fine-tuning the clipped value is necessary to balance performance, which complicates 512 the use of clipped updates with tuning a more hyperparameter. In contrast, as we proved in in 513 Theorem 2 in Section 3, when we minimize the maximal update of Adam via $\beta_1 = \beta_2$, loss spikes 514 are completely disappears, which are further verify the correctness of our analyses in Section 2 and Section 3. 515



Figure 7: The Loss spikes phenomenons during Training GPT-2 (345M) on OpenWebText using AdamW.

530 6 CONCLUSION AND DISCUSSION

531 In this paper, we thoroughly examine the strengths and weaknesses of the widely-used optimizer 532 Adam. Building on our analysis, we propose S3, an innovative optimizer that integrates three pivotal improvements over Adam. Comprehensive experiments spanning vision and language tasks 534 showcase S3's accelerated training efficiency and superior inference capabilities. Furthermore, we 535 challenge the conventional belief that Adam's effectiveness stems from simplifying second-order 536 descent, showing instead that its success relies on sign-like descent. This insight paves the way for 537 developing more advanced optimizers. Additionally, We also provide the first theoretical proof of adaptive optimizer convergence from the perspective of sign descent. Most notably, we identify the 538 root cause of loss spikes during LLM training and propose a solution, offering significant benefits for the community in the LLM era.

540 REFERENCES 541

548

554

558

559

560

561

565

566

567

542	Yossi Arjevani, Yair Carmon, John C Duchi, Dylan J Foster, Nathan Srebro, and Blake Woodworth.
543	Lower bounds for non-convex stochastic optimization. Mathematical Programming, 199(1-2):
544	165–214, 2023.

- Jeremy Bernstein, Yu-Xiang Wang, Kamyar Azizzadenesheli, and Animashree Anandkumar. Signs-546 gd: Compressed optimisation for non-convex problems. In International Conference on Machine 547 Learning, pp. 560-569, 2018.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, 549 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are 550 few-shot learners. Advances in neural information processing systems, 33:1877–1901, 2020. 551
- 552 Lizhang Chen, Bo Liu, Kaizhao Liang, and Qiang Liu. Lion secretly solves constrained optimiza-553 tion: As lyapunov predicts. arXiv preprint arXiv:2310.05898, 2023a.
- Xiangning Chen, Chen Liang, Da Huang, Esteban Real, Kaiyuan Wang, Yao Liu, Hieu Pham, Xu-555 anyi Dong, Thang Luong, Cho-Jui Hsieh, et al. Symbolic discovery of optimization algorithms. 556 arXiv preprint arXiv:2302.06675, 2023b.
 - Xiangyi Chen, Sijia Liu, Ruoyu Sun, and Mingyi Hong. On the convergence of a class of adam-type algorithms for non-convex optimization. arXiv preprint arXiv:1808.02941, 2018.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: 562 Scaling language modeling with pathways. Journal of Machine Learning Research, 24(240): 563 1-113, 2023.
 - Soham De, Anirbit Mukherjee, and Enayat Ullah. Convergence guarantees for rmsprop and adam in non-convex optimization and an empirical comparison to nesterov acceleration. arXiv preprint arXiv:1807.06766, 2018.
- 568 Alexandre Défossez, Léon Bottou, Francis Bach, and Nicolas Usunier. A simple convergence proof 569 of adam and adagrad. arXiv preprint arXiv:2003.02395, 2020. 570
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas 571 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An im-572 age is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arX-573 iv:2010.11929, 2020. 574
- 575 Timothy Dozat. Incorporating nesterov momentum into adam. 2016.
- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and 577 stochastic optimization. Journal of machine learning research, 12(7), 2011. 578
- 579 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-580 nition. In IEEE Conference on Computer Vision and Pattern Recognition, pp. 770–778, 2016. 581
- Geoffrey Hinton, Nitish Srivastava, and Kevin Swersky. Neural networks for machine learning 582 lecture 6a overview of mini-batch gradient descent. Cited on, 14(8):2, 2012. 583
- 584 Yusu Hong and Junhong Lin. On convergence of adam for stochastic optimization under relaxed 585 assumptions. arXiv preprint arXiv:2402.03982, 2024. 586
- Hyunjik Kim, George Papamakarios, and Andriy Mnih. The lipschitz constant of self-attention. In 587 International Conference on Machine Learning, pp. 5562–5571, 2021. 588
- 589 Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In International 590 Conference on Learning Representations (ICLR), 2015. 591
- Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete 592 Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. arXiv preprint arXiv:2304.02643, 2023.

594 595 596	Lingkai Kong and Molei Tao. Stochasticity of deterministic gradient descent: Large learning rate for multiscale objective function. <i>Advances in Neural Information Processing Systems</i> , 33:2625–2638, 2020.
597 598	Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convo-
599	lutional neural networks. Communications of the ACM, 60(6):84-90, 2017.
600	Frederik Kunstner, Jacques Chen, Jonathan Wilder Lavington, and Mark Schmidt. Noise is not the
601	main factor behind the gap between sgd and adam on transformers, but sign descent might be.
602	arXiv preprint arXiv:2304.13960, 2023.
603	Haochuan Li, Alexander Rakhlin, and Ali Jadbabaie. Convergence of adam under relaxed assump-
605	tions. Advances in Neural Information Processing Systems, 36, 2023.
606	Hong Liu, Zhivuan Li, David Hall, Percy Liang, and Tengyu Ma. Sonhia: A scalable stochastic
607	second-order optimizer for language model pre-training. arXiv preprint arXiv:2305.14342, 2023.
608	Liyuan Liu, Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Jiawei
610	Han. On the variance of the adaptive learning rate and beyond. arXiv preprint arXiv:1908.03265,
611	2019.
612	Liyuan Liu, Xiaodong Liu, Jianfeng Gao, Weizhu Chen, and Jiawei Han. Understanding the diffi-
613	culty of training transformers. arXiv preprint arXiv:2004.08249, 2020.
614	Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie
615	A convnet for the 2020s. In <i>Proceedings of the IEEE/CVF conference on computer vision and</i>
616	pattern recognition, pp. 11976–11986, 2022.
617	Ilve Leshahiley and Frank Hutter Descupled weight descy regularization an Viv preprint any
618	iv:1711.05101. 2017.
619	
621	Igor Molybog, Peter Albert, Moya Chen, Zachary De Vito, David Esiobu, Naman Goyal, Punit Singh Kaura, Sharan Narang, Andrew Paultan, Puan Silva, et al. A theory on adam instability in large
622	scale machine learning, arXiv preprint arXiv:2304.09871, 2023
623	
624	Yurii Nesterov. A method of solving a convex programming problem with convergence rate o\bigl(k ² \bigr). In <i>Doklady Akademii Nauk</i> . Russian Academy of Sciences, 1983.
626 627	Yurii Nesterov. <i>Introductory lectures on convex optimization: A basic course</i> , volume 87. Springer Science & Business Media, 2013.
628	Lorenzo Noci, Sotiris Anagnostidis, Luca Biggio, Antonio Orvieto, Sidak Pal Singh, and Aurelien
629 630	Lucchi. Signal propagation in transformers: Theoretical perspectives and the role of rank collapse. In Advances in Neural Information Processing Systems, pp. 27198–27211, 2022.
631	Vissbieg O: Jissen Ware, Viber Chen, Value Chi and Lei Zhang. Lingformer, Interducing ling
632 633	chitz continuity to vision transformers. <i>arXiv preprint arXiv:2304.09856</i> , 2023.
634 635	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. <i>OpenAI blog</i> , 1(8):9, 2019.
636	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal
638	Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
639 640	models from natural language supervision. In <i>International conference on machine learning</i> , pp. 8748–8763, 2021.
641 642	Sashank J Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. In <i>International Conference on Learning Representations</i> , 2018.
643 644 645	Herbert Robbins and Sutton Monro. A stochastic approximation method. <i>The annals of mathematical statistics</i> , pp. 400–407, 1951.
646 647	Teven Le Scao, Thomas Wang, Daniel Hesslow, Lucile Saulnier, Stas Bekman, M Saiful Bari, Stella Biderman, Hady Elsahar, Niklas Muennighoff, Jason Phang, et al. What language model to train if you have one million gpu hours? <i>arXiv preprint arXiv:2210.15424</i> , 2022.

