
Robustness to Unbounded Smoothness of Generalized SignSGD

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Traditional analyses in non-convex optimization typically rely on the smoothness
2 assumption, namely requiring the gradients to be Lipschitz. However, recent
3 evidence shows that this smoothness condition does not capture the properties of
4 some deep learning objective functions, including the ones involving Recurrent
5 Neural Networks and LSTMs. Instead, they satisfy a much more relaxed condition,
6 with potentially unbounded smoothness. Under this relaxed assumption, it has been
7 theoretically and empirically shown that the gradient-clipped SGD has an advantage
8 over the vanilla one. **In this paper, we show that clipping is not indispensable
9 for Adam-type algorithms** in tackling such scenarios: we theoretically prove
10 that a generalized SignSGD algorithm can obtain similar convergence rates as
11 SGD with clipping but does not need explicit clipping at all. This family of
12 algorithms on one end recovers SignSGD and on the other end closely resembles the
13 popular Adam algorithm. **Our analysis underlines the critical role that momentum
14 plays in analyzing SignSGD-type and Adam-type algorithms:** it not only reduces
15 the effects of noise, thus removing the need for large mini-batch in previous
16 analyses of SignSGD-type algorithms, but it also substantially reduces the effects
17 of unbounded smoothness and gradient norms. We also compare these algorithms
18 with popular optimizers on a set of deep learning tasks, observing that we can
19 match the performance of Adam while beating the others.

20 1 Introduction

21 Recent years have witnessed a surge in non-convex machine learning models, with a focus on deep
22 neural networks [26]. DNNs have achieved tremendous progress in a variety of tasks, including
23 computer vision [25, 17, 22], natural language processing [10, 46], and a lot more. Despite their huge
24 empirical success, the theoretical analyses of non-convex optimization [20] prove to be fundamentally
25 more challenging than the established convex optimization theory [4]. Among the numerous literature,
26 many of them assume smoothness of the objective function, namely requiring the gradients to be
27 Lipschitz. Under this scenario, past works have succeeded in proving the convergence rates for a
28 number of algorithms, e.g., Stochastic Gradient Descent [13], AdaGrad [48, 29], and STORM [7].

29 Nevertheless, it was recently observed that the smoothness assumption does not capture the training
30 of LSTMs [19]: the Hessian can grow with the size of the gradients [52]. Inspired by this, Zhang et
31 al. [52] proposed a relaxed smoothness assumption, named (L_0, L_1) smoothness:

$$\|\nabla^2 F(\mathbf{x})\| \leq L_0 + L_1 \|\nabla F(\mathbf{x})\|. \quad (1)$$

32 They also showed that the well-known gradient clipping technique can ensure SGD's convergence
33 in such scenarios. Later, their results were improved to show that SGD with clipping can be made
34 unaffected by the L_1 in (1) and is able to recover the optimal convergence rate of SGD under the
35 original smoothness setting [51, 21].

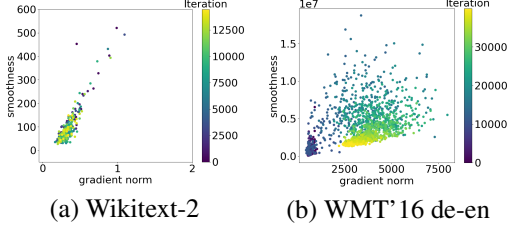


Figure 1: Local gradient Lipschitz constant vs. Gradient norm on training (a) a 2-layer Transformer Encoder on Wikitext-2 (b) a 6-layer Transformer on WMT’16 Multimodal Machine Translation de-en dataset. The colorbar indicates #Iterations in training. Details in Section 5.2.

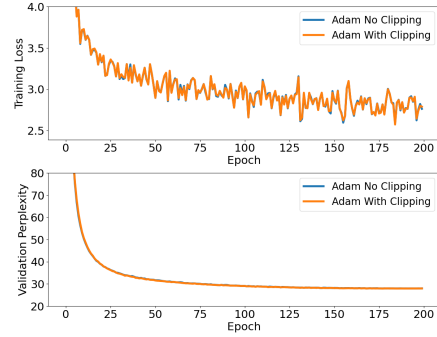


Figure 2: Training GPT-2 on Wikitext-103 using Adam with or without gradient clipping.

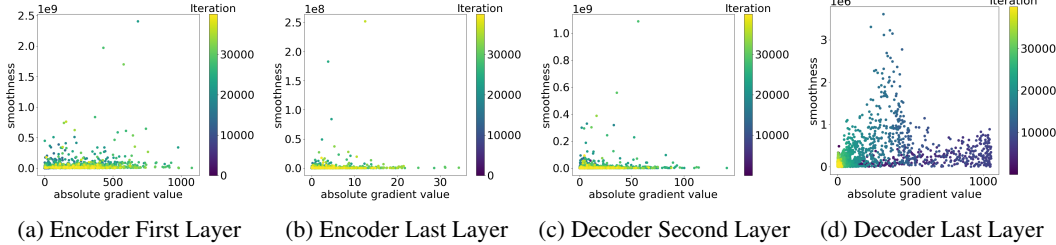


Figure 3: Local gradient Lipschitz constant vs. absolute gradient value on training a Transformer on WMT’16 Multimodal Translation de-en dataset. Each figure represents a randomly picked coordinate in corresponding layers. The colorbar indicates #Iterations during training. Details in Section 5.2.

36 Nevertheless, the (L_0, L_1) condition has not yet been empirically verified beyond LSTMs. Therefore,
 37 our **first contribution** lies in studying the applicability and generalization of the (L_0, L_1) condition.
 38 In particular, we have empirically verified that the popular Transformer [46] model also seems to
 39 satisfy this assumption, see Figure 1. Yet, we noticed that different coordinates, especially when they
 40 are in different layers of the model, exhibit very distinct L_0 and L_1 values as shown in Figure 3. Hence,
 41 we propose to refine the (L_0, L_1) assumption in (1) to a coordinate-wise version (Assumption 2) and
 42 consider this to better capture the loss surface when training deep neural networks like Transformers.

43 Given that we assume (a generalization) of the (L_0, L_1) assumption, it would be natural to use some
 44 clipping procedure. However, we found out that the use of clipping on Adam [24], while carried out
 45 in common practice [e.g., 49], *has no effect on the training and testing performance on optimizing*
 46 *a large transformer model as shown in Figure 2*. In retrospect, this might not be surprising: It is
 47 known that Adam has an implicit clipping behavior due to the normalization by the estimated second
 48 moment of the gradients. Indeed, Adam can be interpreted as a variant of SignSGD [2].

49 Inspired by this, our **second contribution** is to propose and analyze a *generalized SignSGD* algorithm
 50 under the relaxed smoothness assumption. It is parameterized in such a way that it on one end
 51 recovers SignSGD while on the other end closely resembles Adam. Apart from the convergence rates,
 52 we also located the critical role the momentum plays in **analyzing Adam-type algorithms**: it not only
 53 reduces the effects of noise but also gives an exponential decaying effect on the unbounded gradient
 54 norms and smoothness. This can partly explain the phenomenon that clipping does not help Adam.

55 The structure of this paper is as follows. Section 2 discusses related works and how our paper builds
 56 upon and distinguishes from them. The settings and assumptions are carried out in Section 3. We will
 57 introduce formally the generalized SignSGD algorithm and its analysis in Section 4, with a detailed
 58 discussion on the bounds and the role of momentum. The experimental results are shown in Section 5,
 59 comparing our algorithm with some popular competitors in deep learning tasks. Finally, we draw
 60 some conclusions and discuss the limitations of our work in Section 6.

