ON THE CONVERGENCE OF ADAM UNDER NON-UNIFORM SMOOTHNESS: SEPARABILITY FROM SGDM AND BEYOND

Anonymous authors

Paper under double-blind review

ABSTRACT

This paper aims to clearly distinguish between Stochastic Gradient Descent with Momentum (SGDM) and Adam in terms of their convergence rates. We demonstrate that Adam achieves a faster convergence compared to SGDM under the condition of non-uniformly bounded smoothness. Our findings reveal that: (1) in deterministic environments, Adam can attain the known lower bound for the convergence rate of deterministic first-order optimizers, whereas the convergence rate of Gradient Descent with Momentum (GDM) has higher order dependence on the initial function value; (2) in stochastic setting, Adam's convergence rate upper bound matches the lower bounds of stochastic first-order optimizers, considering both the initial function value and the final error, whereas there are instances where SGDM fails to converge with any learning rate. These insights distinctly differentiate Adam and SGDM regarding their convergence rates. Additionally, by introducing a novel stopping-time based technique, we further prove that if we consider the minimum gradient norm during iterations, the corresponding convergence rate can match the lower bounds across all problem hyperparameters. The technique can also help proving that Adam with a specific hyperparameter scheduler is parameter-agnostic, which hence can be of independent interest.

033

006

013

015

016

017

018

019

021

025

026

027

028

1 INTRODUCTION

034 Among various optimization techniques, the Adam optimizer Kingma & Ba (2014); Loshchilov & Hutter (2019) stands out due to its empirical success in a wide range of deep learning applications, 035 especially for pre-training large foundation models with enormous data Touvron et al. (2023); Brown 036 et al. (2020); Zhang et al. (2022a); Rae et al. (2021); Chowdhery et al. (2022); Du et al. (2021). 037 This popularity of Adam can be attributed to its adaptive learning rate mechanism, which smartly adjusts the step size for each parameter, allowing flexible and robust learning rate choices. Adam's versatility is further highlighted by its consistent performance in training various kinds of models, 040 making it a preferred optimizer in both academic and industrial settings Schneider et al. (2022). 041 Its empirical success extends beyond standard benchmarks to real-world challenges, where it often 042 delivers state-of-the-art results. This track record solidifies Adam's position as a fundamental tool for 043 deep learning practitioners.

044 Exploring the theoretical foundations of the Adam optimizer, particularly why it often outperforms traditional optimizers like Stochastic Gradient Descent with Momentum (SGDM), is an intriguing 046 yet complex task. Understanding Adam's convergence behavior is challenging, especially in settings 047 defined by standard convergence rate analysis. In these settings, assumptions include uniformly 048 bounded smoothness and finite gradient noise variance. Current research indicates that under these conditions, SGDM can attain the lower bound of the convergence rate for all first-order optimizers Carmon et al. (2017). This finding implies that, theoretically, Adam's convergence rate should not 051 exceed that of SGDM. This theoretical result contrasts with practical observations where Adam frequently excels, presenting a fascinating challenge for researchers. It highlights the need for 052 more refined theoretical models that can bridge the gap between Adam's empirical success and its theoretical understanding.

Recent research by Zhang et al. (2019) has provided valuable insights into the complexity of neural network optimization, particularly challenging the assumption of uniform bounded smoothness. Their observations indicate that smoothness often varies, showing a positive correlation with the norm of the gradient and experiencing considerable fluctuations during the optimization process. Building on this, they introduce the (L_0, L_1) -smooth condition (detailed in our Assumption 1), which posits that local smoothness can be bounded in relation to the gradient norm. This concept presents an exciting opportunity to theoretically demonstrate that Adam could potentially converge faster than SGDM. However, even in the relatively simpler deterministic settings, no study has yet conclusively shown this to be the case.

063 To effectively compare the convergence rates of Adam and Stochastic Gradient Descent with Mo-064 mentum (SGDM), it's essential to establish an upper bound on Adam's convergence rate and a lower bound for SGDM, and then prove Adam's superiority. This endeavor faces several challenges. First, 065 the known lower bound for SGDM's convergence rate is only available in deterministic settings 066 without momentum Zhang et al. (2019); Crawshaw et al. (2022). Moreover, this result is based on a 067 scenario where the counter-example objective function is selected after fixing the learning rate. This 068 procedure deviates from more common practices where the learning rate is adjusted after defining 069 the objective function Drori & Shamir (2020); Carmon et al. (2017); Arjevani et al. (2022), casting doubts on the standard applicability of this lower bound. Secondly, for Adam, the current assumptions 071 required to derive an upper bound for its convergence rate are quite strict. These include assumptions like bounded adaptive learning rates or deterministically bounded noise Wang et al. (2022); Li et al. 073 (2023a). However, even under these constraints, the convergence rates obtained for Adam are weaker 074 than those of algorithms like clipped SGDM Zhang et al. (2019).

These complexities hinder a straightforward comparison between the convergence rates of Adam and
 SGDM, highlighting a significant gap in the theoretical understanding that remains to be bridged.

Ora Our contributions. In this paper, we aim to bridge the gap and summarize our contributions as follows.

- We separate the convergence rate of Adam and SGDM under (L_0, L_1) -smooth condition both in the deterministic setting and in the stochastic setting.
 - In the deterministic setting, for the first time, we prove that under the (L_0, L_1) -smooth condition, the convergence rate of the Adam optimizer can match the existing lower bound for first-order deterministic optimizers, up to numerical constants. Additionally, we establish a new lower bound for the convergence rate of GDM, where one is allowed to tune the learning rate and the momentum coefficient after the problem is fixed. The lower bound exhibits a higher order dependence on the initial function value gap compared to the upper bound of Adam. This distinction clearly separates Adam and GDM for the deterministic setting.
- In the stochastic setting, for the first time, we prove that under the (L_0, L_1) -smooth condition, the convergence rate of Adam matches the existing lower bound for first-order stochastic optimizers regarding the initial function value $f(w_1) f^*$ and the final error ε . In contrast, counterexamples exist where SGDM fails to converge, irrespective of the learning rate and momentum coefficient. These findings distinctly separate the convergence properties of Adam and SGDM in stochastic settings.
- With the aid of a novel stopping time based technique, we further demonstrate that the convergence rate of minimum error point of Adam can match the lower bound across all problem hyperparameters. We demonstrate that such a technique can be of independent interest by proving that Adam with specific scheduler is parameter-agnostic based on the stopping time.
- 099 100 101

081 082

084

085

090

092

093

095

096

098

2 RELATED WORKS

102

103 **Convergence analysis under non-uniform smoothness.** Observations from empirical studies 104 on deep neural network training indicate that local smoothness can vary significantly throughout 105 the optimization process. In response to this, Zhang et al. (2019) introduced the (L_0, L_1) -smooth 106 condition, which posits that local smoothness can be bounded by a linear function of the gradient norm. 107 Subsequent works have extended this concept by generalizing the linear function to polynomials Chen 108 et al. (2023); Li et al. (2023a), or to more general functions Mei et al. (2021). Under non-uniform smoothness, convergence properties of various optimizers have been studied. For instance, upper bounds on the convergence rate have been established for optimizers such as Clipped SGDM Zhang et al. (2020), sign-based optimizers Jin et al. (2021); Hübler et al. (2023); Sun et al. (2023), AdaGrad Faw et al. (2023); Wang et al. (2023b), variance-reduction methods Reisizadeh et al. (2023); Chen et al. (2023), and trust-region methods Xie et al. (2023). However, research on lower bounds has been comparatively limited, with results primarily focusing on Gradient Descent.

Convergence analysis of Adam. The development of convergence analysis for Adam has been quite tortuous. While Adam was originally proposed with a convergence guarantee Kingma & Ba (2014), subsequent analysis by Reddi et al. (2018) pointed out flaws in this initial analysis and provided counterexamples claiming that Adam could fail to converge. Only recently, Shi et al. (2021) and Zhang et al. (2022b) have shown that the counterexamples in Reddi et al. (2018) only rule out the possibility that Adam can converge problem-agnostically, and it is still possible that Adam can converge with problem-dependent hyperparameters.

121 So far, several works have established the convergence of Adam under the L-smooth condition. 122 Zaheer et al. (2018) proved that Adam without momentum can converge to the neighborhood of 123 stationary points by additionally assuming that λ is large. De et al. (2018) showed that Adam 124 without momentum can converge to stationary points but under the strong assumption that the sign of gradients does not change during the optimization. Zou et al. (2019), Défossez et al. (2022), and Guo 125 et al. (2021) derived the convergence of Adam by assuming the stochastic gradient is bounded. Shi 126 et al. (2021) and Zhang et al. (2022b) characterized the convergence of random-reshuffling Adam but 127 suffer from sub-optimal rates. He et al. (2023) studied the non-ergodic convergence of Adam under 128 a bounded gradient assumption, while Hong & Lin (2023) provided high-probability guarantees 129 for Adam under a deterministically bounded noise assumption. A concurrent work by Wang et al. 130 (2023a) shows that Adam can achieve the lower bound of first-order optimizers with respect to the 131 final error ε under standard assumptions, but it is unknown whether Adam can match the lower bound 132 with respect to other problem specifics. 133

On the other hand, closely related to our work, there are only two works studying the convergence of Adam under non-uniform smoothness Wang et al. (2022); Li et al. (2023a), both with restricted assumptions and results. We will provide a detailed discussion in Section 4.