657

658

659

660

- Frank Seide, Hao Fu, Jasha Droppo, Gang Li, and Dong Yu. 1-bit stochastic gradient descent and its application to data-parallel distributed training of speech dnns. In *Fifteenth annual conference* of the international speech communication association, 2014.
- Noam Shazeer and Mitchell Stern. Adafactor: Adaptive learning rates with sublinear memory cost. In *International Conference on Machine Learning*, pp. 4596–4604, 2018.
- Tao Sun, Qingsong Wang, Dongsheng Li, and Bao Wang. Momentum ensures convergence of SIGNSGD under weaker assumptions. In *International Conference on Machine Learning*, pp. 33077–33099, 2023.
 - Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971, 2023.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. Advances in neural information processing systems, 30, 2017.
- Sanghyun Woo, Shoubhik Debnath, Ronghang Hu, Xinlei Chen, Zhuang Liu, In So Kweon, and Saining Xie. Convnext v2: Co-designing and scaling convnets with masked autoencoders. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16133–16142, 2023.
- Kingyu Xie, Pan Zhou, Huan Li, Zhouchen Lin, and Shuicheng Yan. Adan: Adaptive nesterov
 momentum algorithm for faster optimizing deep models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- Ruibin Xiong, Yunchang Yang, Di He, Kai Zheng, Shuxin Zheng, Chen Xing, Huishuai Zhang,
 Yanyan Lan, Liwei Wang, and Tieyan Liu. On layer normalization in the transformer architecture.
 In *International Conference on Machine Learning*, pp. 10524–10533. PMLR, 2020.
- Aiyuan Yang, Bin Xiao, Bingning Wang, Borong Zhang, Ce Bian, Chao Yin, Chenxu Lv, Da Pan,
 Dian Wang, Dong Yan, et al. Baichuan 2: Open large-scale language models. *arXiv preprint arXiv:2309.10305*, 2023.
- Greg Yang, Jeffrey Pennington, Vinay Rao, Jascha Sohl-Dickstein, and Samuel S Schoenholz. A
 mean field theory of batch normalization. *arXiv preprint arXiv:1902.08129*, 2019.
- Matthew D Zeiler. Adadelta: an adaptive learning rate method. arXiv preprint arXiv:1212.5701, 2012.
- Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu,
 Wendi Zheng, Xiao Xia, et al. Glm-130b: An open bilingual pre-trained model. *arXiv preprint arXiv:2210.02414*, 2022.
- 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*, 35:28386–28399, 2022.
- Zhongwang Zhang and Zhi-Qin John Xu. Loss spike in training neural networks. *arXiv preprint arXiv:2305.12133*, 2023.
- Pan Zhou, Xingyu Xie, and YAN Shuicheng. Win: Weight-decay-integrated nesterov acceleration for adaptive gradient algorithms. In *International Conference on Learning Representations*, 2023.
- Libin Zhu, Chaoyue Liu, Adityanarayanan Radhakrishnan, and Mikhail Belkin. Catapults in sgd:
 spikes in the training loss and their impact on generalization through feature learning. *arXiv* preprint arXiv:2306.04815, 2023.
- Juntang Zhuang, Tommy Tang, Yifan Ding, Sekhar C Tatikonda, Nicha Dvornek, Xenophon Pa pademetris, and James Duncan. Adabelief optimizer: Adapting stepsizes by the belief in observed gradients. In *Advances in neural information processing systems*, pp. 18795–18806, 2020.

⁷⁰² Appendix

703 704 705

706

A RELATED WORK

Optimizers in Deep Learning. Nowadays, Adam has become the dominant optimizer in deep 707 learning. The adaptivity strategy in Adam traces its roots back to earlier optimizers such as Ada-708 grad (Duchi et al. (2011)), RMSprop (Hinton et al. (2012)), and Adadelta (Zeiler (2012)). Beyond 709 Adam, a wide range of variants are proposed (Dozat (2016); Reddi et al. (2018); Loshchilov & Hut-710 ter (2017); Zhuang et al. (2020); Shazeer & Stern (2018)). SignSGD, the first sign descent method, 711 was proposed to reduce communication costs in distributed learning (Seide et al. (2014)). Subse-712 quently, (Bernstein et al. (2018); Sun et al. (2023)) provided theoretical convergence for SignSGD 713 and introduced an enhanced version. Chen et al. (2023b) applied an auto ML method to discover 714 the sign descent optimizer Lion. This optimizer demonstrated improved performance with a faster convergence rate on various tasks compared to Adam. Recently, (Liu et al. (2023)) introduced an 715 effective second-order optimizer for LLM pre-training. 716

717 Nesterov's Accelerated Gradient (NAG). Theoretical demonstrations by (Nesterov (1983; 2013)) 718 indicate that NAG can achieve faster convergence on convex optimization problems compared to 719 vanilla gradient descent, leveraging gradient information at an extrapolation point to anticipate fu-720 ture trends. NAdam by (Dozat (2016)) was the first to incorporate NAG into adaptive optimizers, 721 modifying the first-order momentum of Adam with NAG. Adan by Xie et al. (2024) integrated a equivalence of NAG into both the first and second momentum of Adam, and Win (Zhou et al. (2023)) 722 applied Nesterov acceleration to the update rather than the first and second momentum. Adan and 723 Win outperformed Adam on various tasks, but they require tuning additional hyperparameters and 724 consume more memory, compared to vanilla Adam. Lion (Chen et al. (2023b)), despite being a 725 sign descent method, exhibits a momentum construction similar to NAG (Chen et al. (2023a)). This 726 resemblance could be a contributing factor to its superior speed and performance over Adam. 727