61 **Notations** We will use $[d]$ to denote the sequence $[1, 2, \dots, d]$ and use bold letters to represent vectors,
 62 e.g., $\mathbf{u} \in \mathbb{R}^d$. The j^{th} coordinate of a vector \mathbf{u} is u_j . Throughout this paper, we study the Euclidean
 63 space \mathbb{R}^d with the inner product $\langle \cdot, \cdot \rangle$. The gradient of F at \mathbf{x} is denoted by $\nabla F(\mathbf{x})$. We use $\mathbb{I}(\cdot)$
 64 to denote the indicator function, $\|\mathbf{u}\|_p$ to denote the p -norm: $\|\mathbf{u}\|_p := (\sum_{j=1}^d |u_j|^p)^{1/p}$ and $\|\mathbf{u}\|_\infty$ the
 65 **maximum norm**: $\|\mathbf{u}\|_\infty := \max\{|u_1|, \dots, |u_d|\}$. We also denote by $\sum_{k=i}^j x_k = 0$ when $i > j$.

66 2 Related Works

67 **Adaptive Gradient Methods** Adaptive gradient methods [32, 11, 24, 18, 40] are popular optimizers
 68 for training deep neural networks. The traditional analysis of adaptive gradient methods is providing
 69 regret bounds under the online convex optimization framework [11, 24, 40]. Recently, there are
 70 some analysis of adaptive gradient methods for nonconvex smooth functions [6, 5, 50, 9, 55]. Zou
 71 et al. [54] introduces an intriguing connection between Adam [24] and SignGD [3] when training
 72 a two-layer neural network in the deterministic setting, where SignGD is an algorithm following
 73 the negative gradient sign direction to perform the update. However, these works cannot be directly
 74 extended to nonconvex functions with unbounded smoothness in the stochastic setting. To the best of
 75 our knowledge, this work is the first one establishing guarantees for **coordinate-wise type optimizers**
 76 **like** generalized SignSGD as well as Adam-type updates under a relaxed smoothness condition.

77 **Gradient Clipping** The algorithm and analysis of gradient clipping can be traced back to [1, 44, 12]
 78 under the assumption that the function is convex and rapidly growing. Hazan et al. [16] considered
 79 gradient clipping in quasi-convex optimization. Mai and Johansson [30] showed the stability and
 80 convergence of stochastic gradient clipping algorithms for convex problems without the smoothness
 81 condition. Gradient clipping is a standard technique in training deep neural networks [35, 36]
 82 such as RNNs and LSTMs. The theoretical analysis of gradient clipping for nonconvex models is
 83 pioneered by [52], in which the authors analyzed the convergence of gradient clipping under the
 84 relaxed smoothness assumption rather than the standard smoothness assumption. Zhang et al. [51]
 85 further improved the convergence rate bound under the same assumption as in [52]. Gradient clipping
 86 is also used when there is a heavy tail noise in the stochastic gradient to establish high probability
 87 convergence rates [8, 14, 53]. Cutkosky and Mehta [7] proved that normalized momentum improves
 88 normalized SGD under a second-order smoothness condition. **A close algorithm is the one in [21]**
 89 **which employs gradient normalization, momentum, and no gradient clipping to tackle the (L_0, L_1)**
 90 **condition (1) and control noise. Yet, their algorithm normalizes each coordinate with the same scale**
 91 **unlike popular optimizers such as Adam [24]. Moreover, we observe empirically that normalized**
 92 **SGD with momentum performs worse than Adam. Motivated by this, we propose a coordinate-wise**
 93 **optimization algorithm which requires new analysis tools compared with [21].**

94 **Employ m_t^2 to compute v_t in Adam** Designed to combine the advantages of Adagrad [11] and
 95 RMSProp [45], the update of Adam [24] employs the ratio between the exponential moving average
 96 of the stochastic gradient (m_t) and the exponential moving average of the squared stochastic gradient
 97 (v_t). Many variants of Adam have been proposed ever since. Among them, one idea is to use m_t^2
 98 to compute v_t instead of g_t^2 . The intuition is that m_t represents a better update direction than g_t
 99 and can thus better capture the second-moment information. Reddi et al. [39] adopted this change to
 100 prove the convergence of Adam in a federated learning setting; yet, they only consider the smooth
 101 setting and require a large ϵ to obtain convergence in contrast to the original Adam. Later, Wang et
 102 al. [47] explored this idea in more detail, but their analyses are still restricted to the smooth setting.

103 3 Settings and Preliminaries

104 In this paper, we focus on the following stochastic optimization problem:

$$\min_{\mathbf{x} \in \mathbb{R}^d} F(\mathbf{x}) := \mathbb{E}_{\xi \sim \mathcal{D}}[f(\mathbf{x}, \xi)],$$

105 where ξ is a random variable representing a randomly selected data sample or random noise following
 106 an unknown distribution \mathcal{D} . We will use the following assumptions.

107 **Assumption 1.** $F : \mathbb{R}^d \rightarrow \mathbb{R}$ is differentiable and bounded from below with infimum F^* .

108 **Assumption 2.** We say that a twice differentiable function $F(\mathbf{x})$ is $(\mathbf{L}_0, \mathbf{L}_1)$ -smooth coordinate-
 109 wisely, if for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ for which $\|\mathbf{x} - \mathbf{y}\|_2 \leq \frac{1}{\|\mathbf{L}_1\|_\infty}$, we have for any $j \in [d]$ that

$$\left| \frac{\partial F}{\partial x_j}(\mathbf{y}) - \frac{\partial F}{\partial x_j}(\mathbf{x}) \right| \leq \left(\frac{L_{0,j}}{\sqrt{d}} + L_{1,j} \left| \frac{\partial F}{\partial x_j}(\mathbf{x}) \right| \right) \|\mathbf{y} - \mathbf{x}\|_2. \quad (2)$$

110 We will denote $\mathbf{L}_0 := [L_{0,1}, L_{0,2}, \dots, L_{0,d}]^T$ and $\mathbf{L}_1 := [L_{1,1}, L_{1,2}, \dots, L_{1,d}]^T$.

111 The original (L_0, L_1) smoothness assumption (1) in [52] was proposed as a generalization of the
 112 more common smoothness assumption, which says that the gradient should be Lipschitz. Indeed,

113 when $L_{1,j}$ are zero, we recover the smoothness assumption. In contrast, when $L_{1,j}$ are non-zero,
 114 the smoothness of the function is potentially *unbounded*. Yet, [52] works with norms and applies
 115 to the global scale, while ours is more fine-grained and applies to each coordinate separately. One
 116 motivation for this assumption comes from [Remark 2.3, 51] where they noted that (1) can be relaxed
 117 to an assumption on gradient differences: $\exists K_0, K_1 > 0$ s.t. for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ with $\|\mathbf{x} - \mathbf{y}\|_2 \leq \frac{1}{K_1}$,

$$\|\nabla F(\mathbf{x}) - \nabla F(\mathbf{y})\|_2 \leq (K_0 + K_1 \|\nabla F(\mathbf{x})\|_2) \|\mathbf{x} - \mathbf{y}\|_2 \quad (3)$$

118 Indeed, our Assumption 2 implies (3) when $L_{0,j} = L_0$ and $L_{1,j} = L_1$ for all $j \in [d]$, up to constants
 119 (See Lemma 12 in the Appendix). Note that the $\frac{1}{\sqrt{d}}$ factor in ours is exactly for easy comparison
 120 with (3). The reason we turn to the current coordinate-wise version is that we observed a vast variance
 121 across different layers in training Transformer models: (1) is still true globally (Figure 1), but each
 122 layer or even each coordinate satisfies a very different (L_0, L_1) pair (Figure 3). The smoothness
 123 assumption has been generalized in orthogonal directions in other work [41, 3, 23].