137 138

139

146 147

148

149

150

151

152

153 154 155

3 PRELIMINARY

Notations. In this paper, we will use asymptotic notations $\mathcal{O}, \Omega, \Theta$ to respectively denote asymptotically smaller, larger, and equivalent. We also use $\tilde{\mathcal{O}}, \tilde{\Omega}, \tilde{\Theta}$ to indicate that there is logarithmic factor hidden. We denote \mathcal{F}_t as the filter given by w_1, \cdots, w_t .

Problem and Algorithm. We study the unconstrained minimization problem $\min_{w} f(w)$. We present the psedo-code of Adam as follows.

Algoriumi I Adam Opumize		Algorithm	1	Adam	0	ptimizer
--------------------------	--	-----------	---	------	---	----------

Input: Stochastic oracle O, learning rate $\eta > 0$, initial point $w_1 \in \mathbb{R}^d$, initial conditioner $\nu_0 \in \mathbb{R}^+$, initial momentum m_0 , momentum parameter β_1 , conditioner parameter β_2 , number of epoch Tfor t = 1 to T do Generate a random z_t , and query stochastic oracle $g_t = O_f(w_t, z_t)$ Calculate $\nu_t = \beta_2 \nu_{t-1} + (1 - \beta_2) g_t^{\odot 2}$ Calculate $m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$ Update $w_{t+1} = w_t - \eta \frac{1}{\lambda + \sqrt{\nu_t}} \odot m_t$ end for

156 157 158

We would like to highlight that all the analysis in this paper is for $\lambda = 0$. This is because $\lambda = 0$ means we do not require the adaptive learning rate to be upper bounded (a restrictive assumption in existing works Li et al. (2023a); Guo et al. (2021)) and is most challenging. The proof can be immediately extended to $\lambda > 0$ without any modification. Meanwhile, we briefly state the SGDM optimizer as follows: with initial point w_1 and initial momentum m_0 , the update of t-th iteration of SGDM is given by

168

169

170 171

$$\boldsymbol{m}_t = \beta \boldsymbol{m}_{t-1} + (1-\beta)\boldsymbol{g}_t, \boldsymbol{w}_{t+1} = \boldsymbol{w}_t - \eta \boldsymbol{m}_t.$$

Assumptions. In this paper, all the analyses are established under the following two standard assumptions.

Assumption 1 ((L_0, L_1) -smooth condition). We assume f is differentiable and lower bounded, and there exist non-negative constants $L_0, L_1 > 0$, such that $\forall w_1, w_2 \in \mathbb{R}^d$ satisfying $||w_1 - w_2|| \leq \frac{1}{L_1}$,

$$\|\nabla f(\boldsymbol{w}_1) - \nabla f(\boldsymbol{w}_2)\| \le (L_0 + L_1 \|\nabla f(\boldsymbol{w}_1)\|) \|\boldsymbol{w}_1 - \boldsymbol{w}_2\|.$$

Assumption 2 (Affine noise variance). We assume that the stochastic noise g_t is unbiased, i.e., $\mathbb{E}^{|\mathcal{F}_t}g_t = G_t$. We further assume g_t has affine variance, i.e., there exists $\sigma_0 \ge 0, \sigma_1 \ge 1$, $\mathbb{E}^{|\mathcal{F}_t}[||g_t||^2] \le \sigma_0^2 + \sigma_1^2 ||\nabla f(w_t)||^2$.

Assumption 1 is a more general form of (L_0, L_1) -smooth condition and is equivalent to the Hessianbound form Zhang et al. (2019) when Hessian exists. Assumption 2 is one of the weakest assumptions on the noise in existing literature, and generalizes bounded variance assumption Li et al. (2023b), bounded gradient assumption Défossez et al. (2022), bounded noise assumption Li et al. (2023a).

180 181

182

186 187

188

4 SEPARATING THE CONVERGENCE RATES OF ADAM AND (S)GD

In this section, we elucidate the disparate convergence rates of Adam and (S)GD under Assumptions
 1 and 2, examining both deterministic and stochastic settings. We commence with the deterministic
 scenario before delving into the stochastic complexities.

4.1 ANALYSIS FOR THE DETERMINISTIC SETTING

As discussed in the introduction section, to discern the differential convergence rates of deterministic 189 Adam and GD, it is necessary to establish not only Adam's upper bound but also GD's lower bound, 190 given a consistent set of assumptions. Crucially, these bounds must be sufficiently tight to ensure 191 that Adam's upper bound is indeed the lesser. To date, only a couple of studies have addressed 192 the convergence of deterministic Adam. The first, referenced in Wang et al. (2022), indicates a 193 convergence rate of $\mathcal{O}(\frac{(f(\boldsymbol{w}_1)-f^*)^2}{\varepsilon^2})$, which is sub-optimal compared to the classical deterministic 194 rate of $\mathcal{O}(\frac{f(\boldsymbol{w}_1)-f^*}{\varepsilon^2})$ Zhang et al. (2019; 2020) regarding the initial function value gap $(f(\boldsymbol{w}_1)-f^*)$. 195 196 The second study, Li et al. (2023a), presents a convergence rate that depends polynomially on $\frac{1}{\lambda}$, 197 where λ is the small constant introduced to prevent the adaptive learning rate from becoming infinity. Therefore, their result is only non-vacuous when λ is large, which deviates from practical settings. Additionally, their bound exhibits an exaggerated dependency on the initial function value gap, 199 yielding $\min_{t \in [T]} \|\nabla f(\boldsymbol{w}_t)\| = \mathcal{O}(\frac{(f(\boldsymbol{w}_1) - f^*)^3}{\epsilon^2})$. As we will see later, such dependencies create 200 upper bounds that surpass the lower bounds of GD, making them unable to serve our purpose. To 201 overcome these limitations and accurately assess the performance of deterministic Adam, we propose 202 a new theorem that establishes an improved convergence rate for deterministic Adam. 203

²⁰⁴ An upper bound for the convergence rate of deterministic Adam.

Theorem 1 (Informal). Let Assumption 1 hold. Then, $\forall \beta_1, \beta_2 \ge 0$ satisfying $\beta_1^2 < \beta_2 < 1$, $\lambda = 0$, and $\varepsilon = \mathcal{O}(L_0/L_1)$, if $T \ge \Theta\left(\frac{L_0(f(\boldsymbol{w}_1) - f^*)}{\varepsilon^2}\right)$, then Algorithm 1 satisfies

$$\frac{1}{T}\sum_{t=1}^{T} \|\nabla f(\boldsymbol{w}_t)\| \leq \varepsilon.$$

213

207 208

Proof. Please see Appendix B.1 for the formal statement of theorem and the proof.

Our result offers a tighter bound than those presented in prior studies Wang et al. (2022); Li et al. (2023a). It is noteworthy that under the uniform smoothness constraint—where the objective function's smoothness is capped at L (that is, when $L_0 = L$ and $L_1 = 0$ as per Assumption 1, referred to 216 as the L-smooth condition in existing literature Arjevani et al. (2022); Carmon et al. (2017); Faw 217 et al. (2022))—Assumption 1 is met with $L_0 = L$ and any $L_1 \ge 0$. Consequently, the established 218 lower bound for all first-order optimizers Carmon et al. (2017) pertaining to the L-smooth condition 219 inherently provides a lower bound for the (L_0, L_1) -smooth condition, which is $\Omega\left(\frac{\sqrt{L_0(f(\mathbf{w}_1)-f^*)}}{\sqrt{T}}\right)$.

220

252 253

266

267

268

269

This coincides with our upper bound up to numerical constants. Such correspondence suggests that 221 222 our proposed bound is, in fact, optimal.