Training instability and Loss Spikes in LLM Training. Training instability and loss spikes are 728 frequently encountered (Zeng et al. (2022); Chowdhery et al. (2023); Touvron et al. (2023); Yang 729 et al. (2023)) during LLM training, posing challenges to further scaling AI systems. To address 730 this issue, practitioners have employed an ad hoc engineering approach, skipping data batches be-731 fore spikes and restarting training from a nearby checkpoint (Chowdhery et al. (2023)). However, 732 this method requires manual monitoring and intervention, leading to resource wastage. Previous 733 attempts to mitigate instability include embedding norm with BF16, but this comes at a significant 734 performance tax (Scao et al. (2022)). Some researchers found that gradient shrink on the embedding 735 layer reduces loss spikes (Zeng et al. (2022)). Others suggest normalizing the output embedding to 736 lower spike risks (Yang et al. (2023)). (Molybog et al. (2023)) argues that the time-domain corre-737 lation between gradient estimates of earlier layers contributes to training loss instability. Mitigation 738 strategies proposed include tuning down the ϵ value of Adam and reducing batch size. However, it is acknowledged that these methods are not silver bullets for a fundamental solution. 739

740 741

742

B ADDITIONAL EXPERIMENTAL RESULTS

743 B.1 TRAINING SETTING 744

We use the PyTorch vision reference codes 3 to implement vision tasks. For data augmentation, we 745 adhere to the recommended settings in the codes, incorporating RepeatedAugment, AutoAugment 746 Policy (magnitude=9), and Mixup(0.2)/CutMix(1.0) with a probability of 0.5. Additionally, label-747 smoothing with a value of 0.11 is applied. The batch size is set to 1024 for ResNet-50 and 4096748 for ViT-B/16. Regarding the learning rate scheme, we linearly increase it to its peak in the initial 30 749 epochs and then apply a cosine decay, decreasing it to 0 in the subsequent epochs. Other customized 750 hyperparameters for SGD and AdamW are well-established in the codes, and the settings for Adan 751 and lion are followed the recommendations to train ResNet-50 and ViT-B/16 in their respective 752 official papers (Xie et al. (2024); Chen et al. (2023b)). Since NAdam is similar to AdamW, so all its 753 hyperparameters are also the same as AdamW. Specifically, we list the hyperparameters of all the 754 optimizers as follows: 755

³https://github.com/pytorch/vision/tree/main/references/classification



⁴https://github.com/NVIDIA/Megatron-LM

• For S3, we set $\beta = 0.95$, p = 3. We conducted a coarse hyperparameters search on ViT-B/16 (Subsection 5.4) and extended the hypermeters to train LLMs without further tuning.

Notably, we followed the weight decay adjustment strategy outlined in the paper (Chen et al. (2023b)). Specifically, we used the product of the peak learning rate (lr_{Adam}) and the weight decay (λ_{Adam}) from AdamW as a constant. For other optimizers, we just determine the peak learning rate, and the weight decay was derived directly using the formula $\lambda = \frac{lr_{Adam}\lambda_{Adam}}{lr}$. Importantly, the the baseline peak learning rates and weight decays of Adam for training our ResNet-50 and ViT-B-16 are also the same with those reported in (Chen et al. (2023b)), while that for GPT-2 are the same with the paper on Llama (Touvron et al. (2023)).

Table 4: Ablation study on validation perplexity of S3 for training GPT-2(345M).

Exp.	large <i>lr</i>	same β	NAG	flexible p	Validation perplexity
1	-	-	-	-	4.78 (AdamW,lr=3e-4)
2	 Image: A set of the set of the	-	-	-	4.97 (AdamW, <i>lr</i> =1e-3)
3	-	×	-	-	4.77 (S3, lr =3e-4, same β , w/o NAG, $p = 2$)
4	 Image: A second s	~	-	-	4.67 (S3, lr =1e-3, same β , w/o NAG, $p = 2$)
5	 Image: A set of the set of the	×	1	-	4.64 (S3, lr =1e-3, same β , w/ NAG, $p = 2$)
6	 Image: A set of the set of the	 Image: A set of the set of the	 ✓ 	 Image: A second s	4.60 (S3, lr =3e-3, same β , w/ NAG, $p = 3$)

B.2 ADDITIONAL ABLATION STUDY

We also implement ablation experiments for training GPT-2(345M). Figure 10 and Table 4 showcase that employing a large learning rate and sharing the same β alone have little or even a negative impact on performance (e.g., Exp. ① vs. Exp. ②, Exp. ① vs. Exp. ③), while their combination results in a notable improvement (e.g., Exp. ②, Exp. ③ vs. Exp. ④). In contrast, the performance gain from the benefits of NAG is more pronounced (e.g., Exp. ④ vs. Exp. ⑤), and harnessing a larger p can also have an individually beneficial effect on performance (e.g., Exp. ⑥ vs. Exp. ⑥),



Figure 9: Ablation study on train loss of S3 for training GPT-2(345M) on OpenWebText.

B.3 DOWNSTREAM EVALUATION FOR LANGUAGE TASKS

To further validate the effectiveness of the proposed optimizers, we conducted evaluation experiments on pre-trained GPT-2 models, specifically GPT-2 (345M) and GPT-2 (7B), using downstream reasoning benchmarks from OpenCompass ⁵. As depicted in Figure 10, the results demonstrate that S3 consistently outperforms Adam across the majority of benchmarks. This superiority is evident in the improved downstream accuracy, indicating that the lower validation loss achieved by S3 translates into enhanced performance on these reasoning tasks.

An interesting observation is that the superiority of S3 becomes more pronounced as the model size becomes large. This suggests that the benefits of S3 extend beyond its effectiveness with smaller models, showcasing its scalability and adaptability to larger and more complex architectures.

⁵https://github.com/open-compass/opencompass



Figure 10: Zero-shot evaluation of the pre-trained GPT-2 (345M) and GPT-2 (7B) with AdamW, Adan, aLion, and S3 on downstream reasoning tasks.

It is essential to acknowledge that the GPT-2 (345M) model, along with its training dataset, is relatively small. Consequently, the pre-trained models may lack the inherent capabilities needed to perform well on downstream benchmarks, regardless of the optimizer used. Consequently, the accuracies achieved by GPT-2 (345M) with these optimizers may exhibit a degree of randomness due to the model's inherent limitations in handling more complex tasks with a smaller scale.

B.4 VISUALIZATION OF GRADIENT NORMS

As shown in Figure 11, the gradient norms can differ by several orders of magnitude across different layers, and within the same layer, the gradient norms can vary by more than 30 times. This significant variation highlights the challenge of maintaining consistent update magnitudes during the training process.



Figure 11: Visualization of gradient norms within different layers in ViT-B/16 at initialization.