124 One merit of Assumption 2 is that it gives us the following descent lemma.

125 **Lemma 1.** *Let F be (L_0, L_1) -smooth coordinate-wisely. Then, for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ for which*
 126 *$\|\mathbf{x} - \mathbf{y}\|_2 \leq \frac{1}{\|L_1\|_\infty}$, we have*

$$F(\mathbf{y}) \leq F(\mathbf{x}) + \langle \nabla F(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \sum_{j=1}^d \frac{\left(\frac{L_{0,j}}{\sqrt{d}} + L_{1,j} \left| \frac{\partial F}{\partial x_j}(\mathbf{x}) \right| \right) \|\mathbf{y} - \mathbf{x}\|_2}{2} |y_j - x_j|.$$

127 Our last assumption is common in the literature studying the (L_0, L_1) smooth condition [52, 51].

128 **Assumption 3.** *For each $j \in [d]$, there exists $\sigma_j > 0$ such that for all $\mathbf{x} \in \mathbb{R}^d$ and $\xi \sim \mathcal{D}$, the noise*
 129 *satisfies $\left| [\nabla f(\mathbf{x}, \xi)]_j - \frac{\partial F}{\partial x_j}(\mathbf{x}) \right| \leq \sigma_j$ with probability 1. We will denote $\boldsymbol{\sigma} := [\sigma_1, \sigma_2, \dots, \sigma_d]^T$.*

130 4 A Generalized SignSGD Algorithm

Algorithm 1 Generalized SignSGD (All operations on vectors are element-wise.)

```

1: Inputs:  $\mathbf{x}_1, \beta_1, \beta_2, \eta$ 
2:  $\mathbf{m}_0 = 0, \mathbf{v}_0 = 0$ 
3: for  $t = 1, \dots, T$  do
4:   Compute an unbiased estimate  $\nabla f(\mathbf{x}_t, \xi_t)$  of  $\nabla F(\mathbf{x}_t)$ , denoted as  $\mathbf{g}_t$ 
5:    $\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t$ 
6:    $\mathbf{v}_t = \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) \mathbf{m}_t^2$ 
7:    $\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\mathbf{m}_t}{\sqrt{\mathbf{v}_t}}$ 
8: end for

```

131 In this section, we present in Algorithm 1 a generalized SignSGD algorithm. This algorithm
 132 encompasses a variety of optimization algorithms.

133 At first sight, it seems very similar to Adam. Indeed, if we employ \mathbf{g}_t^2 in computing \mathbf{v}_t instead of \mathbf{m}_t^2 ,
 134 then it is exactly Adam, except for the bias correction terms. We would like to clarify that the idea
 135 of this change has been explored before, as detailed in Section 2. In this paper, the motivation for
 136 adopting this idea comes from the known effect of momentum on reducing the influence of noises [7].
 137 Indeed, in our analysis the difference between \mathbf{m}_t and $\nabla F(\mathbf{x}_t)$ is much more controllable than
 138 between \mathbf{g}_t and $\nabla F(\mathbf{x}_t)$. Thus, we consider employing \mathbf{m}_t in computing \mathbf{v}_t a better choice.

139 On the other end, the careful reader might observe that Algorithm 1 recovers the SignSGD with
 140 Momentum algorithm, also called SIGNUM in [3], when setting $\beta_2 = 0$. Sign-based algorithms
 141 are naturally suited to distributed learning [28] and the idea dated back to at least Rprop [42]. The
 142 convergence to a stationary point (with ℓ_1 norm) under a coordinate-wise smoothness condition
 143 has been established for SignSGD with/without the momentum in [3] though they necessitate large
 144 mini-batches to control the variance of the noise. Yet, we are more interested in their property of the
 145 update size being bounded without the need for explicit clipping.

146 Note that both SignSGD and Adam are good candidates for optimization algorithms whose update
 147 must be bounded on functions that satisfy the $(\mathbf{L}_0, \mathbf{L}_1)$ condition. Indeed, SignSGD can be seen as
 148 an extreme form of gradient clipping. On the other hand, as said in the introduction, Adam does not
 149 seem to require gradient clipping at all when used to train the large Transformer model in Figure 2.

150 Hence, we expect our algorithm, a generalization of SignSGD and a close resemblance to Adam, can
 151 enjoy the merits of both and be robust to the unbounded smoothness in the $(\mathbf{L}_0, \mathbf{L}_1)$ scenario. In the
 152 next section, we will formalize this claim by presenting the theoretical analysis of Algorithm 1.

153 4.1 Theoretical Convergence Analysis

154 **Theorem 1.** *Under Assumptions 1, 2, and 3, assume $M_j := \sup \left\{ \left| \frac{\partial F}{\partial x_j}(\mathbf{x}) \right| : F(\mathbf{x}) \leq F(\mathbf{x}_1) \right\}$ is
 155 finite for each $j \in [d]$, let Δ be any upper bound on $F(\mathbf{x}_1) - F^*$, $\alpha = \min \left(\frac{\sqrt{\|\mathbf{L}_0\|_1 \sqrt{\Delta}}}{\|\boldsymbol{\sigma}\|_1 \sqrt{T}}, 1 \right)$, $\beta_1 =$
 156 $1 - \alpha$, $\frac{\sqrt{\beta_2}}{\beta_1} < 1$, $\rho = 1 - \frac{\sqrt{\beta_2}}{\beta_1}$, $\eta = \frac{\sqrt{\Delta \alpha}}{\sqrt{\|\mathbf{L}_0\|_1 \sqrt{T}}}$, for $T \geq \max \left(\frac{100d\Delta \|\mathbf{L}_1\|_\infty^2}{(1-\beta_2)\rho^2 \|\mathbf{L}_0\|_1}, \frac{10000d^2 \Delta \|\boldsymbol{\sigma}\|_1^2 \|\mathbf{L}_1\|_\infty^4}{(1-\beta_2)^2 \rho^4 \|\mathbf{L}_0\|_1^3} \right)$,
 157 Algorithm 1 guarantees, with probability at least $1 - \delta$, that*

$$\min_{t \in [T]} \|\nabla F(\mathbf{x}_t)\|_1 = \mathcal{O} \left(\frac{\sqrt{\log(dT/\delta)} \|\mathbf{L}_0\|_1^{1/4} \Delta^{1/4} \|\boldsymbol{\sigma}\|_1^{1/2}}{\rho \sqrt{1 - \beta_2} T^{1/4}} + \frac{\log(dT/\delta) \sqrt{\|\mathbf{L}_0\|_1 \Delta}}{\rho \sqrt{T}} \right) \\ + \mathcal{O} \left(\frac{\|\mathbf{M}\|_1 + \|\boldsymbol{\sigma}\|_1}{\rho} \exp \left(-\frac{\sqrt{1 - \beta_2} \|\mathbf{L}_0\|_1^{3/4}}{\sqrt{d} \|\mathbf{L}_1\|_\infty \|\boldsymbol{\sigma}\|_1^{1/2} \Delta^{1/4}} T^{1/4} \right) + \frac{\|\nabla F(\mathbf{x}_1)\|_1}{T} \right).$$

158 Furthermore, for the case when $\beta_2 = 0$, we have the following refined guarantee:

$$\min_{t \in [T]} \|\nabla F(\mathbf{x}_t)\|_1 = \mathcal{O} \left(\frac{\sqrt{\log(dT/\delta)} \|\mathbf{L}_0\|_1^{1/4} \Delta^{1/4} \|\boldsymbol{\sigma}\|_1^{1/2}}{T^{1/4}} + \frac{\log(dT/\delta) \sqrt{\|\mathbf{L}_0\|_1 \Delta}}{\sqrt{T}} \right) \\ + \mathcal{O} \left(\frac{\|\nabla F(\mathbf{x}_1)\|_1}{\sqrt{T}} \left(\frac{1}{\sqrt{T}} + \frac{\|\boldsymbol{\sigma}\|_1}{\sqrt{\|\mathbf{L}_0\|_1 \Delta}} \right) + \frac{\|\boldsymbol{\sigma}\|_1}{T} \right).$$