223 Our proof strategy utilizes a distinctive Lyapunov function, $f(w_t) + \frac{\beta_1}{2(1-\beta_1)\sqrt[4]{\beta_2}} \eta \frac{||\boldsymbol{m}_{t-1}||^2}{\lambda + \sqrt{\nu_{t-1}}}$, which 224 draws inspiration from the current analysis of Gradient Descent with Momentum (GDM) under the 225 L-smooth condition Sun et al. (2019). However, we have introduced significant modifications to 226 accommodate the integration of an adaptive learning rate. This carefully crafted Lyapunov function 227 enables us to effectively control the deviation between the momentum term and the current gradient, 228 even under (L_0, L_1) -smooth condition. Through this approach, we successfully establish the final 229 optimal bound. 230

Remark 1 (On the comparison with AdaGrad). Our result also suffices to separate Adam from 231 AdaGrad. It is important to note that the convergence rate of AdaGrad under the (L_0, L_1) -smooth 232 condition in a deterministic setting, as reported in Wang et al. (2023b), is $\frac{(f(w_1) - f^*)^2}{c^2}$. This rate 233 is outperformed by that of $Adam^1$. In Appendix B.3, we show that the rate in Wang et al. (2023b) 234 is tight by providing a counterexample. The comparatively slower convergence rate of AdaGrad 235 can be attributed to that (L_0, L_1) -smooth condition demands the update norm to be bounded by 236 $\mathcal{O}(1)$ to prevent the local smoothness from exponentially increasing. This, in turn, necessitates a 237 learning rate of $\mathcal{O}(1)$. However, the adaptive conditioner in AdaGrad, which accumulates over time, 238 causes the adaptive learning rate to become excessively small during later training stages, resulting 239 in reduced convergence speed. Conversely, Adam utilizes an exponential moving average for its 240 adaptive learning rate, which prevents the conditioner from accumulating excessively. Consequently, 241 Adam does not suffer from the aforementioned issue. 242

A lower bound for the convergence rate of GDM 243

244 With Adam's upper bound, we then move on to a lower bound for the convergence rate of GDM. In 245 fact, there has already been such lower bounds for GD in the existing literature Zhang et al. (2019); 246 Crawshaw et al. (2022), which we restate as follows: 247

Proposition 1 (Theorem 2, Crawshaw et al. (2022)). Fix ε , L_0 , L_1 , and Δ_1 , with learning rate η , 248 there exists objective function f satisfying (L_0, L_1) -smooth condition and $f(w_1) - f^* = \Delta_1$, such 249 that the minimum step T of GD to achieve final error ε (i.e., let $\{w_t\}_{t=1}^{\infty}$ be the iterates of GD, and 250 $T \triangleq \min\{t : \|\nabla f(\boldsymbol{w}_t)\| < \varepsilon\}$) satisfies 251

$$T = \tilde{\Omega} \left(\frac{L_1^2 \Delta_1^2 + L_0 \Delta_1}{\varepsilon^2} \right).$$

254 However, the proposition presents a limitation: the counter-example is chosen after the learning rate 255 has been determined. This approach is inconsistent with standard practices, where hyperparameters 256 are usually adjusted based on the specific task, and deviates from conventional lower bounds Carmon 257 et al. (2017); Arjevani et al. (2022) that offer assurances for optimally-tuned hyperparameters. This 258 type of result does not eliminate the possibility that, if the learning rate were adjusted after selecting the objective function—as is common practice—Gradient Descent (GD) could potentially achieve a 259 markedly faster convergence rate. This misalignment raises concerns about the appropriateness of 260 the proposition's methodology. Moreover, this proposition does not take momentum into account, a 261 technique that is commonly employed in conjunction with GD in practice. 262

263 To address these shortcomings, we introduce a new lower bound for GDM. This lower bound is 264 applicable under the standard practice of adjusting hyperparameters after the objective function has 265 been selected. Moreover, it encompasses scenarios where momentum is incorporated.

Theorem 2 (Informal). Fixing ε , L_0 , L_1 , and Δ_1 , there exists an objective function f satisfying (L_0, L_1) -smooth condition and $f(w_1) - f^* = \Delta_1$, such that for any learning rate $\eta > 0$ and

¹The state-of-art rate of AdaGrad under (L_0, L_1) -smooth condition and stochastic setting is $\frac{(f(w_1)-f^*)^2}{c^4}$, which is also worse than the rate of Adam established latter in Theorem 3.

 $\beta \in [0,1]$, the minimum step T of GDM to achieve final error ε satisfies

$$T = \tilde{\Omega} \left(\frac{L_1^2 \Delta_1^2 + L_0 \Delta_1}{\varepsilon^2} \right).$$

272 273

270

271

274 275 276

277

Proof. Please see Appendix B.2 for the formal statement of theorem and the proof.

It should be noticed in the above theorem, the hyperparameters (i.e., the learning rate and the momentum coefficient) are chosen after the objective function is determined, which agrees with practice and the settings of common lower bounds, and overcomes the shortcoming of Proposition 1. Moreover, as shown in Zhang et al. (2019), it is easy to prove that the upper bound of GD's convergence rate is also $\mathcal{O}\left(\frac{L_1^2\Delta_1^2+L_0\Delta_1}{\varepsilon^2}\right)$, which indicates such a lower bound is optimal.

The proof addresses two primary challenges outlined above. The first challenge involves handling momentum. To tackle this, we extend the counterexample provided in Proposition 1 for cases where the momentum coefficient β is small. Additionally, we introduce a new counterexample for situations with a large β , demonstrating how large momentum can bias the optimization process and decelerate convergence. The second challenge is how to derive a universal counterexample such that every hyperparameter setting will lead to slow convergence. We overcome this by a simple but effective trick: we independently put counterexamples for different hyperparameters in Proposition 1 over different coordinates and make it a whole counterexample. Therefore, for different hyperparameters, there will be at least one coordinate converge slowly, which leads to the final result.

292 Separating deterministic Adam and GDM. Upon careful examination of Theorem 1 and Theorem 2, 293 it becomes apparent that the convergence rate of GDM is inferior to that of Adam since $\frac{\sum_{t=1}^{T} ||G_t||}{T} \ge$ 294 min_{t∈[T]} $||G_t||$. Notably, GDM exhibits a more pronounced dependency on the initial function value 296 gap in comparison to Adam. This implies that with a sufficiently poor initial point, the convergence of 297 GDM can be significantly slower than that of Adam. The underlying reason for this disparity can be 298 attributed to GDM's inability to adeptly manage varying degrees of sharpness within the optimization 299 landscape. Consequently, GDM necessitates a learning rate selection that is conservative, tailored to 299 the most adverse sharpness encountered—often present during the initial optimization stages.

300 301 302

4.2 ANALYSIS FOR THE STOCHASTIC SETTING

303 Transitioning to the more complex stochastic setting, we extend our analysis beyond the deterministic 304 framework. As with our previous approach, we start by reviewing the literature to determine if 305 the existing convergence rates for Adam under the (L_0, L_1) -smooth condition can delineate a clear 306 distinction between the convergence behaviors of Adam and Stochastic Gradient Descent with 307 Momentum (SGDM). In fact, the only two studies that delve into this problem are the ones we discussed in Section 4.1, i.e., Wang et al. (2022); Li et al. (2023a). However, these results pertaining 308 to Adam are contingent upon rather stringent assumptions. Wang et al. (2022) postulates that 309 stochastic gradients not only conform to the (L_0, L_1) -smooth condition but are also limited to a 310 finite set of possibilities. These assumptions are more restrictive than merely assuming that the true 311 gradients satisfy the (L_0, L_1) -smooth condition, and such strong prerequisites are seldom employed 312 outside of the analysis of variance-reduction algorithms. Meanwhile, Li et al. (2023a) aligns its 313 findings on stochastic Adam with those on deterministic Adam, leading to a polynomial dependency 314 on $1/\lambda$, which deviates from practical scenarios as discussed in Section 4.1. Furthermore, it presumes 315 an a.s. bounded difference between stochastic gradients and true gradients, an assumption that closely 316 resembles the boundedness of stochastic gradients and is more limiting than the standard assumption 317 of bounded variance for stochastic gradients.

These more restricted and non-standard assumptions cast challenges in establishing a lower bound for the convergence of SGDM in the relevant contexts, let alone attempting a comparison between SGDM and Adam. In addition to the fact that these upper bounds fail to facilitate a clear comparison between Adam and SGDM, there are also concerns regarding their convergence rates. Wang et al. (2022) reports a convergence rate of $\frac{(f(w_1) - f^*)^2}{\varepsilon^8}$, which has a higher-order dependence on the initial function value gap and the final error than the $\frac{(f(w_1) - f^*)}{\varepsilon^4}$ rate established for Clipped SGDM under the (L_0, L_1) -smooth condition Zhang et al. $(2020)^2$. Furthermore, Li et al. (2023a) indicates a convergence rate of $\mathcal{O}(\frac{(f(w_1)-f^*)^4 \operatorname{poly}(1/\lambda)}{\varepsilon^4})$, which, aside from the previously mentioned dependency issues on $1/\lambda$, shows a significantly stronger dependence over the initial function value gap compared to the analysis of Clipped SGDM. This naturally leads to the question of whether such rates for Adam can be improved to match Clipped SGDM.

330 To tackle these obstacles, we present the following upper bound for Adam.

An upper bound for the convergence rate of Adam.