C THEORETICAL PROOFS

C.1 PROOF OF THEOREM 1

Proof. Recalling Eq. (2), we know

$$\boldsymbol{m}_{t}^{(j)} = \frac{1 - \beta_{1}}{1 - \beta_{1}^{t}} \sum_{k=1}^{t} \beta_{1}^{t-k} \boldsymbol{g}_{k}^{(j)}$$

$$\boldsymbol{v}_{t}^{(j)} = \frac{1 - \beta_{2}}{1 - \beta_{2}^{t}} \sum_{k=1}^{t} \beta_{2}^{t-k} (\boldsymbol{g}_{k}^{(j)})^{2}.$$
(11)

Then, $\frac{|\boldsymbol{m}_{t}^{(j)}|}{\sqrt{\boldsymbol{v}_{t}^{(j)}}} = \frac{(1-\beta_{1})\sqrt{1-\beta_{2}^{t}}}{(1-\beta_{1}^{t})\sqrt{1-\beta_{2}}} \cdot \frac{|\sum_{k=1}^{t}\beta_{1}^{t-k}\boldsymbol{g}_{k}^{(j)}|}{\sqrt{\sum_{k=1}^{t}\beta_{2}^{t-k}(\boldsymbol{g}_{k}^{(j)})^{2}}}$ $\stackrel{(i)}{\leq} \frac{(1-\beta_1)\sqrt{1-\beta_2^t}}{(1-\beta_1^t)\sqrt{1-\beta_2}} \cdot \frac{\sum_{k=1}^t \beta_1^{t-k} |\boldsymbol{g}_k^{(j)}|}{\sqrt{\sum_{k=1}^t \beta_2^{t-k} (\boldsymbol{g}_k^{(j)})^2}}$ $\stackrel{(ii)}{\leq} \frac{(1-\beta_1)\sqrt{1-\beta_2^t}}{(1-\beta_1^t)\sqrt{1-\beta_2}} \cdot \frac{\sqrt{\sum_{k=1}^t \beta_2^{t-k} (\boldsymbol{g}_k^{(j)})^2} \sqrt{\sum_{k=1}^t \frac{\beta_1^{2(t-k)}}{\beta_2^{t-k}}}}{\sqrt{\sum_{k=1}^t \beta_2^{t-k} (\boldsymbol{g}_k^{(j)})^2}}$ (12) $=\frac{(1-\beta_1)\sqrt{1-\beta_2^t}}{(1-\beta_1^t)\sqrt{1-\beta_2}}\cdot\sqrt{\sum_{k=1}^t \frac{\beta_1^{2(t-k)}}{\beta_2^{t-k}}}$ $\stackrel{(iii)}{=} \frac{(1-\beta_1)\sqrt{1-\beta_2^t}\sqrt{1-(\frac{\beta_1^2}{\beta_2})^t}}{(1-\beta_1^t)\sqrt{1-\beta_2}\sqrt{1-\frac{\beta_1^2}{\beta_2}}}$ $\simeq \frac{1-\beta_1}{\sqrt{1-\beta_2}\sqrt{1-\frac{\beta_1^2}{\beta_2}}},$

where (i) holds due to the fact $|\boldsymbol{a}^{(j)} + \boldsymbol{b}^{(j)}| \leq |\boldsymbol{a}^{(j)}| + |\boldsymbol{b}^{(j)}|$; (ii) holds due to Cauchy-Schiwaz inequality; (iii) holds due to $\beta_1^2 \leq \beta_2$. $\frac{|\boldsymbol{m}_t^{(j)}|}{\sqrt{\boldsymbol{v}_t^{(j)}}}$ reach to the largest value if the signs of $\{\boldsymbol{g}_t^{(j)}, \boldsymbol{g}_{t-1}^{(j)}, \ldots\}$ are the same and $|\boldsymbol{g}_t^{(j)}| = \frac{\beta_2 |\boldsymbol{g}_{t-1}^{(j)}|}{\beta_1} = \frac{\beta_2^2 |\boldsymbol{g}_{t-2}^{(j)}|}{\beta_1^2} = \ldots = \frac{\beta_2^k |\boldsymbol{g}_{t-k}^{(j)}|}{\beta_1^k}$

C.2 PROOF OF THEOREM 2

Proof. (1). According to S3 in Algorithm 1, we have

$$egin{split} egin{aligned} egin{aligned} eta \ eta \$$

$$\mathbf{b}_{t}^{(j)}(p) = \left((1 - \beta_2) \left(\sum_{k=1}^{t} \beta_2^{t-k+1} |\mathbf{g}_k^{(j)}|^p + |\mathbf{g}_t^{(j)}|^p \right) \right)^{1/p}$$

 $n_{t}^{(}$

(13)

Then, assuming q satisfies $\frac{1}{p} + \frac{1}{q} = 1$, we obtain

$$\frac{|\boldsymbol{n}_{t}^{(j)}|}{\boldsymbol{b}_{t}^{(j)}(p)} = \frac{1-\beta_{1}}{(1-\beta_{2})^{1/p}} \cdot \frac{|\sum_{k=1}^{t}\beta_{1}^{t-k+1}\boldsymbol{g}_{k}^{(j)} + \boldsymbol{g}_{t}^{(j)}|}{\left(\sum_{k=1}^{t}\beta_{2}^{t-k+1}|\boldsymbol{g}_{k}^{(j)}|^{p} + |\boldsymbol{g}_{t}^{(j)}|^{p}\right)^{1/p}} \\ \stackrel{(i)}{\leq} \frac{1-\beta_{1}}{(1-\beta_{2})^{1/p}} \cdot \frac{\sum_{k=1}^{t}\beta_{1}^{t-k+1}|\boldsymbol{g}_{k}^{(j)}| + |\boldsymbol{g}_{t}^{(j)}|}{\left(\sum_{k=1}^{t}\beta_{2}^{t-k+1}|\boldsymbol{g}_{k}^{(j)}|^{p} + |\boldsymbol{g}_{t}^{(j)}|^{p}\right)^{1/p}}$$

 $\stackrel{(ii)}{\leq} \frac{1-\beta_{1}}{(1-\beta_{2})^{1/p}} \cdot \frac{\left(\sum_{k=1}^{t} \beta_{2}^{t-k+1} | \boldsymbol{g}_{k}^{(j)} |^{p} + | \boldsymbol{g}_{t}^{(j)} |^{p}\right)^{1/p} \left(\sum_{k=1}^{t} \left(\frac{\beta_{1}^{(t-k+1)}}{\beta_{2}^{\frac{1}{p}(t-k+1)}}\right)^{q} + 1\right)^{1/q}}{\left(\sum_{k=1}^{t} \beta_{2}^{t-k+1} | \boldsymbol{g}_{k}^{(j)} |^{p} + | \boldsymbol{g}_{t}^{(j)} |^{p}\right)^{1/p}} = \frac{1-\beta_{1}}{(1-\beta_{2})^{1/p}} \cdot \left(\sum_{k=1}^{t} \left(\frac{\beta_{1}^{(t-k+1)}}{\beta_{2}^{\frac{1}{p}(t-k+1)}}\right)^{q} + 1\right)^{1/q}}{\left(\sum_{k=1}^{t} \left(\frac{\beta_{1}^{(t-k+1)}}{\beta_{2}^{\frac{1}{p}(t-k+1)}}\right)^{q} + 1\right)^{1/q}}$ $\stackrel{(iii)}{\leq} \frac{1-\beta_{1}}{(1-\beta_{2})^{1/p} \left(1-\frac{\beta_{1}}{\beta_{2}^{\frac{q}{p}}}\right)^{1/q}},$ (14)

where (i) holds due to the fact $|a + b| \le |a| + |b|$; (ii) holds due to Hölder inequality $\sum_{i=1}^{s} a_i b_i \le (\sum_{i=1}^{s} a_i^p)^{1/p} (\sum_{i=1}^{s} b_i^q)^{1/q}$ if $\frac{1}{p} + \frac{1}{q} = 1$ and $p, q \ge 1$; (iii) holds due to $\beta_1 \le \beta_2^{1/p}$.