159 Here, M_j denotes the maximum absolute value of the partial derivative of F for coordinate j among
 160 the sub-level set of $F(\mathbf{x}_1)$, namely any point \mathbf{x} with $F(\mathbf{x}) \leq F(\mathbf{x}_1)$. In other words, we assume
 161 gradients to be bounded in the sub-level set of $F(\mathbf{x}_1)$; yet, we do not make any restriction on gradients
 162 outside of this set. We believe this is not a strong assumption, for example, when the sub-level
 163 set of $F(\mathbf{x}_1)$ is bounded, then by the assumed continuity of gradients it trivially holds. Also, we
 164 just require an upper bound and it can even be exponentially large as we have an exponentially
 165 decaying coefficient to counteract it: notice how the term $\|\mathbf{M}\|_1$ is multiplied by a term that decays
 166 exponentially with T . Better still, when $\beta_2 = 0$, we no longer even need this assumption and the
 167 algorithm is entirely free of the influence of $\|\mathbf{M}\|_1$. To see why this is good, we show a refined lower
 168 bound of Gradient Descent under the relaxed smoothness scenario below which is originally in [52].

169 **Theorem 2.** *Fix $\epsilon > 0$, $L_0 > 0$, $L_1 > 0$, $M \geq \max(\frac{L_0}{L_1}, \epsilon)$, and $x_0 \in \mathbb{R}$. Pick any constant learning
 170 rate η for GD, with the knowledge of the above constants. Then, there exists a 1-d (L_0, L_1) -smooth
 171 function, bounded from below by f^* (finite), and such that $\sup\{|f'(x)| : f(x) \leq f(x_0)\} \leq M$ on
 172 which the number of iterations T of GD with learning rate η to guarantee $|f'(x_T)| < \epsilon$ is at least*

$$ML_1(f(x_0) - f^* - \frac{15\epsilon^2}{16L_0}) / 2\epsilon^2 \left(\ln \frac{ML_1}{L_0} + 1 \right).$$

173 Theorem 2 shows that in the relaxed smoothness setting, GD with any constant step size will suffer
 174 from a linear term depending on $L_1 M$. On a side note, it is a fixed version of the lower bound in
 175 [52]: we provide in Appendix an explanation of errors in their lower bound and our corrected proof.

176 Compared with GD, our algorithm only has an exponentially decaying dependence on $L_1 M$. We
 177 consider this to be substantial merit of our algorithm. Furthermore, when $\beta_2 = 0$ in which case we
 178 recover the SignSGD with Momentum algorithm, we can even show that it completely removes the
 179 effects of the unbounded gradient norms. Also notice that in such case we actually no longer need the
 180 assumption of $M_j := \sup \left\{ \left| \frac{\partial F}{\partial x_j}(\mathbf{x}) \right| : F(\mathbf{x}) \leq F(\mathbf{x}_1) \right\}$ being finite for each $j \in [d]$ anymore, and
 181 the $\|\mathbf{L}_1\|_\infty$ term does not appear in the final bound anymore.

182 We also would like to point out that this bound closely resembles the one achieved by SGD with
 183 gradient clipping algorithm [51] except that we consider the coordinate-wise setting: take the setting
 184 of $\beta_2 = 0$ for example, we need at most $\mathcal{O}\left(\Delta \max\left\{\frac{\|\sigma\|_1^2 \|\mathbf{L}_0\|_1}{\epsilon^4}, \frac{d^2 \|\sigma\|_1^2 \|\mathbf{L}_1\|_\infty^4}{\|\mathbf{L}_0\|_1^3}, \frac{d \|\mathbf{L}_1\|_\infty^2}{\|\mathbf{L}_0\|_1}\right\}\right)$ to get a
 185 point \mathbf{x} with $\|\nabla F(\mathbf{x})\|_1 \leq \epsilon$ with high probability.

186 **Remark 1** The almost surely bounded assumption 3 can be relaxed to sub-gaussian noise, using
 187 standard extensions of Freedman inequality [e.g., 15].

188 **Remark 2** When $\beta_2 = 0$, we can prove an average-iterate complexity bound (see Proof of Theorem 1
 189 for $\beta_2 = 0$ in Appendix A.3); yet, we use the min form for consistency between the two cases.

190 **Remark 3** Our bound is incomparable with the one in [51, Theorem 3.2]. Yet, as we said, if
 191 $L_{0,j} = L_0$ and $L_{1,j} = L_1$ for all $j \in [d]$, then the function satisfies (3). In this case, assuming the
 192 noise vector and the gradient vector to be dense to be able to compare the ℓ_1 -norm and the ℓ_2 -norm,
 193 we recover the same bound of [51, Theorem 3.2] in terms of dependencies on L_1 , L_0 , and T . Instead,
 194 in the more general case when $L_{0,j}$ and $L_{1,j}$ are not uniform vectors, our bound allows a finer control
 195 of the unbounded smoothness.

196 The proof of the theorem is highly technical and it uses recent advancements in the analysis of
 197 momentum methods [7], key techniques to deal with the (L_0, L_1) assumption [51], as well as a novel
 198 and essential inductive argument to control the norm of past gradients. **We want to stress that the**
 199 **difficulty mainly comes from analyzing Adam-type updates when $\beta_2 > 0$, while for the other case of**
 200 **$\beta_2 = 0$ the proof is significantly simpler.** The full proof is in the Appendix, but here we present a
 201 proof sketch that underlines the main steps. First, we list some key lemmas we used but move their
 202 proofs to the appendix due to space constraints.

203 **Lemma 4.** *With notations in Algorithm 1, for $\tau \leq \bar{\tau} = \frac{\sqrt{1-\beta_2}}{\eta\sqrt{d}\|\mathbf{L}_1\|_\infty}$, we have $\|\mathbf{x}_{t-\tau} - \mathbf{x}_t\|_2 \leq \frac{1}{\|\mathbf{L}_1\|_\infty}$.*
 204 Lemma 4 limits our focus to the most recent $\bar{\tau}$ steps on which Assumption 2 and Lemma 1 can apply.

205 **Lemma 6.** *Assume Assumption 3. With the notation of Algorithm 1, let $j \in [d]$ and $\beta_1 \leq 1$. Then,*
 206 *with probability at least $1 - 3\delta$, for any $t_0 \in [t]$, we have*

$$\left| \sum_{\tau=1}^{t_0} \beta_1^{t-\tau} \left(g_{\tau,j} - \frac{\partial F}{\partial x_j}(\mathbf{x}_\tau) \right) \right| \leq 3\sigma_j \max(1, \log(1/\delta)) + \frac{3}{\sqrt{1-\beta_1^2}} \sqrt{\sigma_j^2 \max(1, \log(1/\delta))} \triangleq E_j.$$

207 Lemma 6 is the major tool we use to handle the noise we incur during drawing stochastic gradients.
 208 It is derived based on Lemma 12 in [8].

209 **Lemma 9.** *With the notation of Algorithm 1 and under the assumptions of Theorem 1, if $\left| \frac{\partial F}{\partial x_j}(\mathbf{x}_\tau) \right| \leq$
 210 M_j holds for all $\tau \leq t$ and $j \in [d]$, and $D > 0$, then, with probability at least $1 - 3t\delta$ we have that,*

$$\text{either } \left| \frac{\partial F}{\partial x_j}(\mathbf{x}_t) \right| < \frac{5B_j}{D} \text{ or } \frac{|m_{t,j}|}{\sqrt{v_{t,j}}} \geq \frac{\rho D}{5\sqrt{1-\beta_2}},$$

211 where $B_j \triangleq \frac{\eta L_{0,j}}{\sqrt{1-\beta_2}(1-\beta_1)} + \beta_1 \bar{\tau} (M_j + \sigma_j) + (1-\beta_1)E_j$ and $D \triangleq 1 - \frac{2\eta\sqrt{d}\|\mathbf{L}_1\|_\infty}{\sqrt{1-\beta_2}(1-\beta_1)}$.