Theorem 3 (Informal). Let Assumptions 1 and 2 hold. Then, $\forall 1 > \beta_1 \ge 0$ and $\lambda = 0$, if $\varepsilon \le \frac{1}{\operatorname{poly}(f(\boldsymbol{w}_1) - f^*, L_0, L_1, \sigma_0, \sigma_1)}$, with a proper choice of learning rate η and momentum hyperparameter β_2 , we have if $T \ge \Theta\left(\frac{(L_0 + L_1\sigma_0)\sigma_0^2\sigma_1^2(f(\boldsymbol{w}_1) - f^*)}{\varepsilon^4}\right)$,

 $\frac{1}{T}\mathbb{E}\sum_{t=1}^{T} \|\nabla f(\boldsymbol{w}_t)\| \leq \varepsilon.$

Proof. Please see Appendix C.1 for the formal statement of theorem and the proof.

Below we include several discussions regarding Theorem 3. To begin with, one can immediately observe that Theorem 3 only requires Assumptions 1 and 2, and the convergence rate with respect to the initial function value gap and the final error $\frac{f(w_1)-f^*}{\varepsilon^4}$ matches that of Clipped SGDM Zhang et al. (2020) even with a weaker noise assumption. Therefore, our result successfully mitigate these barriers raised above. Indeed, to the best of our knowledge, it is for the first time that an algorithm is shown to converge with rate $\mathcal{O}\left(\frac{f(w_1)-f^*}{\varepsilon^4}\right)$ only requiring Assumptions 1 and 2, showcasing the advantage of Adam.

We briefly sketch the proof here before moving on to the result of SGDM. Specifically, the proof is inspired by recent analysis of Adam under *L*-smooth condition Wang et al. (2023a), but several challenges arise during the proof:

- The first challenge lies in the additional error introduced by the (L_0, L_1) -smooth condition. We address this by demonstrating that the telescoping sum involving the auxiliary function $\frac{\|G_t\|^2}{\sqrt{\nu_{t-1}}}$, as employed in Wang et al. (2023a), can bound this additional error when the adaptive learning rate is upper bounded. Although the adaptive learning rate in the Adam algorithm is not inherently bounded, we establish that the deviation incurred by employing a bounded surrogate adaptive learning rate is manageable;
- The second challenge involves deriving the desired dependence on the initial function value gap. Wang et al. (2023a) introduces two distinct proof strategies for bounding the conditioner ν_t and determining the final convergence rate. However, one strategy introduces an additional logarithmic dependence on ε , while the other exhibits sub-optimal dependence on the initial function value gap. We propose a novel two-stage divide-and-conquer approach to surmount this issue. In the first stage, we bound ν_t effectively. Subsequently, we leverage this bound within the original descent lemma to achieve the optimal dependence on $f(w_1) f^*$.

Remark 2 (On the limitations). Although Theorem 3 addresses certain deficiencies identified in prior 367 studies Wang et al. (2022); Li et al. (2023a), it is not without its limitations. As noted by Arjevani et al. 368 (2022), the established lower bound for the convergence rate of first-order optimization algorithms 369 under the L₀-smooth condition with bounded noise variance (specifically, $\sigma_0 = \sigma_0$ and $\sigma_1 = 1$ as 370 stated in Assumption 2) is $\mathcal{O}(\frac{(f(w_1)-f^*)L_0\sigma_0^2}{4})$. This sets a benchmark for the performance under 371 Assumptions 1 and 2. The upper bound of Adam's convergence rate as presented in Theorem 3 falls 372 short when compared to this benchmark, exhibiting a weaker noise scale dependency (σ_0^3 as opposed 373 to σ_0^2) and additional dependencies on L_1 and σ_1 . 374

To address these issues, we demonstrate in the subsequent section that by focusing on the convergence of the minimum gradient norm, $\mathbb{E}\min_{t\in[T]} \|\nabla f(\boldsymbol{w}_t)\|$, we can attain an improved convergence rate

377

333 334

335 336

341

353

354

355

357

359

360

361

362

364

²While Zhang et al. (2020) also assumes an a.s. bounded gap between stochastic gradients and true gradients.

of $\mathcal{O}(\frac{(f(\boldsymbol{w}_1)-f^*)L_0\sigma_0^2}{\varepsilon^4})$. This rate aligns with the aforementioned lower bound across all the problem hyperparameters.

We now establish the lower bound of SGDM. This is, however, more challenging than the deterministic case as to the best of our knowledge, there is no such a lower bound in existing literature (despite that the lower bounds of GD Zhang et al. (2019); Crawshaw et al. (2022) naturally offer a lower bound of SGD, which is considerably loose given the factor of $1/\varepsilon^2$). Intuitively, stochasticity can make the convergence of GDM even worse, as random fluctuations can inadvertently propel the iterations towards regions characterized by high smoothness even with a good initialization. We formulate this insight into the following theorem.

388 A lower bound for the convergence rate of SGDM.

Theorem 4 (Informal). Fix L_0, L_1 , and Δ_1 , there exists objective function f satisfying (L_0, L_1) smooth condition and $f(w_1) - f^* = \Delta_1$, and a gradient noise oracle satisfying Assumption 2, such that for any learning rate $\eta > 0$ and $\beta \in [0, 1]$, for all T > 0,

$$\min_{\mathbf{t}\in[T]} \mathbb{E} \|\nabla f(\boldsymbol{w}_t)\| = \|\nabla f(\boldsymbol{w}_1)\| \ge L_1 \Delta_1.$$

Proof. Please see Appendix C.2 for the formal statement of theorem and the proof.

Theorem 4 provides concrete evidence for the challenges inherent in the convergence of SGDM. It shows that there are instances that comply with Assumption 1 and Assumption 2 for which SGDM fails to converge, regardless of the chosen learning rate and momentum coefficient. This outcome confirms our earlier hypothesis: the stochastic elements within SGDM can indeed adversely affect its convergence properties under non-uniform smoothness.

Our proof is founded upon a pivotal observation: an objective function that escalates rapidly can effectively convert non-heavy-tailed noise into a "heavy-tailed" one. In particular, under the (L_0, L_1) smooth condition, the magnitude of the gradient is capable of exponential growth. As a result, even if the density diminishes exponentially, the expected value of the gradient norm may still become unbounded. This situation mirrors what occurs under the *L*-smooth condition when faced with heavy-tailed noise. Such a dynamic can lead to the non-convergence of SGDM.

407 Separating Adam and SGDM. Considering that Adam can achieve convergence under Assumptions 408 1 and 2, while SGD cannot, the superiority of Adam over SGDM becomes evident. It is important to 409 note, however, a recent study by Li et al. (2023b), which demonstrates that SGD can converge with 410 high probability under the same assumptions, provided the noise variance is bounded. We would like 411 to contextualize this finding in relation to our work as follows: First, this result does not conflict with 412 our Theorem 4, since our theorem pertains to bounds in expectation rather than with high probability. Second, our comparison of Adam and SGDM within an in-expectation framework is reasonable and 413 aligns with the convention of most existing lower bounds in the literature Carmon et al. (2017); Drori 414 & Shamir (2020); Arjevani et al. (2022). Moreover, establishing high-probability lower bounds is 415 technically challenging, and there are few references to such bounds in the existing literature. Lastly, 416 while we have not derived a corresponding high-probability lower bound for SGD, the upper bound 417 provided by Li et al. (2023b) is $\mathcal{O}(\frac{(f(\boldsymbol{w}_1)-f^*)^4}{c^4})$, which indicates a less favorable dependency on the 418 initial function value gap compared to the bound for Adam. 419

420

392 393 394

395

421 422

423

5 CAN ADAM REACH THE LOWER BOUND OF THE CONVERGENCE RATE UNDER (L_0, L_1) -SMOOTH CONDITION?

As we mentioned in Remark 2, although Theorem 3 matches the lower bound established by Arjevani et al. (2022) with respect to the initial function value gap $f(w_1) - f^*$, the final error ε , and the smoothness coefficient L_0 , it exhibits sub-optimal dependence on the noise scale σ_0 and additional dependence on L_1 and σ_1 . One may wonder whether these dependencies are inherently unavoidable or if they stem from technical limitations in our analysis.

429 Upon revisiting the proof, we identified that the sub-optimal dependencies arise from our strategy of 430 substituting the original adaptive learning rate with a bounded surrogate. For example, the correlation 431 between stochastic gradient and adaptive learning rate will introduce an error term $\eta \frac{\sigma_0^2(1-\beta_2)||g_t||^2}{\sqrt{\beta_2\nu_{t-1}\nu_t}}$, 432 detailed in Eq. (8). To bound this term, we add a constant λ to $\beta_2 \nu_{t-1}$, allowing us to upper bound $\frac{1}{\sqrt{\beta_2 \nu_{t-1} + \lambda}}$. Consequently, the term $\eta \frac{\sigma_0^2 (1 - \beta_2) \| \boldsymbol{g}_t \|^2}{\sqrt{\beta_2 \nu_{t-1} + \lambda} \nu_t}$ can be bounded by $\eta \frac{\sigma_0^2 (1 - \beta_2) \| \boldsymbol{g}_t \|^2}{\sqrt{\lambda} \nu_t}$, which has the same order as a second-order Taylor expansion. To control the error introduced by 433 434 435 adding λ , we cannot choose a value for λ that is too large. The optimal choice of λ for balancing 436 the new error against the original error is $(1 - \beta_2)\sigma_0^2$. This selection results in the original error 437 term $\eta \frac{\sigma_0 \sqrt{1-\beta_2} \|g_t\|^2}{\nu_t}$, which induces an additional σ_0 factor, ultimately leading to the sub-optimal 438 dependence on σ_0 . Therefore, we need to explore alternative methods to handle the error term to 439 440 eliminate the sub-optimal dependence on σ_0 .