(2). The upper bound of each element of the update $\frac{n_t^{(j)}}{b_t^{(j)}}$ can be

$$\frac{1-\beta_{1}}{(1-\beta_{2})^{1/p}\left(1-\frac{\beta_{1}^{q}}{\beta_{2}^{q/p}}\right)^{1/q}} \stackrel{(i){=}{=} \frac{1-\beta_{1}}{\frac{1}{p}(1-\beta_{2})+\frac{1}{q}(1-\frac{\beta_{1}^{q}}{\beta_{2}^{q/p}})}$$

$$\frac{(ii){=} \frac{1-\beta_{1}}{1-(\frac{\beta_{2}}{p}+\frac{\beta_{1}^{q}}{q\beta_{2}^{q/p}})}$$

$$\stackrel{(iii){=}{=} \frac{1-\beta_{1}}{1-\beta_{1}^{1/p}}$$

$$=1,$$
(15)

where (i) holds due to Young's inequality $\frac{a}{p} + \frac{b}{q} \ge a^{1/p}b^{1/q}$; (ii) holds owing to $\frac{1}{p} + \frac{1}{q} = 1$; (iii) holds resulting from Young's inequality again.

Notably, according to the property of Young's inequality, the equality in (i) and (iii) can be reached, if and only if

$$\beta_2 = \frac{\beta_1^q}{\beta_2^{q/p}} \quad \Rightarrow \quad \beta_1^q = \beta_2^{1+\frac{q}{p}} \quad \Rightarrow \quad \beta_1^q = \beta_2^{1+(1-\frac{1}{q})q} \quad \Rightarrow \quad \beta_1^q = \beta_2^q \quad \Rightarrow \quad \beta_1 = \beta_2 \quad (16)$$

1025 Therefore, when $\beta_1 = \beta_2$, the upper bound of each element of the update $\frac{n_t^{(j)}}{b_t^{(j)}}$ reaches to the smallest 1. 1026 (3). Following S3 in Algorithm 1, we have 1027 1028 $\boldsymbol{b}_{t}^{(j)}(p_{1}) = \left((1-\beta) \left(\sum_{k=1}^{t} \beta^{t-k+1} |\boldsymbol{g}_{k}^{(j)}|^{p_{1}} + |\boldsymbol{g}_{t}^{(j)}|^{p_{1}} \right) \right)^{p_{1}}$ 1029 1030 (17)1031 $\boldsymbol{b}_{t}^{(j)}(p_{2}) = \left((1-\beta) \left(\sum_{k=1}^{t} \beta^{t-k+1} |\boldsymbol{g}_{k}^{(j)}|^{p_{2}} + |\boldsymbol{g}_{t}^{(j)}|^{p_{2}} \right) \right)^{1/p_{2}}.$ 1032 1033 1034 Denoting $r = \frac{p_2}{p_1}$, we then obtain 1035 1036 $(\boldsymbol{b}_t^{(j)}(p_1))^{p_2} = (\boldsymbol{b}_t^{(j)}(p_1))^{rp_1} = \left((1-\beta) \left(\sum_{k=1}^t \beta^{t-k+1} |\boldsymbol{g}_k^{(j)}|^{p_1} + |\boldsymbol{g}_t^{(j)}|^{p_1} \right) \right)^r$ 1037 1038 1039 $\leq (1-\beta) \left(\sum_{i=1}^{t} \beta^{t-k+1} | \boldsymbol{g}_{k}^{(j)} |^{rp_{1}} + | \boldsymbol{g}_{t}^{(j)} |^{rp_{1}} \right)$ 1040 (18)1041 $= (1 - \beta) \left(\sum_{k=1}^{t} \beta^{t-k+1} |\boldsymbol{g}_{k}^{(j)}|^{p_{2}} + |\boldsymbol{g}_{t}^{(j)}|^{p_{2}} \right)$ 1043 1044 1045 $=(\boldsymbol{b}_{t}^{(j)}(p_{2}))^{p_{2}},$ 1046 where the inequality holds due to Jensen's inequality and the fact $(1 - \beta)(\sum_{k=1}^{t} \beta^{t-k+1} + 1) < 1$. 1047 1048 1049 C.3 **PROOF OF THEOREM 3** 1050 Proof. We first deduce from ASGD(I) to ASGD(II). According to (I), we have 1051 1052 $\tilde{\boldsymbol{x}}_{t+1} = \tilde{\boldsymbol{x}}_t - \gamma \boldsymbol{m}_t$ (19)1053 $= \tilde{\boldsymbol{x}}_{t} - \gamma(\beta \boldsymbol{m}_{t-1} + (1-\beta)\nabla f(\tilde{\boldsymbol{x}}_{t} - \gamma\beta \boldsymbol{m}_{t-1}; \zeta_{t}))$ 1054 1055 Subtracting $\gamma \beta m_t$ on both sides, we obtain 1056 $\tilde{\boldsymbol{x}}_{t+1} - \gamma \beta \boldsymbol{m}_t = \tilde{\boldsymbol{x}}_t - \gamma \beta \boldsymbol{m}_{t-1} - \gamma (\beta \boldsymbol{m}_t - (1-\beta)\nabla f(\tilde{\boldsymbol{x}}_t - \gamma \beta \boldsymbol{m}_{t-1}; \zeta_t))$ (20)1057 1058 Setting $\boldsymbol{x}_t = \tilde{\boldsymbol{x}}_t - \gamma \beta \boldsymbol{m}_{t-1}$, we further have 1059 $\boldsymbol{x}_{t+1} = \boldsymbol{x}_t - \gamma(\beta \boldsymbol{m}_t + (1-\beta)\nabla f(\boldsymbol{x}_t;\zeta_t))$ (21)1060 1061 Thus, ASGD(I) becomes ASGD(II). 1062 1063 Then, we deduce from ASGD(III). Denoting $n_t = \beta m_t + (1 - \beta)g_t$, we have 1064 $\boldsymbol{n}_t - \beta \boldsymbol{n}_{t-1} = \beta \boldsymbol{m}_t + (1-\beta)\boldsymbol{g}_t - \beta \boldsymbol{n}_{t-1}$ $= (1-\beta)\boldsymbol{g}_t + \beta(\beta\boldsymbol{m}_{t-1} + (1-\beta)\boldsymbol{g}_t) - \beta\boldsymbol{n}_{t-1}$ $=(1-\beta)\boldsymbol{g}_t + \beta(\beta\boldsymbol{m}_{t-1} + (1-\beta)\boldsymbol{g}_t) - \beta(\beta\boldsymbol{m}_{t-t} + (1-\beta)\boldsymbol{g}_{t-1})$ (22)1067 $=(1-\beta)\boldsymbol{g}_t+\beta(1-\beta)(\boldsymbol{g}_t-\boldsymbol{g}_{t-1})$ 1068 1069 $= (1 - \beta)(\boldsymbol{q}_t + \beta(\boldsymbol{q}_t - \boldsymbol{q}_{t-1})).$ 1070 1071 It indicates $n_t = m_t + \beta r_t$ where $m_t = \beta m_{t-1} + (1-\beta)g_t$ and $r_t = \beta r_{t-1} + (1-\beta)(g_t - g_{t-1})$. 1072 Therefore, ASGD (II) is equivalent to ASGD (III). 1073 1074 C.4 AUXILIARY LEMMA 1075 **Lemma 5.** Under Assumption 2, for any $x, y \in \mathbb{R}^d$ with $||x - y||_2 \leq R$, the function obeys 1076 1077 $F(\boldsymbol{y}) \leq F(\boldsymbol{x}) + \langle \nabla F(\boldsymbol{x}), \boldsymbol{y} - \boldsymbol{x} \rangle + \frac{L_0 + L_1 \|\nabla F(\boldsymbol{x})\|_2}{2} \|\boldsymbol{y} - \boldsymbol{x}\|_2^2.$ 1078 (23)1079