212 Lemma 9 is similar to Lemma A.2 in [54] which considered the deterministic and smooth setting;
 213 in contrast, our proof is much more challenging in that we need to tackle both the noise and the
 214 unbounded smoothness. With this lemma, we know that either the true gradient is small or that the
 215 update of our Algorithm 1 can be lower bounded.

216 **Lemma 11.** *Under Assumptions 1, 2, and 3, using the hyperparameters in Theorem 1, denoting*
 217 *$\alpha = 1 - \beta_1$ and $\epsilon_t = \mathbf{m}_t - \nabla F(\mathbf{x}_t)$, for all t and $j \in [d]$ we have, with probability at least $1 - 3\delta$,*

$$|\epsilon_{t+1,j}| \leq (1-\alpha)^t \left(\alpha\sigma_j + (1-\alpha) \left| \frac{\partial F}{\partial x_j}(\mathbf{x}_1) \right| \right) + \frac{\eta L_{0,j}}{\alpha} + \alpha E_j + (1-\alpha)\eta\sqrt{d}L_{1,j} \sum_{\tau=0}^{t-1} (1-\alpha)^\tau \left| \frac{\partial F}{\partial x_j}(\mathbf{x}_{t-\tau}) \right|.$$

218 Lemma 11 shows how the use of momentum can help control the noise by choosing β_1 wisely. It is
 219 adapted from the proof of Theorem 2 in [8] but with the added difficulty of unbounded smoothness.

220 *Proof sketch of Theorem 1.* Observing the formula of setting β_1 , we can see that when $\|\sigma\|_1 \leq$
 221 $\sqrt{\|\mathbf{L}_0\|_1 \Delta} / \sqrt{T}$, $\beta_1 = 0$. As $\beta_2 < \beta_1$, Algorithm 1 reduces to SignSGD. In this case, the key
 222 component is Lemma 11 using which we are able to show that $\sum_{t=1}^T \left| m_{t,j} - \frac{\partial F}{\partial x_j}(\mathbf{x}_t) \right|$ can be

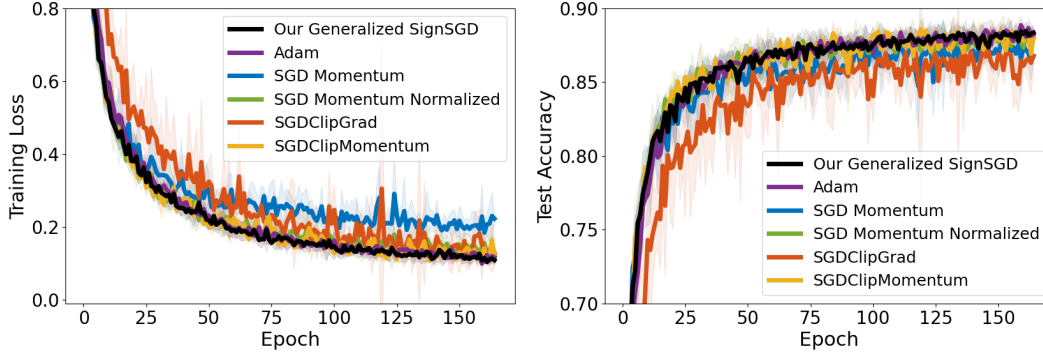


Figure 4: Training a 20-layer Resnet on CIFAR10. The shading of each curve represents the 95% confidence interval computed across 5 independent runs from different random seeds.

223 controlled as $C_1 \sum_{t=1}^T \left| \frac{\partial F}{\partial x_j}(\mathbf{x}_t) \right| + C_2$. The summation of true gradients over time can then be
 224 offsetted by choosing η and β_1 wisely when we invoke the descent lemma 1. The rest is standard.

225 Now for the other case in which $\|\sigma\|_1 > \sqrt{\|\mathbf{L}_0\|_1 \Delta} / \sqrt{T}$, we take a different route.

226 First, notice that Assumption 2 and the Descent Lemma 1 only hold when two points are not too far
 227 away. Thus, we need to restrict our attention to the recent updates (Lemma 4), beyond which we
 228 would have no control. This means we want the influence of those updates too long ago to not have
 229 a big effect on the current one. To make this happen, one natural idea is to use a bounded gradient
 230 assumption, then with the use of exponential averaging, their effect would be quickly reduced. Yet,
 231 assuming directly that all gradients are bounded would trivialize the $(\mathbf{L}_0, \mathbf{L}_1)$ assumption. Thus,
 232 we pose a much weaker condition, assuming that $M_j := \sup \left\{ \left| \frac{\partial F}{\partial x_j}(\mathbf{x}) \right| : F(\mathbf{x}) \leq F(\mathbf{x}_1) \right\}$ being
 233 finite for each $j \in [d]$. Then, we prove that M_j will provide an upper bound to all the true gradients
 234 the algorithm see. We prove it using induction, analyzing separately the case that either the true
 235 gradient is already very small and we have reached the proximity of a stationary point, or the objective
 236 function is monotonically non-increasing and the gradient remains bounded.

237 Having controlled the past gradients, we prove in Lemma 9 that the update of Algorithm 1 is either
 238 very small that we can pass or having a constant lower bound that we can use in the Descent Lemma 1.

239 Also, considering that this is the stochastic setting, noise typically slows down convergence or can
 240 even cause the algorithm to diverge if the hyperparameters are not chosen wisely. To handle this, we
 241 invoke Freedman’s inequality to show that the addition of adjacent stochastic noise almost cancels
 242 out each other and the absolute value of the sum remains controlled (Lemma 6).

243 Yet, we still need another block to handle the difference between the true gradient and the momentum
 244 as we are updating in the direction of the momentum instead of the true gradient. Turns out that we
 245 can prove that $\text{sign}(m_{t,j}) = \text{sign} \left(\frac{\partial F}{\partial x_j}(\mathbf{x}_t) \right)$ when $\left| \frac{\partial F}{\partial x_j}(\mathbf{x}_t) \right|$ is not too small. As before, in the case
 246 $\left| \frac{\partial F}{\partial x_j}(\mathbf{x}_t) \right|$ is small, we have converged on that coordinate.

247 Combining all these blocks together, we are able to arrive at the final results. □

248 5 Experiments

249 We conducted our experiments using PyTorch [37] on Nvidia V100 GPUs.

250 5.1 Comparison with Other Optimizers

251 To validate the efficacy of our Algorithm 1, we compare it with Adam [24], SGD [43], **SGD Momen-**
 252 **tum Normalized** [21], SGDClipGrad, and SGDClipMomentum. The latter two are from Algorithm
 253 1 in [51] where SGDClipGrad corresponds to the case where $\nu = 0$ and SGDClipMomentum
 254 corresponds to the case when $\nu = 1$.