We begin our analysis by observing that the term $\frac{(1-\beta_2)||g_t||^2}{\sqrt{\beta_2\nu_{t-1}\nu_t}}$ can in fact be bounded by an "approximate telescoping" series of $\frac{1}{\sqrt{\nu_t}}$ (noting an additional coefficient $\frac{1}{\sqrt{\beta_2}}$ in comparison to standard telescoping):

445 446 447

448

$$\frac{(1-\beta_2)\|\boldsymbol{g}_t\|^2}{\sqrt{\beta_2\boldsymbol{\nu}_{t-1}}\boldsymbol{\nu}_t} \leq \mathcal{O}\left(\frac{1}{\sqrt{\beta_2\boldsymbol{\nu}_{t-1}}} - \frac{1}{\sqrt{\boldsymbol{\nu}_t}}\right).$$

449 Accordingly, summing $\eta \frac{\sigma_0^2(1-\beta_2)||g_t||^2}{\sqrt{\beta_2\nu_{t-1}\nu_t}}$ over t yields a bound of $\mathcal{O}(\eta\sigma_0^2\sum_t(1-\beta_2)\frac{1}{\sqrt{\nu_t}})$. However, this term could potentially be unbounded since $\sqrt{\nu_t}$ is not lower bounded. To circumvent this issue, 450 451 we consider the first-order Taylor's expansion of the descent lemma, which, gives $-\sum_t \eta \frac{\|\nabla f(w_t)\|^2}{\sqrt{\nu_t}}$. 452 453 Intuitively, if any $\|\nabla f(\boldsymbol{w}_t)\|^2$ is of the order $\mathcal{O}(\sigma_0^2(1-\beta_2))$, our proof would be completed since we choose $1-\beta_2 = \Theta(\varepsilon^4)$. In the other case, the term $\mathcal{O}(\eta\sigma_0^2\sum_t(1-\beta_2)\frac{1}{\sqrt{\nu_t}})$ can be offset by the 454 455 negative term $-\sum_t \eta \frac{\|\nabla f(\boldsymbol{w}_t)\|^2}{\sqrt{\nu_t}}$. However, formalizing this intuition into a proof is challenging in 456 the context of stochastic analysis, where the randomness across iterations complicates the analysis. 457 Specifically, if we condition on the event that "no gradient norm is as small as $\sigma_0^2(1-\beta_2)$," which is 458 supported over the randomness of all iterations, it becomes difficult to express many expected values 459 (such as those from the first-order Taylor expansion) in closed form. 460

461 We address this difficulty by introducing a stopping time $\tau \triangleq \min\{t : \|\nabla f(\boldsymbol{w}_{t+1})\|^2 \leq \mathcal{O}(\sigma_0^2(1 - \beta_2))\}$. By applying the optimal stopping theorem Durrett (2019), we can maintain closed-form expressions for the expected values up to the stopping time, allowing the problematic error term to be absorbed within this interval. Building on this methodology, we formulate the following theorem.

Theorem 5 (Informal). Let Assumptions 1 and 2 hold. Then, $\forall 1 > \beta_1 \ge 0$, if $\varepsilon \le \frac{1}{\operatorname{Poly}(L_0, L_1, \sigma_0, \sigma_1, \frac{1}{1-\beta_1}, f(w_1) - f^*)}$, with a proper choice of learning rate η and momentum hyperparameter β_2 , we have that if $T \ge \Theta(\frac{L_0 \sigma_0^2(f(w_1) - f^*)}{\varepsilon^4})$

$$\mathbb{E}\min_{t\in[1,T]} \|\nabla f(\boldsymbol{w}_t)\| \leq \varepsilon.$$

470 471 472

473

Proof. Please see Appendix D.1 for the formal statement of theorem and the proof.

474 One can easily see that the convergence rate of Theorem 5 matches the lower bound in Arjevani et al. 475 (2022) with respect to all problem hyperparameters up to numerical constants even under the weaker 476 (L_0, L_1) -smooth condition. Therefore, such a rate is optimal and provides an affirmative answer to 477 the question raised in the beginning of this section.

478 One may notice that in the construction of the stopping time, we set the threshold for the squared gradient norm to be $\mathcal{O}(1-\beta_2)$. As we set $1-\beta_2 = \Theta(\varepsilon^4)$, the threshold is actually much smaller than what we aim for, since our goal is to have $\|\nabla f(\mathbf{w}_t)\|^2 \leq \varepsilon^2$. Therefore, based on the stopping-time 479 480 technique, we can actually show that Adam can converge with an optimal rate of $\mathcal{O}(\varepsilon^{-4})$ when 481 $1 - \beta_2 = \varepsilon^2$, or $1/\sqrt{T}$ if expressed in terms of the iteration number T. To the best of our knowledge, 482 483 this is the first time that Adam has been shown to converge with an optimal rate under the condition that $1 - \beta_2 = \Omega(1/T)$, which greatly enlarges the hyperparameter range. We show in Appendix 484 D.2 that based on this technique, we can show Adam is hyperparameter agnostic even under the 485 (L_0, L_1) -smooth condition.

486 CONCLUSION 6 487

488 In this paper, we have conducted a mathematical examination of the performance of the Adam 489 optimizer and SGDM within the context of non-uniform smoothness. Our convergence analysis 490 reveals that Adam exhibits a faster rate of convergence compared to SGDM under these conditions. 491 Moreover, we introduce a novel stopping time technique that demonstrates Adam's capability to 492 achieve the existing lower bounds for convergence rates. This finding underscores the robustness of Adam in complex optimization landscapes and contributes to a deeper understanding of its theoretical 493 494 properties.

REFERENCES

495 496

497

527

529

- Yossi Arjevani, Yair Carmon, John C Duchi, Dylan J Foster, Nathan Srebro, and Blake Woodworth. 498 Lower bounds for non-convex stochastic optimization. Mathematical Programming, pp. 1–50, 499 2022. 500
- 501 Larry Armijo. Minimization of functions having lipschitz continuous first partial derivatives. Pacific 502 Journal of mathematics, 16(1):1–3, 1966.
- 504 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, 505 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, 506 Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott 507 Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya 508 Sutskever, and Dario Amodei. Language models are few-shot learners, 2020. 509
- 510 Yair Carmon, John C Duchi, Oliver Hinder, and Aaron Sidford. Lower bounds for finding stationary 511 points i. arXiv preprint arXiv:1710.11606, 2017. 512
- 513 Ziyi Chen, Yi Zhou, Yingbin Liang, and Zhaosong Lu. Generalized-smooth nonconvex optimization is as efficient as smooth nonconvex optimization. arXiv preprint arXiv:2303.02854, 2023. 514
- 515 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam 516 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, 517 Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam 518 Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James 519 Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin 521 Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. 522 Dai, Thanumalayan Sankaranarayana Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon 523 Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark 524 Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, 525 Slav Petrov, and Noah Fiedel. Palm: Scaling language modeling with pathways, 2022. 526
- Michael Crawshaw, Mingrui Liu, Francesco Orabona, Wei Zhang, and Zhenxun Zhuang. Robustness 528 to unbounded smoothness of generalized signSGD. arXiv preprint arXiv:2208.11195, 2022.
- 530 Ashok Cutkosky and Harsh Mehta. Momentum improves normalized SGD. In International conference on machine learning, pp. 2260–2268. PMLR, 2020. 531
- 532 Soham De, Anirbit Mukherjee, and Enayat Ullah. Convergence guarantees for RMSProp and ADAM 533 in non-convex optimization and an empirical comparison to Nesterov acceleration. arXiv preprint 534 arXiv:1807.06766, 2018. 535
- Alexandre Défossez, Leon Bottou, Francis Bach, and Nicolas Usunier. A simple convergence proof 537 of Adam and Adagrad. Transactions on Machine Learning Research, 2022.
- Yoel Drori and Ohad Shamir. The complexity of finding stationary points with stochastic gradient descent. In International Conference on Machine Learning, pp. 2658–2667. PMLR, 2020.