Proof. For any $x, y \in \mathbb{R}^d$ with $||x - y||_2 \leq R$, we have $F(\boldsymbol{y}) = F(\boldsymbol{x}) + \int_{a}^{1} \langle \nabla F(\boldsymbol{x} + t(\boldsymbol{y} - \boldsymbol{x})), \boldsymbol{y} - \boldsymbol{x} \rangle dt$ $=F(\boldsymbol{x}) + \langle \nabla F(\boldsymbol{x}), \boldsymbol{y} - \boldsymbol{x} \rangle + \int^{1} \langle \nabla F(\boldsymbol{x} + t(\boldsymbol{y} - \boldsymbol{x})) - \nabla F(\boldsymbol{x}), \boldsymbol{y} - \boldsymbol{x} \rangle dt$ $\overset{(i)}{\leq} F(\boldsymbol{x}) + \langle \nabla F(\boldsymbol{x}), \boldsymbol{y} - \boldsymbol{x} \rangle + \int_{0}^{1} \|\nabla F(\boldsymbol{x} + t(\boldsymbol{y} - \boldsymbol{x})) - \nabla F(\boldsymbol{x})\|_{2} \|\boldsymbol{y} - \boldsymbol{x}\|_{2} dt$ (24) $\overset{(ii)}{\leq} F(\boldsymbol{x}) + \langle \nabla F(\boldsymbol{x}), \boldsymbol{y} - \boldsymbol{x} \rangle + (L_0 + L_1 \nabla F(\boldsymbol{x})) \| \boldsymbol{y} - \boldsymbol{x} \|_2^2 \int_0^1 t d_t$ $=F(\boldsymbol{x})+\langle \nabla F(\boldsymbol{x}),\boldsymbol{y}-\boldsymbol{x}\rangle+\frac{L_0+L_1\|\nabla F(\boldsymbol{x})\|}{2}\|\boldsymbol{y}-\boldsymbol{x}\|_2^2,$ where (i) holds due to Cauchy-Schwarz inequality, and (ii) holds due to Assumption 2. C.5 PROOF OF THEOREM 4 **Proof.** Following Lemma 5 with $x_{t+1} \rightarrow y$ and $x_t \rightarrow x$, we have $F(\boldsymbol{x}_{t+1}) \leq F(\boldsymbol{x}_{t}) + \langle \nabla F(\boldsymbol{x}_{t}), \boldsymbol{x}_{t+1} - \boldsymbol{x}_{t} \rangle + \frac{L_{0} + L_{1} \|\nabla F(\boldsymbol{x}_{t})\|_{2}}{2} \|\boldsymbol{x}_{t+1} - \boldsymbol{x}_{t}\|_{2}^{2}$ (25)Recalling the update rule $x_{t+1} = x_t - \gamma \frac{n_t}{b_t} = x_t - \gamma \frac{|n_t|}{b_t} \circ \frac{n_t}{|n_t|} = x_t - \gamma u_t \circ \operatorname{Sign}(n_t)$, we further obtain $F(\boldsymbol{x}_{t+1}) \leq F(\boldsymbol{x}_t) - \langle \nabla F(\boldsymbol{x}_t), \gamma \boldsymbol{u}_t \circ \operatorname{Sign}(\boldsymbol{n}_t) \rangle + \frac{(L_0 + L_1 \|\nabla F(\boldsymbol{x}_t)\|_2)\gamma^2}{2} \|\boldsymbol{u}_t\|_2^2$ $=F(\boldsymbol{x}_{t}) - \langle \nabla F(\boldsymbol{x}_{t}), \gamma \boldsymbol{u}_{t} \circ \operatorname{Sign}(\nabla F(\boldsymbol{x}_{t})) \rangle + \underbrace{\langle \nabla F(\boldsymbol{x}_{t}), \gamma \boldsymbol{u}_{t} \circ (\operatorname{Sign}(\nabla F(\boldsymbol{x}_{t})) - \operatorname{Sign}(\boldsymbol{n}_{t})) \rangle}_{T_{t}}$ + $\frac{(L_0 + L_1 \|\nabla F(\boldsymbol{x}_t)\|)\gamma^2}{2} \|\boldsymbol{u}_t\|_2^2$ (26)There are two cases for each element of T_1 . If $\operatorname{Sign}(\nabla F(\boldsymbol{x}_t)^{(j)}) = \operatorname{Sign}(\boldsymbol{n}_t^{(j)}), \nabla F(\boldsymbol{x}_t)^{(j)} \cdot \boldsymbol{u}_t^{(j)}$. $\left(\operatorname{Sign}(\nabla F(\boldsymbol{x}_t))^{(j)} - \operatorname{Sign}(\boldsymbol{n}_t^{(j)})\right) = 0.$ If $\operatorname{Sign}(\nabla F(\boldsymbol{x}_t)^{(j)}) \neq \operatorname{Sign}(\boldsymbol{n}_t^{(j)}), \nabla F(\boldsymbol{x}_t)^{(j)} \cdot \boldsymbol{u}_t^{(j)}$ $\left(\operatorname{Sign}(\nabla F(\boldsymbol{x}_t))^{(j)} - \operatorname{Sign}(\boldsymbol{n}_t^{(j)})\right) = 2\boldsymbol{u}_t^{(j)} |\nabla F(\boldsymbol{x}_t)^{(j)}| \le 2\boldsymbol{u}_t^{(j)} |\nabla F(\boldsymbol{x}_t)^{(j)} - \boldsymbol{n}_t^{(j)}|, \text{ hence } T_1 = 0$ $2\sum_{i=1}^{d} u_t^{(j)} |\nabla F(x_t)^{(j)} - n_t^{(j)}|.$ Rearranging Eq. (26), we have $F(\boldsymbol{x}_{t+1}) \leq F(\boldsymbol{x}_t) - \langle \nabla F(\boldsymbol{x}_t), \gamma \boldsymbol{u}_t \circ \operatorname{Sign}(\nabla F(\boldsymbol{x}_t)) \rangle + 2\gamma \sum_{t=1}^d \boldsymbol{u}_t^{(j)} |\nabla F(\boldsymbol{x}_t)^{(j)} - \boldsymbol{n}_t^{(j)}|$ + $\frac{(L_0 + L_1 \|\nabla F(\boldsymbol{x}_t)\|_2)\gamma^2}{2} \|\boldsymbol{u}_t\|_2^2$ $\leq F(\boldsymbol{x}_t) - \gamma u_{\min} \|\nabla F(\boldsymbol{x}_t)\|_1 + 2\gamma \|\boldsymbol{n}_t - \nabla F(\boldsymbol{x}_t)\|_1 + \frac{\gamma^2 d(L_0 + L_1 \|\nabla F(\boldsymbol{x}_t)\|_2)}{2}$ $\leq F(\boldsymbol{x}_t) - \gamma u_{\min} \|\nabla F(\boldsymbol{x}_t)\|_1 + 2\gamma \sqrt{d} \|\boldsymbol{n}_t - \nabla F(\boldsymbol{x}_t)\|_2 + \frac{\gamma^2 d(L_0 + L_1 \|\nabla F(\boldsymbol{x}_t)\|_1)}{2}$ (27)where the second inequality holds due to $0 < u_{\min} \le u_t^{(j)} \le 1$, and the third inequality holds owing