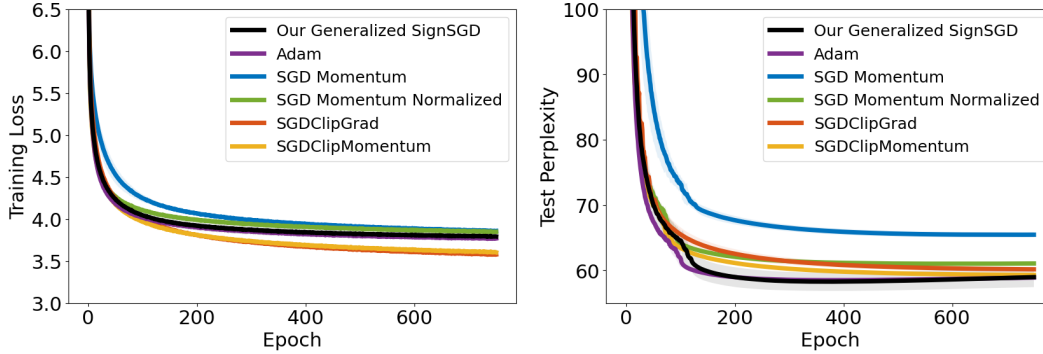


Figure 5: Training an AWD-LSTM to do language modeling (word level) on Penn Treebank. The shading of each curve represents the 95% confidence interval over 5 independent runs.

Table 1: Average final training loss and test accuracy achieved by each method when optimizing respective models on each dataset. The \pm shows 95% confidence intervals of the mean loss/accuracy/perplexity value over 5 runs starting from different random seeds.

Methods	CIFAR10		Penn Treebank	
	Training loss	Test accuracy	Training loss	Test perplexity
SGD Momentum	0.2226 \pm 0.0169	0.8674 \pm 0.0048	3.8587 \pm 0.0058	65.4622 \pm 0.3842
SGD Momentum Normalized	0.1262 \pm 0.0170	0.8795 \pm 0.0086	3.8487 \pm 0.0073	61.0558 \pm 0.3224
SGDClipGrad	0.1288 \pm 0.0403	0.8677 \pm 0.0106	3.5774 \pm 0.0081	60.1604 \pm 0.2797
SGDClipMomentum	0.1220 \pm 0.0162	0.8809 \pm 0.0022	3.6038 \pm 0.0102	59.3052 \pm 0.2798
Adam	0.1161 \pm 0.0111	0.8823 \pm 0.0041	3.7692 \pm 0.0062	58.9005 \pm 0.3058
Our Algorithm 1	0.1086 \pm 0.0129	0.8835 \pm 0.0032	3.7928 \pm 0.0425	58.9661 \pm 1.5218

255 **Training** Unless otherwise specified, we use grid-search to fine-tune the initial learning rate for all
 256 optimizers, as well as the clipping threshold for SGDClipGrad and SGDClipMomentum, and β_2 for
 257 Adam and our algorithm, to select the one giving the best validation performance on a separated
 258 validation set. We then employ the best performing hyperparameters to train the model over all
 259 training data and report the testing performance. The testing is repeated with random seeds 5 times to
 260 eliminate the influence of stochasticity. For more details, please refer to Section A.4.

261 **Resnet for Image Classification on CIFAR-10** We employ the 20-layer Residual Network model [17]
 262 to do image classification on the CIFAR-10 dataset. Images are normalized per channel using the
 263 means and standard deviations computed from all training images. We adopt the data augmentation
 264 technique following [27] (for training only): 4 pixels are padded on each side of an image and a $32 \times$
 265 32 crop is randomly sampled from the padded image or its horizontal flip. The mini-batch size is 128
 266 and we train all algorithms for 164 epochs. We do not employ any learning rate decay schedule in
 267 order to focus on the comparison of the optimizers themselves. We fixed the weight decay value to be
 268 0.0001 and the momentum parameter (β_1) to be 0.9. Figure 4 and Table 1 report the training and
 269 testing performance for each algorithm, showing that ours is among the best of all.

270 **LSTM for Language Modeling on Penn Treebank** We adopt a 3-layer AWD-LSTM [33] to do
 271 language modeling on the Penn Treebank (PTB) dataset [31](word level). The mini-batch size is 40
 272 and we trained each algorithm for 750 epochs. Apart from the hyperparameters we stated above, we
 273 further fine-tuned the weight decay value for all algorithms noticing its significant influence on the
 274 performance. We choose the set of hyperparameters that give the smallest final validation perplexity.
 275 We report the results in Figure 5 and Table 1. It can be seen that we can match the performance of
 276 Adam while beating the others.

277 5.2 Transformers Observe (L_0, L_1) -smoothness

278 For Figure 1 which verifies the original form (1) of the (L_0, L_1) condition using the norm, we
 279 followed the method in Section H.3 of [52]. Specifically, given \mathbf{x}_t and \mathbf{x}_{t+1} , denote $\mathbf{d} := \mathbf{x}_{t+1} - \mathbf{x}_t$.
 280 We estimate the smoothness at \mathbf{x}_t by

$$\hat{L}_t = \max_{\gamma \in \{\delta_1, \delta_2, \dots, \delta_N\}} \frac{\|\nabla F(\mathbf{x}_t + \gamma \mathbf{d}) - \nabla F(\mathbf{x}_t)\|_2}{\|\gamma \mathbf{d}\|_2},$$

281 where $\{\delta_1, \delta_2, \dots, \delta_N\}$ denotes the sample locations and we use $\{\frac{1}{6}, \frac{2}{6}, \frac{3}{6}, \frac{4}{6}, \frac{5}{6}\}$.

282 For Figure 3 verifying the coordinate-wise version (2) of the $(\mathbf{L}_0, \mathbf{L}_1)$ condition, note that the
283 equation is symmetric in that if we just swap \mathbf{x} and \mathbf{y} it shall still holds. Thus, during plotting, we
284 compare $\left| \frac{\partial F}{\partial x_j}(\mathbf{x}_{t+1}) - \frac{\partial F}{\partial x_j}(\mathbf{x}_t) \right| / |x_{t+1,j} - x_{t,j}|$ vs. $\min \left(\left| \frac{\partial F}{\partial x_j}(\mathbf{x}_t) \right|, \left| \frac{\partial F}{\partial x_j}(\mathbf{x}_{t+1}) \right| \right)$.

285 Figure 1(a) is on training a 2-layer Transformer Encoder to do language modeling on the Wikitext-2
286 dataset. The implementation, settings, and parameter choices follow this.¹ We only plot the first 5
287 training epochs. Figure 1(b) and 3 are on training a 6-layer Transformer [46] to do machine translation
288 on the WMT’16 Multimodal Machine Translation Task German-English dataset. The implementation
289 of the transformer is forked from here² and we also follow their default settings. The mini-batch size
290 is 256 and we trained for 400 epochs using Adam and report the whole training trajectory.

291 5.3 Clipping does not Affect Adam’s Performance

292 We compare clipping and non-clipping for Adam optimizer on the Wikitext-103 (103 million tokens,
293 180MB) [34] language modeling task, with a 16-layer GPT-2 transformer model [38]. This GPT-2
294 model has an input length of 256 tokens, 410-dimension word embedding, 16 Attention layers with
295 10 Attention heads and 2100 hidden dimensions. Model size is 201.58 MB. The vocabulary size is
296 28996. We use the hyper-parameter settings prescribed in [49]: batch size 256, warm up learning rate
297 from 0 to 2.5×10^{-4} in the first 64000 samples (i.e., 250 iterations) and then cosine-anneal learning
298 rate to zero, on top of an Adam optimizer. It takes about 40 hours to train 200 epochs on 8 V100
299 GPUs. We use clipping threshold max_norm 0.25 for the entire model as prescribed in the literature
300 [49]. We also count that with this clipping scheme, clipping occurs in every single batch. As we can
301 see from Figure 2, neither training loss (2.79 vs 2.76) nor perplexity score (27.92 vs 27.97) differs
302 much in the clipping and non-clipping case, which is consistent with our theory that Adam naturally
303 achieves gradients clipping effect.