540 541 542	Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. All NLP tasks are generation tasks: A general pretraining framework. <i>CoRR</i> , abs/2103.10360, 2021. URL https://arxiv.org/abs/2103.10360.
544	Rick Durrett. Probability: theory and examples, volume 49. Cambridge university press, 2019.
545 546 547 548	Matthew Faw, Isidoros Tziotis, Constantine Caramanis, Aryan Mokhtari, Sanjay Shakkottai, and Rachel Ward. The power of adaptivity in SGD: Self-tuning step sizes with unbounded gradients and affine variance. In <i>Conference on Learning Theory</i> , pp. 313–355. PMLR, 2022.
549 550	Matthew Faw, Litu Rout, Constantine Caramanis, and Sanjay Shakkottai. Beyond uniform smoothness: A stopped analysis of adaptive sgd. <i>arXiv preprint arXiv:2302.06570</i> , 2023.
551 552 553	Zhishuai Guo, Yi Xu, Wotao Yin, Rong Jin, and Tianbao Yang. A novel convergence analysis for algorithms of the Adam family. <i>arXiv preprint arXiv:2112.03459</i> , 2021.
554 555 556 557	Meixuan He, Yuqing Liang, Jinlan Liu, and Dongpo Xu. Convergence of adam for non-convex objectives: Relaxed hyperparameters and non-ergodic case. <i>arXiv preprint arXiv:2307.11782</i> , 2023.
558 559	Yusu Hong and Junhong Lin. High probability convergence of adam under unbounded gradients and affine variance noise. <i>arXiv preprint arXiv:2311.02000</i> , 2023.
560 561 562	Florian Hübler, Junchi Yang, Xiang Li, and Niao He. Parameter-agnostic optimization under relaxed smoothness. <i>arXiv preprint arXiv:2311.03252</i> , 2023.
563 564 565	Jikai Jin, Bohang Zhang, Haiyang Wang, and Liwei Wang. Non-convex distributionally robust optimization: Non-asymptotic analysis. <i>Advances in Neural Information Processing Systems</i> , 34: 2771–2782, 2021.
566 567	Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2014.
568 569 570	Haochuan Li, Ali Jadbabaie, and Alexander Rakhlin. Convergence of Adam under relaxed assumptions. <i>arXiv preprint arXiv:2304.13972</i> , 2023a.
571 572	Haochuan Li, Jian Qian, Yi Tian, Alexander Rakhlin, and Ali Jadbabaie. Convex and non-convex optimization under generalized smoothness. <i>arXiv preprint arXiv:2306.01264</i> , 2023b.
573 574 575	Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In International Confer- ence on Learning Representations, 2019.
576 577 578	Jincheng Mei, Yue Gao, Bo Dai, Csaba Szepesvari, and Dale Schuurmans. Leveraging non-uniformity in first-order non-convex optimization. In <i>International Conference on Machine Learning</i> , pp. 7555–7564. PMLR, 2021.
579 580	Yurii Nesterov et al. Lectures on convex optimization, volume 137. Springer, 2018.
581 582 583	Francesco Orabona and Dávid Pál. Coin betting and parameter-free online learning. Advances in Neural Information Processing Systems, 29, 2016.
584 585 586	Francesco Orabona and Tatiana Tommasi. Training deep networks without learning rates through coin betting. <i>Advances in Neural Information Processing Systems</i> , 30, 2017.
587 588 589	Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, F. Song, John Aslanides, Sarah Henderson, R. Ring, S. Young, et al. Scaling language models: Methods, analysis & insights from training gopher. arXiv preprint arXiv:2112.11446, 2021.
590 591	Sashank J. Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. In <i>International Conference on Learning Representations</i> , 2018.
592 593	Amirhossein Reisizadeh, Haochuan Li, Subhro Das, and Ali Jadbabaie. Variance-reduced clipping for non-convex optimization. <i>arXiv preprint arXiv:2303.00883</i> , 2023.

594 595 596	Frank Schneider, Zachary Nado, Naman Agarwal, George E. Dahl, and Philipp Hennig. HITY work- shop poll, NeurIPS 2022. https://github.com/fsschneider/HITYWorkshopPoll, 2022.
597	
598 599	Naichen Shi, Dawei Li, Mingyi Hong, and Ruoyu Sun. RMSprop converges with proper hyper- parameter. In <i>International Conference on Learning Representations</i> , 2021.
600	
601	Tao Sun, Penghang Yin, Dongsheng Li, Chun Huang, Lei Guan, and Hao Jiang. Non-ergodic
602	Artificial Intelligence, volume 33, pp. 5033–5040, 2019.
603	Tao Sun, Congliang Chen, Peng Oiao, Li Shen, Xinwang Liu, and Dongsheng Li, Rethinking sign
604 605	training: Provable nonconvex acceleration without first-and second-order gradient lipschitz. arXiv preprint arXiv:2310.14616, 2023
606	
607 608 609	Hugo Touvron, Thibault Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothee Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> , 2023.
610	
611 612	Provable adaptivity in Adam. <i>arXiv preprint arXiv:2208.09900</i> , 2022.
613	Bohan Wang Jingwen Fu Huishuai Zhang Nanning Zheng and Wei Chen. Closing the gap between
614 615	the upper bound and lower bound of adam's iteration complexity. In <i>Thirty-seventh Conference on</i> <i>Neural Information Processing Systems</i> , 2023a.
616	
617	Bohan Wang, Huishuai Zhang, Zhiming Ma, and Wei Chen. Convergence of adagrad for non-convex
618	objectives: Simple proofs and relaxed assumptions. In <i>The Thirty Sixth Annual Conference on Learning Theory</i> , pp. 161–190. PMLR, 2023b.
600	Rachel Ward Xiaoxia Wu and Leon Bottou. Adagrad stensizes: Sharp convergence over nonconvex
621	landscapes. The Journal of Machine Learning Research, 21(1):9047–9076, 2020.
622	Chenghan Xie, Chenxi Li, Chuwen Zhang, Qi Deng, Dongdong Ge, and Yinyu Ye. Trust region
623 624	methods for nonconvex stochastic optimization beyond lipschitz smoothness. <i>arXiv preprint arXiv:2310.17319</i> , 2023.
625 626 627	Yu Xing, Xingkang He, et al. On the convergence of msgd and adagrad for stochastic optimization. In International Conference on Learning Representations, 2021.
628 629	Junchi Yang, Xiang Li, Ilyas Fatkhullin, and Niao He. Two sides of one coin: the limits of untuned sgd and the power of adaptive methods. <i>arXiv preprint arXiv:2305.12475</i> , 2023.
630	Manzil Zahaar Sashank Daddi Davandra Sashan Satuan Vala and Sartin Varman Adartin weeks the
631 632	for nonconvex optimization. Advances in neural information processing systems, 31, 2018.
633	Bohang Zhang, Jikai Jin, Cong Fang, and Liwei Wang. Improved analysis of clipping algorithms for
634 635	non-convex optimization. Advances in Neural Information Processing Systems, 33:15511–15521, 2020.
636	
637	Jingznao Zhang, Hanxing He, Suvrit Sra, and Ali Jadoabale. Why gradient clipping accelerates
638 639	Representations, 2019.
640	Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Ming-Wei Chen, Shuohui Chen. Christo-
641 642	pher Dewan, Mona Diab, Xiaodong Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language models. <i>arXiv preprint arXiv:2205.01068</i> , 2022a.
643	Vishing Thomas Congligner Chan Meigher Shi David State 1711 One I. A david
644 645	rusnun Znang, Congliang Chen, Naichen Shi, Ruoyu Sun, and Zhi-Quan Luo. Adam can converge without any modification on update rules. <i>arXiv preprint arXiv:2208.09632</i> , 2022b.
646 647	Fangyu Zou, Li Shen, Zequn Jie, Weizhong Zhang, and Wei Liu. A sufficient condition for convergences of Adam and RMSProp. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pp. 11127–11135, 2019.

A AUXILIARY LEMMAS

In this section, we provide auxiliary results which will be used in subsequent results.

Lemma 1. We have
$$\forall t \geq 1$$
, $\|\boldsymbol{w}_{t+1} - \boldsymbol{w}_t\| \leq \eta \frac{1-\beta_1}{\sqrt{1-\beta_2}\sqrt{1-\frac{\beta_1^2}{\beta_2}}}$.

Proof. We have that

$$\begin{split} \|\boldsymbol{w}_{t+1} - \boldsymbol{w}_t\| &= \eta \left| \frac{\boldsymbol{m}_t}{\sqrt{\boldsymbol{\nu}_t}} \right| \le \eta \frac{\sum_{i=0}^{t-1} (1-\beta_1) \beta_1^i \|\boldsymbol{g}_{t-i}\|}{\sqrt{\sum_{i=0}^{t-1} (1-\beta_2) \beta_2^i} \|\boldsymbol{g}_{t-i}\|^2 + \beta_2^t \boldsymbol{\nu}_0} \\ \le \eta \frac{1-\beta_1}{\sqrt{1-\beta_2}} \frac{\sqrt{\sum_{i=0}^{t-1} \beta_2^i} \|\boldsymbol{g}_{t-i}\|^2}{\sqrt{\sum_{i=0}^{t-1} \beta_2^i} \|\boldsymbol{g}_{t-i}\|^2} \le \eta \frac{1-\beta_1}{\sqrt{1-\beta_2} \sqrt{1-\beta_1^2}}. \end{split}$$

Here the second inequality is due to Cauchy's inequality. The proof is completed.