to the fact $\|\boldsymbol{a}\|_2 \leq \|\boldsymbol{a}\|_1 \leq \sqrt{d} \|\boldsymbol{a}\|_2$ for any $\boldsymbol{a} \in \mathbb{R}^d$.

Summing over 1 to T and taking expectation on it, we have

$$\frac{(u_{\min} - \frac{\gamma dL_1}{2})}{T} \sum_{t=1}^{T} \mathbb{E}[\|\nabla F(\boldsymbol{x}_t)\|_1] \le \frac{F(\boldsymbol{x}_1) - F(\boldsymbol{x}^*)}{\gamma T} + \frac{2\sqrt{d}}{T} \sum_{t=1}^{T} \mathbb{E}[\|\boldsymbol{n}_t - \nabla F(\boldsymbol{x}_t)\|_2] + \frac{\gamma dL_0}{2}$$
(28)

Recalling $\boldsymbol{m}_t = \beta \boldsymbol{m}_{t-1} + (1 - \beta) \boldsymbol{g}_t$, we obtain

Utilizing recursion, we further have

$$m_{t} - \nabla F(\boldsymbol{x}_{t}) = -\beta^{t} \nabla F(\boldsymbol{x}_{1}) + (1-\beta) \sum_{k=1}^{t} \beta^{t-k} (\boldsymbol{g}_{k} - \nabla F(\boldsymbol{x}_{k})) - \sum_{k=1}^{t} \beta^{t-k+1} (\nabla F(\boldsymbol{x}_{k}) - \nabla F(\boldsymbol{x}_{k-1})),$$
(30)

where $m_1 - \nabla F(x_1) = -\beta_1 \nabla F(x_1) + (1 - \beta_1)(g_1 - \nabla F(x_1))$ due to $m_0 = 0$. Hence,

$$\boldsymbol{n}_{t} - \nabla F(\boldsymbol{x}_{t}) = \beta(\boldsymbol{m}_{t} - \nabla F(\boldsymbol{x}_{t})) + (1 - \beta)(\boldsymbol{g}_{t} - \nabla F(\boldsymbol{x}_{t}))$$

$$= -\beta^{t+1} \nabla F(\boldsymbol{x}_{1}) + (1 - \beta) \left(\sum_{k=1}^{t} \beta^{t-k+1}(\boldsymbol{g}_{k} - \nabla F(\boldsymbol{x}_{k}) + (\boldsymbol{g}_{t} - \nabla F(\boldsymbol{x}_{t}))) \right)$$

$$-\beta^{2} \sum_{k=1}^{t} \beta^{t-k} (\nabla F(\boldsymbol{x}_{k}) - \nabla F(\boldsymbol{x}_{k-1})),$$
(31)

Then, we obtain

$$\frac{1170}{1171} \qquad \frac{1}{T} \sum_{t=1}^{T} \mathbb{E} \left[\| \boldsymbol{n}_{t} - \nabla F(\boldsymbol{x}_{t}) \|_{2} \right] \leq \underbrace{\frac{\beta}{T} \sum_{t=1}^{T} \beta^{t} \| \nabla F(\boldsymbol{x}_{1}) \|_{2}}_{T_{2}}}_{T_{2}} + \underbrace{\frac{1 - \beta}{T} \sum_{t=1}^{T} \mathbb{E} \left[\left\| \sum_{k=1}^{t} \beta^{t-k+1} (\boldsymbol{g}_{k} - \nabla F(\boldsymbol{x}_{k})) \right\|_{2} + \| \boldsymbol{g}_{t} - \nabla F(\boldsymbol{x}_{t}) \|_{2} \right]}_{T_{3}}}_{T_{3}} + \underbrace{\frac{\beta^{2}}{T} \sum_{t=1}^{T} \mathbb{E} \left[\left\| \sum_{k=1}^{t} \beta^{t-k} (\nabla F(\boldsymbol{x}_{k}) - \nabla F(\boldsymbol{x}_{k-1})) \right\|_{2} \right]}_{T_{4}}}_{T_{4}}$$
(32)

In terms of T_2 , we obtain

$$T_2 = \frac{\beta}{T} \sum_{t=1}^T \beta^t \left\| \nabla F(\boldsymbol{x}_1) \right\|_2 \le \frac{\beta}{(1-\beta)T} \left\| \nabla F(\boldsymbol{x}_1) \right\|_2$$
(33)

As for
$$T_3$$
, we have

$$T_3 = \frac{1 - \beta}{T} \sum_{l=1}^{T} \mathbb{E} \left[\left\| \sum_{k=1}^{t} \beta^{t-k+1} (g_k - \nabla F(x_k)) \right\|_2^2 + \|g_l - \nabla F(x_l)\|_2^2 \right]$$

$$\stackrel{(i)}{\leq} \frac{1 - \beta}{T} \sum_{l=1}^{T} \sqrt{\mathbb{E} \left[\left\| \sum_{k=1}^{t} \beta^{t-k+1} (g_k - \nabla F(x_k)) \right\|_2^2 + \|g_l - \nabla F(x_l)\|_2^2 \right]}$$

$$\stackrel{(ii)}{=} \frac{1 - \beta}{T} \sum_{l=1}^{T} \sqrt{\sum_{k=1}^{t} \beta^{2(t-k+1)} \mathbb{E} \left[\|g_k - \nabla F(x_k)\|_2^2 + \|g_l - \nabla F(x_l)\|_2^2 \right]}$$