304 6 Conclusion and Limitations

305 Smoothness has been a widely adopted condition for proving convergence rates of algorithms in
306 the non-convex optimization scenario. Yet, it has been found that this assumption does not capture
307 losses when employing some deep learning models including RNNs and LSTMs. In light of this, a
308 relaxed smoothness assumption was proposed that aligns well with the practice. We observed that
309 the loss surface of training using Transformers also exhibits this relaxed smoothness. Under this
310 assumption, SGD with clipped gradient has been proven to work well. However, we found that
311 clipping is not necessary for achieving convergence in such a setting. Indeed, we showed that a
312 generalized SignSGD algorithm does not require explicit clipping but can almost guarantee the same
313 bound as SGD with clipping. In the analyses, we identified the key effect of using momentum in
314 **analyzing Adam-type algorithms**, that it reduces both the noise and the unbounded gradient norms.
315 Finally, we conducted a variety of deep learning tasks showing that our algorithm can match Adam’s
316 performance while exceeding others.

317 **Limitations** The current work is in no way a perfect one and there are many directions worth
318 exploring beyond it. First of all, though our algorithm could be seen as a close resemblance to the
319 original Adam algorithm, they are still not equal. Considering the huge popularity of Adam and
320 its established effectivity in practice, it is worth studying whether Adam in its original form can
321 converge in the relaxed smooth setting. Second, **while our Theorem 1 are upper bounds and cannot be**
322 **directly compared between the two cases of β_2 , it does suggests that $\beta_2 = 0$ minimizes the worst-case**
323 **convergence rate. However, it still does not fully explain the phenomenon that a choice of β_2 close to**
324 **1 yields better performance in using our Algorithm 1 as well as Adam in practice. Third,** despite there
325 are lower bounds showing that, for example, GD with a constant step size can be arbitrarily worse
326 than GD with clipping, it would be more meaningful to study whether the relaxed smooth condition
327 is inherently more difficult, possibly by establishing a lower bound for all first-order optimization
328 algorithms. Finally, we did show that Transformers observe the relaxed smoothness condition, but we
329 consider it more beneficial to research in-depth what properties or structures make a model satisfy
330 such conditions.

¹https://pytorch.org/tutorials/beginner/transformer_tutorial.html

²<https://github.com/jadore801120/attention-is-all-you-need-pytorch>

331 **References**

- 332 [1] Ya I Alber, Alfredo N. Iusem, and Mikhail V. Solodov. On the projected subgradient method for
333 nonsmooth convex optimization in a Hilbert space. *Mathematical Programming*, 81(1):23–35,
334 1998.
- 335 [2] Lukas Balles and Philipp Hennig. Dissecting Adam: The sign, magnitude and variance of
336 stochastic gradients. In *International Conference on Machine Learning*, pages 404–413. PMLR,
337 2018.
- 338 [3] Jeremy Bernstein, Yu-Xiang Wang, Kamyar Azizzadenesheli, and Animashree Anandkumar.
339 signSGD: Compressed optimisation for non-convex problems. In *International Conference on*
340 *Machine Learning*, pages 560–569. PMLR, 2018.
- 341 [4] Stephen Boyd and Lieven Vandenbergh. *Convex optimization*. Cambridge university press,
342 2004.
- 343 [5] Xiangyi Chen, Sijia Liu, Ruoyu Sun, and Mingyi Hong. On the convergence of a class of
344 Adam-type algorithms for non-convex optimization. In *International Conference on Learning*
345 *Representations*, 2019.
- 346 [6] Zaiyi Chen, Zhuoning Yuan, Jinfeng Yi, Bowen Zhou, Enhong Chen, and Tianbao Yang.
347 Universal stagewise learning for non-convex problems with convergence on averaged solutions.
348 In *International Conference on Learning Representations*, 2019.
- 349 [7] Ashok Cutkosky and Harsh Mehta. Momentum improves normalized SGD. In *International*
350 *Conference on Machine Learning*, pages 2260–2268. PMLR, 2020.
- 351 [8] Ashok Cutkosky and Harsh Mehta. High-probability bounds for non-convex stochastic opti-
352 mization with heavy tails. *Advances in Neural Information Processing Systems*, 34, 2021.
- 353 [9] Alexandre Défossez, Léon Bottou, Francis Bach, and Nicolas Usunier. A simple convergence
354 proof of Adam and Adagrad. *arXiv preprint arXiv:2003.02395*, 2020.
- 355 [10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of
356 deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*,
357 2018.
- 358 [11] John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning
359 and stochastic optimization. *Journal of Machine Learning Research*, 12(Jul):2121–2159, 2011.
- 360 [12] Yuri Ermoliev. Stochastic quasigradient methods. In *Numerical techniques for stochastic*
361 *optimization*, number 10 in Springer Series in Computational Mathematics, pages 141–185.
362 Springer, 1988.
- 363 [13] Saeed Ghadimi and Guanghui Lan. Stochastic first-and zeroth-order methods for nonconvex
364 stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013.
- 365 [14] Eduard Gorbunov, Marina Danilova, and Alexander Gasnikov. Stochastic optimization with
366 heavy-tailed noise via accelerated gradient clipping. *arXiv preprint arXiv:2005.10785*, 2020.
- 367 [15] Nicholas JA Harvey, Christopher Liaw, Yaniv Plan, and Sikander Randhawa. Tight analyses for
368 non-smooth stochastic gradient descent. In *Conference on Learning Theory*, pages 1579–1613.
369 PMLR, 2019.
- 370 [16] Elad Hazan, Kfir Y Levy, and Shai Shalev-Shwartz. Beyond convexity: Stochastic quasi-convex
371 optimization. In *Advances in Neural Information Processing Systems*, 2015.
- 372 [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
373 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
374 pages 770–778, 2016.

- 375 [18] Geoffrey Hinton, Li Deng, Dong Yu, George Dahl, Abdel-rahman Mohamed, Navdeep Jaitly,
376 Andrew Senior, Vincent Vanhoucke, Patrick Nguyen, and Brian Kingsbury. Deep neural
377 networks for acoustic modeling in speech recognition. *IEEE Signal processing magazine*, 29,
378 2012.
- 379 [19] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*,
380 9(8):1735–1780, 1997.
- 381 [20] Prateek Jain and Purushottam Kar. Non-convex optimization for machine learning. *Found.*
382 *Trends Mach. Learn.*, 10:142–336, 2017.
- 383 [21] Jikai Jin, Bohang Zhang, Haiyang Wang, and Liwei Wang. Non-convex distributionally robust
384 optimization: Non-asymptotic analysis. In *Advances in Neural Information Processing Systems*,
385 volume 34, pages 2771–2782. Curran Associates, Inc., 2021.
- 386 [22] Michael Kamp, Linara Adilova, Joachim Sicking, Fabian Hüger, Peter Schlicht, Tim Wirtz, and
387 Stefan Wrobel. Efficient decentralized deep learning by dynamic model averaging. In *Joint*
388 *European Conference on Machine Learning and Knowledge Discovery in Databases*, pages
389 393–409. Springer, 2018.
- 390 [23] Ahmed Khaled and Peter Richtárik. Better theory for SGD in the nonconvex world. *arXiv*
391 *preprint arXiv:2002.03329*, 2020.
- 392 [24] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Internat-*
393 *ional Conference on Learning Representations*, 2015.
- 394 [25] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep
395 convolutional neural networks. In *Advances in neural information processing systems*, pages
396 1097–1105, 2012.
- 397 [26] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436–444,
398 2015.
- 399 [27] Chen-Yu Lee, Saining Xie, Patrick Gallagher, Zhengyou Zhang, and Zhuowen Tu. Deeply-
400 supervised nets. In *Artificial intelligence and statistics*, pages 562–570. PMLR, 2015.
- 401 [28] Mu Li, David G Andersen, Jun Woo Park, Alexander J Smola, Amr Ahmed, Vanja Josifovski,
402 James Long, Eugene J Shekita, and Bor-Yiing Su. Scaling distributed machine learning with
403 the parameter server. In *OSDI*, volume 14, pages 583–598, 2014.
- 404 [29] Xiaoyu Li and Francesco Orabona. On the convergence of stochastic gradient descent with
405 adaptive stepsizes. In *The 22nd International Conference on Artificial Intelligence and Statistics*,
406 pages 983–992. PMLR, 2019.
- 407 [30] Vien V Mai and Mikael Johansson. Stability and convergence of stochastic gradient clipping:
408 Beyond lipschitz continuity and smoothness. In *International Conference on Machine Learning*,
409 pages 7325–7335. PMLR, 2021.
- 410 [31] Mitchell P. Marcus, Mary Ann Marcinkiewicz, and Beatrice Santorini. Building a large
411 annotated corpus of English: The Penn Treebank. *Comput. Linguist.*, 19(2):313–330, June
412 1993.
- 413 [32] H Brendan McMahan and Matthew Streeter. Adaptive bound optimization for online convex
414 optimization. *arXiv preprint arXiv:1002.4908*, 2010.
- 415 [33] Stephen Merity, Nitish Shirish Keskar, and Richard Socher. Regularizing and optimizing LSTM
416 language models. In *International Conference on Learning Representations*, 2018.
- 417 [34] Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture
418 models. *ICLR*, 2017.
- 419 [35] Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. Understanding the exploding gradient
420 problem. corr abs/1211.5063 (2012). *arXiv preprint arXiv:1211.5063*, 2012.