The following lemma provides a novel descent lemma under (L_0, L_1) -smooth condition. **Lemma 2.** Let Assumption 1 hold. Then, for any three points $w^1, w^2, w^3 \in \mathcal{X}$ satisfying $||w^1 - w^2|| \leq \frac{1}{2L_1}$ and $||w^1 - w^3|| \leq \frac{1}{2L_1}$, we have

$$f(\boldsymbol{w}^2) \le f(\boldsymbol{w}^3) + \langle \nabla f(\boldsymbol{w}^1), \boldsymbol{w}^2 - \boldsymbol{w}^3 \rangle + \frac{1}{2} (L_0 + L_1 \| \nabla f(\boldsymbol{w}^1) \|) \| \boldsymbol{w}^2 - \boldsymbol{w}^3 \| (\| \boldsymbol{w}^1 - \boldsymbol{w}^3 \| + \| \boldsymbol{w}^1 - \boldsymbol{w}^2 \|).$$

Proof. By the Fundamental Theorem of Calculus, we have

$$\begin{split} f(\boldsymbol{w}^{2}) =& f(\boldsymbol{w}^{3}) + \int_{0}^{1} \langle \nabla f(\boldsymbol{w}^{3} + a(\boldsymbol{w}^{2} - \boldsymbol{w}^{3})), \boldsymbol{w}^{2} - \boldsymbol{w}^{3} \rangle \mathrm{d}a \\ =& f(\boldsymbol{w}^{3}) + \langle \nabla f(\boldsymbol{w}^{1}), \boldsymbol{w}^{2} - \boldsymbol{w}^{3} \rangle + \int_{0}^{1} \langle \nabla f(\boldsymbol{w}^{3} + a(\boldsymbol{w}^{2} - \boldsymbol{w}^{3})) - \nabla f(\boldsymbol{w}^{1}), \boldsymbol{w}^{2} - \boldsymbol{w}^{3} \rangle \mathrm{d}a \\ \leq& f(\boldsymbol{w}^{3}) + \langle \nabla f(\boldsymbol{w}^{1}), \boldsymbol{w}^{2} - \boldsymbol{w}^{3} \rangle + \int_{0}^{1} \|\nabla f(\boldsymbol{w}^{3} + a(\boldsymbol{w}^{2} - \boldsymbol{w}^{3})) - \nabla f(\boldsymbol{w}^{1})\| \|\boldsymbol{w}^{2} - \boldsymbol{w}^{3}\| \mathrm{d}a \\ \leq& f(\boldsymbol{w}^{3}) + \langle \nabla f(\boldsymbol{w}^{1}), \boldsymbol{w}^{2} - \boldsymbol{w}^{3} \rangle + \int_{0}^{1} (L_{0} + L_{1} \|\nabla f(\boldsymbol{w}^{1})\|) \|a(\boldsymbol{w}^{2} - \boldsymbol{w}^{1}) + (1 - a)(\boldsymbol{w}^{3} - \boldsymbol{w}^{1})\| \|\boldsymbol{w}^{2} - \boldsymbol{w}^{3}\| \mathrm{d}a \\ \leq& f(\boldsymbol{w}^{3}) + \langle \nabla f(\boldsymbol{w}^{1}), \boldsymbol{w}^{2} - \boldsymbol{w}^{3} \rangle + \frac{1}{2} (L_{0} + L_{1} \|\nabla f(\boldsymbol{w}^{1})\|) \|\boldsymbol{w}^{2} - \boldsymbol{w}^{3}\| (\|\boldsymbol{w}^{1} - \boldsymbol{w}^{3}\| + \|\boldsymbol{w}^{1} - \boldsymbol{w}^{2}\|), \end{split}$$

where Inequality (\star) is because due to

$$\|\boldsymbol{w}^3 + a(\boldsymbol{w}^2 - \boldsymbol{w}^3) - \boldsymbol{w}^1\| = \|a(\boldsymbol{w}^2 - \boldsymbol{w}^1) + (1 - a)(\boldsymbol{w}^3 - \boldsymbol{w}^1)\| \le \frac{1}{L_1},$$

the definition of (L_0, L_1) -smooth condition can be applied.

690691 The proof is completed.

The following lemma is helpful when bounding the second-order term.

Lemma 3. Assume we have $0 < \beta_1^2 < \beta_2 < 1$ and a sequence of real numbers $(a_n)_{n=1}^{\infty}$. Let $b_0 > 0$, $b_n = \beta_2 b_{n-1} + (1 - \beta_2) a_n^2$, $c_0 = 0$, and $c_n = \beta_1 c_{n-1} + (1 - \beta_1) a_n$. Then, we have

$$\sum_{n=1}^{T} \frac{|c_n|^2}{b_n} \le \frac{(1-\beta_1)^2}{(1-\frac{\beta_1}{\sqrt{\beta_2}})^2(1-\beta_2)} \left(\ln\left(\frac{b_T}{b_0}\right) - T\ln\beta_2 \right).$$

Proof. This is a lemma commonly adopted in the literature of the convergence of Adam Défossez et al. (2022); Wang et al. (2023a). We invite interesting readers to see (Lemma A.2, Défossez et al. (2022)) for the proof.

Lemma 4. If $\beta_2 \geq \beta_1$, then we have

$$\frac{\|\boldsymbol{m}_t\|^2}{(\sqrt{\boldsymbol{\nu}_t})^3} \le 4(1-\beta_1) \left(\sum_{s=1}^t \sqrt[4]{\beta_1^{t-s}} \frac{2}{1-\beta_2} \left(\frac{1}{\sqrt{\beta_2 \boldsymbol{\nu}_{s-1}}} - \frac{1}{\sqrt{\boldsymbol{\nu}_s}} \right) \right).$$

Proof. To begin with, we have

$$\frac{\|\boldsymbol{m}_t\|}{\sqrt[4]{\boldsymbol{\nu}_t^3}} \le (1-\beta_1) \sum_{s=1}^t \frac{\beta_2^{t-s} \|\boldsymbol{g}_s\|}{\sqrt[4]{\boldsymbol{\nu}_t^3}} \le (1-\beta_1) \sum_{s=1}^t \frac{\beta_1^{t-s} \|\boldsymbol{g}_s\|}{\sqrt[4]{\beta_2^{3(t-s)}}} \frac{1}{\sqrt[4]{\boldsymbol{\nu}_s^3}}.$$

Here in the last inequality we use $\nu_t \ge \beta_2^{t-s} \nu_s$.

714 By further applying Cauchy-Schwartz inequality, we obtain

$$\frac{\|\boldsymbol{m}_t\|^2}{\sqrt{\boldsymbol{\nu}_t^3}} \le (1-\beta_1)^2 \left(\sum_{s=1}^t \frac{\beta_1^{t-s} \|\boldsymbol{g}_s\|^2}{\sqrt[4]{\beta_2^{3(t-s)}}} \sqrt{\boldsymbol{\nu}_s^3}\right) \left(\sum_{s=1}^t \frac{\beta_1^{t-s}}{\sqrt[4]{\beta_2^{3(t-s)}}}\right)$$

$$\leq \frac{(1-\beta_1)^2}{1-\frac{\beta_1}{\sqrt[4]{\beta_2^3}}} \left(\sum_{s=1}^t \frac{\beta_1^{t-s} \|\boldsymbol{g}_s\|^2}{\sqrt[4]{\beta_2^{3(t-s)}}} \sqrt{\boldsymbol{\nu}_s^3} \right)$$

$$\leq 4(1-\beta_1) \left(\sum_{s=1}^t \frac{\beta_1^{t-s} \|\boldsymbol{g}_s\|^2}{\sqrt[4]{\beta_2^{3(t-s)}}} \sqrt{\boldsymbol{\nu}_s^3} \right).$$

As
$$\frac{\|\boldsymbol{g}_s\|^2}{\sqrt{\boldsymbol{\nu}_s^3}} \leq \frac{2\|\boldsymbol{g}_s\|^2}{\sqrt{\boldsymbol{\nu}_s}\sqrt{\beta_2\boldsymbol{\nu}_{s-1}}(\sqrt{\boldsymbol{\nu}_s}+\sqrt{\beta_2\boldsymbol{\nu}_{s-1}})} = \frac{2}{1-\beta_2}\left(\frac{1}{\sqrt{\beta_2\boldsymbol{\nu}_{s-1}}}-\frac{1}{\sqrt{\boldsymbol{\nu}_s}}\right)$$
, the proof is completed. \Box