$$\stackrel{(iii)}{\leq} \frac{1 - \beta}{T} \sum_{l=1}^{T} \sqrt{\sum_{k=1}^{t+1} \beta^{2(t-k+1)} \sigma^2}$$

$$\stackrel{(iii)}{\leq} \frac{1 - \beta}{\sqrt{1 - \beta^2}} \sigma$$

$$\leq \sqrt{1 - \beta \sigma},$$
where (i) holds ue to the fact $(\mathbb{E}[Z])^2 \leq \mathbb{E}[Z^2]$; (ii) holds owing to $\mathbb{E}[g_k - \nabla F(x_k)]|_2^2 = \sigma^2$ according to
Assumption 3; (iii) holds resulting from $\mathbb{E} \left[\|g_k - \nabla F(x_k)\|_2^2 \right] \leq \sigma^2$ according to
Assumption 4.
Now we turn attention to T_4 , i.e.,

$$T_4 = \frac{\beta^2}{T} \sum_{t=1}^{T} \sum_{k=1}^{t} \beta^{t-k} \mathbb{E} \left[\|\nabla F(x_k) - \nabla F(x_{k-1})\|_2 \right]$$

$$\stackrel{(ii)}{\leq} \frac{\beta^2}{T} \sum_{t=1}^{T} \sum_{k=1}^{t} \beta^{t-k} \mathbb{E} \left[|\nabla F(x_k) - \nabla F(x_{k-1})||_2 \right]$$

$$\stackrel{(ii)}{\leq} \frac{\beta^2}{T} \sum_{t=1}^{T} \sum_{k=1}^{t} \beta^{t-k} \mathbb{E} \left[|\nabla F(x_k)||_2 \right] \|u_{t-1}\|_2 \right]$$

$$\stackrel{(iii)}{\leq} \frac{\beta^2}{T} \sum_{t=1}^{T} \sum_{k=1}^{t} \beta^{t-k} \mathbb{E} \left[|\nabla F(x_k)\|_2 \right] \|u_{t-1}\|_2 \right]$$

$$\stackrel{(ii)}{\leq} \frac{\beta^2 L_0 \gamma \sqrt{d}}{1 - \beta} + \beta^2 L_1 \gamma \sqrt{d} \sum_{k=1}^{T} \mathbb{E} \left[\|\nabla F(x_k)\|_2 \right] \sum_{k=k}^{T} \beta^{t-k} \\$$

$$\leq \frac{\beta^2 L_0 \gamma \sqrt{d}}{1 - \beta} + \frac{\beta^2 L_1 \gamma \sqrt{d}}{(1 - \beta)T} \sum_{t=1}^{T} \mathbb{E} \left[\|\nabla F(x_k)\|_2 \right]$$

$$\stackrel{(i))}{=} \frac{\beta^2 L_0 \gamma \sqrt{d}}{1 - \beta} + \frac{\beta^2 L_1 \gamma \sqrt{d}}{(1 - \beta)T} \sum_{t=1}^{T} \mathbb{E} \left[\|\nabla F(x_k)\|_2 \right]$$

where (i) holds due to the fact $||a + b||_2 \le ||a||_2 + ||b||_2$; (ii) holds owing to Assumption 2; (iii) holds due to the update rule; (iv) holds depends on $u^{(j)} \le 1$ according to Theorem 2; (v) holds resulting from the fact that $\sum_{i=1}^{n} \sum_{j=1}^{i} a_{i,j} = \sum_{j=1}^{n} \sum_{i=j}^{n} a_{i,j}$; (vi) holds due to the fact $||a||_2 \le ||a||_1$. 1242 Combining Eq.(32) - Eq.(35), we have

$$\frac{1}{T}\sum_{t=1}^{T} \mathbb{E}\left[\|\boldsymbol{n}_{t} - \nabla F(\boldsymbol{x}_{t})\|_{2}\right] \leq \frac{\beta}{(1-\beta)T} \|\nabla F(\boldsymbol{x}_{1})\|_{2} + \sqrt{1-\beta}\sigma + \frac{\beta^{2}L_{0}\gamma\sqrt{d}}{1-\beta} + \frac{\beta^{2}L_{1}\gamma\sqrt{d}}{(1-\beta)T}\sum_{t=1}^{T} \mathbb{E}[\|\nabla F(\boldsymbol{x}_{t})\|_{1}] \leq \frac{\beta}{(1-\beta)T} \|\nabla F(\boldsymbol{x}_{1})\|_{2} + \sqrt{1-\beta}\sigma + \frac{\beta^{2}L_{0}\gamma\sqrt{d}}{1-\beta} + \frac{\beta^{2}L_{1}\gamma\sqrt{d}}{(1-\beta)T}\sum_{t=1}^{T} \mathbb{E}[\|\nabla F(\boldsymbol{x}_{t})\|_{1}]$$
(36)

Combining Eq.(28) and Eq.(36), we obtain

$$\frac{u_{\min} - \frac{\gamma dL_1}{2} - \frac{2\beta^2 L_1 \gamma \sqrt{d}}{1 - \beta}}{T} \sum_{t=1}^T \mathbb{E}[\|\nabla F(\boldsymbol{x}_t)\|_1] \leq \frac{F(\boldsymbol{x}_1) - F(\boldsymbol{x}^*)}{\gamma T} + \frac{2\beta \sqrt{d}}{(1 - \beta)T} \mathbb{E}[\|\nabla F(\boldsymbol{x}_1)\|_2] + 2\sqrt{(1 - \beta)d\sigma} + \frac{2\gamma \beta^2 L_0 d}{1 - \beta} + \frac{\gamma dL_0}{2}.$$
(37)

1263
1264 Let
$$\gamma = \frac{1}{L_0 T^{3/4}}, 1 - \beta = \frac{1}{T^{1/2}}$$
. When $T \ge \max\{(\frac{2dL_1}{L_0 u_{\min}})^{4/3}, (\frac{8\beta^2 \sqrt{d}L_1}{(1-\beta)L_0 u_{\min}})^4\}$, we can guarantee

$$u_{\min} - \frac{\gamma dL_1}{2} - \frac{2\gamma \sqrt{dL_1}}{1-\beta} \ge \frac{u_{\min}}{2}.$$
 (38)

1269 Then, setting $U_{\text{max}} = \frac{1}{u_{\text{min}}}$, we reformulate Eq. (37) as

$$\frac{1}{T} \sum_{t=1}^{T} \mathbb{E}[\|\nabla F(\boldsymbol{x}_{t})\|_{1}] \leq \frac{2L_{0}U_{\max}(F(\boldsymbol{x}_{1}) - F(\boldsymbol{x}^{*}))}{T^{1/4}} + \frac{4\beta U_{\max}\sqrt{d}\mathbb{E}\left[\|\nabla F(\boldsymbol{x}_{1})\|_{2}\right]}{T^{1/2}} + \frac{4U_{\max}\sqrt{d}\sigma}{T^{1/4}} + \frac{4\beta^{2}U_{\max}d}{T^{1/4}} + \frac{U_{\max}d}{T^{3/4}}.$$
(39)