- 421 [36] Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent
422 neural networks. In *International conference on machine learning*, pages 1310–1318. PMLR,
423 2013.
- 424 [37] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,
425 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative
426 style, high-performance deep learning library. In *Advances in Neural Information Processing*
427 *Systems*, pages 8024–8035, 2019.
- 428 [38] Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language
429 models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 430 [39] Sashank J. Reddi, Zachary Charles, Manzil Zaheer, Zachary Garrett, Keith Rush, Jakub Konečný,
431 Sanjiv Kumar, and Hugh Brendan McMahan. Adaptive federated optimization. In *International*
432 *Conference on Learning Representations*, 2021.
- 433 [40] Sashank J Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of Adam and beyond. In
434 *International Conference on Learning Representations*, 2018.
- 435 [41] Peter Richtárik and Martin Takáč. Iteration complexity of randomized block-coordinate descent
436 methods for minimizing a composite function. *Mathematical Programming*, 144:1–38, 2014.
- 437 [42] Martin Riedmiller and Heinrich Braun. A direct adaptive method for faster backpropagation
438 learning: The rprop algorithm. In *IEEE international conference on neural networks*, pages
439 586–591. IEEE, 1993.
- 440 [43] Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of*
441 *mathematical statistics*, pages 400–407, 1951.
- 442 [44] Naum Zuselevich Shor. *Minimization methods for non-differentiable functions*, volume 3.
443 Springer Science & Business Media, 2012.
- 444 [45] Tijmen Tieleman, Geoffrey Hinton, et al. Lecture 6.5-rmsprop: Divide the gradient by a
445 running average of its recent magnitude. *COURSERA: Neural networks for machine learning*,
446 4(2):26–31, 2012.
- 447 [46] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
448 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information*
449 *processing systems*, 30, 2017.
- 450 [47] Yizhou Wang, Yue Kang, Can Qin, Huan Wang, Yilun Xu, Yulun Zhang, and Yun Raymond
451 Fu. Rethinking adam: A twofold exponential moving average approach. *arXiv preprint*
452 *arXiv:2106.11514*, 2021.
- 453 [48] Rachel Ward, Xiaoxia Wu, and Leon Bottou. Adagrad stepsizes: Sharp convergence over
454 nonconvex landscapes. In *International Conference on Machine Learning*, pages 6677–6686.
455 PMLR, 2019.
- 456 [49] Thomas Wolf. *Transfer Learning in Natural Language Processing*, 2019. Available at https://github.com/huggingface/naacl_transfer_learning_tutorial.
457
- 458 [50] Manzil Zaheer, Sashank Reddi, Devendra Sachan, Satyen Kale, and Sanjiv Kumar. Adaptive
459 methods for nonconvex optimization. *Advances in neural information processing systems*, 31,
460 2018.
- 461 [51] Bohang Zhang, Jikai Jin, Cong Fang, and Liwei Wang. Improved analysis of clipping algorithms
462 for non-convex optimization. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin,
463 editors, *Advances in Neural Information Processing Systems*, volume 33, pages 15511–15521.
464 Curran Associates, Inc., 2020.
- 465 [52] Jingzhao Zhang, Tianxing He, Suvrit Sra, and Ali Jadbabaie. Why gradient clipping accelerates
466 training: A theoretical justification for adaptivity. In *International Conference on Learning*
467 *Representations*, 2020.

- 468 [53] Jingzhao Zhang, Sai Praneeth Karimireddy, Andreas Veit, Seungyeon Kim, Sashank Reddi,
 469 Sanjiv Kumar, and Suvrit Sra. Why are adaptive methods good for attention models? *Advances*
 470 *in Neural Information Processing Systems*, 33:15383–15393, 2020.
- 471 [54] Difan Zou, Yuan Cao, Yuanzhi Li, and Quanquan Gu. Understanding the generalization of
 472 Adam in learning neural networks with proper regularization. *arXiv preprint arXiv:2108.11371*,
 473 2021.
- 474 [55] Fangyu Zou, Li Shen, Zequn Jie, Weizhong Zhang, and Wei Liu. A sufficient condition
 475 for convergences of Adam and RMSProp. In *Proceedings of the IEEE/CVF Conference on*
 476 *Computer Vision and Pattern Recognition*, pages 11127–11135, 2019.

477 Checklist

- 478 1. For all authors...
- 479 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
 480 contributions and scope? [Yes]
- 481 (b) Did you describe the limitations of your work? [Yes] See Section 6
- 482 (c) Did you discuss any potential negative societal impacts of your work? [N/A] This
 483 paper is mainly a theoretical one analyzing machine learning algorithms, we do not
 484 foresee any potential negative societal impacts of our work.
- 485 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
 486 them? [Yes]
- 487 2. If you are including theoretical results...
- 488 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See Section 3.
- 489 (b) Did you include complete proofs of all theoretical results? [Yes] See the Appendix.
- 490 3. If you ran experiments...
- 491 (a) Did you include the code, data, and instructions needed to reproduce the main ex-
 492 perimental results (either in the supplemental material or as a URL)? [Yes] In the
 493 supplemental material.
- 494 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
 495 were chosen)? [Yes] See Section 5.1 and A.4.
- 496 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
 497 ments multiple times)? [Yes]
- 498 (d) Did you include the total amount of compute and the type of resources used (e.g., type
 499 of GPUs, internal cluster, or cloud provider)? [Yes] See Section 5
- 500 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 501 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 502 (b) Did you mention the license of the assets? [Yes] Included in the code folder of the
 503 supplemental materials.
- 504 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 505 (d) Did you discuss whether and how consent was obtained from people whose data you’re
 506 using/curating? [No] Data use follows the permission of corresponding licenses.
- 507 (e) Did you discuss whether the data you are using/curating contains personally identifiable
 508 information or offensive content? [No]
- 509 5. If you used crowdsourcing or conducted research with human subjects...
- 510 (a) Did you include the full text of instructions given to participants and screenshots, if
 511 applicable? [N/A]
- 512 (b) Did you describe any potential participant risks, with links to Institutional Review
 513 Board (IRB) approvals, if applicable? [N/A]
- 514 (c) Did you include the estimated hourly wage paid to participants and the total amount
 515 spent on participant compensation? [N/A]