Lemma 5. Under the same set of assumptions in Theorem 11, if $\beta_2 \ge \beta_1$, then we have

$$\frac{\|\boldsymbol{m}_t\|^2 \|\boldsymbol{G}_t\|^2}{\boldsymbol{\nu}_t \sqrt{\beta_2 \boldsymbol{\nu}_{t-1}}} \le 4(1-\beta_1) \left(\sum_{s=1}^t \frac{\sqrt[8]{\beta_1^{t-s}} \|\boldsymbol{g}_s\|^2 \|\boldsymbol{G}_s\|^2}{\boldsymbol{\nu}_s \sqrt{\beta_2 \boldsymbol{\nu}_{s-1}}} \right) + 8\frac{1-\beta_1}{1-\beta_2} \frac{L_1^2}{L_0^2} \left(\sum_{s=1}^t \sqrt[8]{\beta_1^{t-s}} \left(\frac{1}{\sqrt{\beta_2 \boldsymbol{\nu}_{s-1}}} - \frac{1}{\sqrt{\boldsymbol{\nu}_s}} \right) \right).$$

Proof. Similar to the proof of Lemma 4, we have

$$\frac{\|\boldsymbol{m}_t\|^2}{\sqrt{\beta_2 \boldsymbol{\nu}_{t-1} \boldsymbol{\nu}_t}} \le 4(1-\beta_1) \left(\sum_{s=1}^t \frac{\beta_1^{t-s} \|\boldsymbol{g}_s\|^2}{\sqrt[4]{\beta_2^{3(t-s)}} \sqrt{\beta_2 \boldsymbol{\nu}_{s-1} \boldsymbol{\nu}_s}} \right).$$
(1)

Meanwhile, according to Assumption 1, we have

$$\begin{aligned} \|\boldsymbol{G}_{t}\|^{2} &\leq \|\boldsymbol{G}_{t-1}\|^{2} + 2\|\boldsymbol{G}_{t-1}\|\|\boldsymbol{G}_{t} - \boldsymbol{G}_{t-1}\| + \|\boldsymbol{G}_{t} - \boldsymbol{G}_{t-1}\|^{2} \\ &\leq \|\boldsymbol{G}_{t-1}\|^{2} + 2\|\boldsymbol{G}_{t-1}\|(\boldsymbol{L}_{0} + \boldsymbol{L}_{1}\|\boldsymbol{G}_{t-1}\|)\|\boldsymbol{w}_{t+1} - \boldsymbol{w}_{t}\| + 2(\boldsymbol{L}_{0}^{2} + \boldsymbol{L}_{1}^{2}\|\boldsymbol{G}_{t-1}\|^{2})\|\boldsymbol{w}_{t+1} - \boldsymbol{w}_{t}\|^{2} \\ &\leq \|\boldsymbol{G}_{t-1}\|^{2} + \frac{1 - \sqrt[8]{\beta_{1}}}{2\sqrt[8]{\beta_{1}}}\|\boldsymbol{G}_{t-1}\|^{2} + \frac{3\sqrt[8]{\beta_{1}}\boldsymbol{L}_{0}^{2}}{\sqrt[8]{\beta_{1}}}\|\boldsymbol{w}_{t+1} - \boldsymbol{w}_{t}\|^{2} + 2\boldsymbol{L}_{1}\|\boldsymbol{G}_{t-1}\|^{2}\|\boldsymbol{w}_{t+1} - \boldsymbol{w}_{t}\| \end{aligned}$$

$$\leq \|\boldsymbol{G}_{t-1}\|^{2} + \frac{1 - \sqrt{\rho_{1}}}{3\sqrt[8]{\beta_{1}}} \|\boldsymbol{G}_{t-1}\|^{2} + \frac{3\sqrt{\rho_{1}L_{0}}}{1 - \sqrt[8]{\beta_{1}}} \|\boldsymbol{w}_{t+1} - \boldsymbol{w}_{t}\|^{2} + 2L_{1} \|\boldsymbol{G}_{t-1}\|^{2} \|\boldsymbol{w}_{t+1} - 2L_{0}\|\boldsymbol{G}_{t-1}\|^{2} \|\boldsymbol{w}_{t+1} - \boldsymbol{w}_{t}\|^{2} + 2L_{1} \|\boldsymbol{G}_{t-1}\|^{2} \|\boldsymbol{w}_{t+1} - \boldsymbol{w}_{t}\|^{2}$$

$$\stackrel{(\star)}{\leq} \|\boldsymbol{G}_{t-1}\|^{2} + \frac{1 - \sqrt[8]{\beta_{1}}}{3\sqrt[8]{\beta_{1}}} \|\boldsymbol{G}_{t-1}\|^{2} + \frac{1 - \sqrt[8]{\beta_{1}}}{2} \frac{L_{0}^{2}}{L_{1}^{2}} + \frac{1 - \sqrt[8]{\beta_{1}}}{3\sqrt[8]{\beta_{1}}} \|\boldsymbol{G}_{t-1}\|^{2}$$

$$+ \frac{1 - \sqrt[8]{\beta_{1}}}{2} \frac{L_{0}^{2}}{L_{1}^{2}} + \frac{1 - \sqrt[8]{\beta_{1}}}{3\sqrt[8]{\beta_{1}}} \|\boldsymbol{G}_{t-1}\|^{2}$$

$$\frac{1}{2} \|\boldsymbol{G}_{t-1}\|^{2} - \frac{1 - \sqrt[8]{\beta_{1}}}{3\sqrt[8]{\beta_{1}}} \|\boldsymbol{G}_{t-1}\|^{2}$$

 $\leq \frac{1}{\sqrt[8]{\beta_1}} \|\boldsymbol{G}_{t-1}\|^2 + (1 - \sqrt[8]{\beta_1}) \frac{L_1^{-1}}{L_0^2}.$

Here inequality (*) is because $||w_{t+1} - w_t|| \le \frac{1 - \frac{8}{6L_1}}{6L_1}$ (According to Lemma 1 and the choice of η and β_2 in Theorem 11, we have $|w_{t+1} - w_t| \le \frac{(1 - \beta_1)\sqrt{1 - \beta_1}}{256\sigma_1^2 L_1}$, and to prove the conclusion, we need to show that $\frac{1-\sqrt[8]{\beta_1}}{6} \ge \frac{(1-\beta_1)\sqrt{1-\beta_1}}{256\sigma_1^2}$. Since $(1-\sqrt[8]{\beta_1})(1+\sqrt[8]{\beta_1})(1+\sqrt[4]{\beta_1})(1+\sqrt[2]{\beta_1}) = (1-\beta_1)$, and $(1+\sqrt[8]{\beta_1})(1+\sqrt[4]{\beta_1})(1+\sqrt[4]{\beta_1})(1+\sqrt[4]{\beta_1}) \le 2 \times 2 \times 2 \le 8$, it follows that $\frac{1-\sqrt[8]{\beta_1}}{6} \ge \frac{1-\beta_1}{48} \ge \frac{(1-\beta_1)\sqrt{1-\beta_1}}{256\sigma_1^2}$. Thus, the claim is proven). Recursively applying the above inequality, we obtain that $\|\boldsymbol{G}_t\|^2 \leq rac{1}{\sqrt[8]{eta_1^{t-s}}} \|\boldsymbol{G}_s\|^2 + \left(\left(rac{1}{\sqrt[8]{eta_1}}
ight)^{t-s} - 1
ight) rac{L_1^2}{L_0^2},$ which by Eq. (1) further gives $\frac{\|\boldsymbol{m}_t\|^2 \|\boldsymbol{G}_t\|^2}{\sqrt{\beta_2 \boldsymbol{\nu}_{t-1} \boldsymbol{\nu}_t}} \leq 4(1-\beta_1) \left(\sum_{s=1}^t \frac{\beta_1^{t-s} \|\boldsymbol{g}_s\|^2 \|\boldsymbol{G}_t\|^2}{\sqrt[4]{\beta_1^{3(t-s)} \boldsymbol{\mu}_s} \sqrt{\beta_2 \boldsymbol{\mu}_{s-1}}} \right)$ $\leq 4(1-\beta_1) \left(\sum_{s=1}^t \frac{\sqrt[8]{\beta_1^{t-s}} \|\boldsymbol{g}_s\|^2 \|\boldsymbol{G}_s\|^2}{\boldsymbol{\nu}_s \sqrt{\beta_2 \boldsymbol{\nu}_{s-1}}} + \sum_{s=1}^t \frac{\sqrt[8]{\beta_1^{t-s}} \|\boldsymbol{g}_s\|^2}{\boldsymbol{\nu}_s \sqrt{\beta_2 \boldsymbol{\nu}_{s-1}}} \frac{L_1^2}{L_0^2} \right)$ $\leq 4(1-\beta_1) \left(\sum_{s=1}^t \frac{\sqrt[8]{\beta_1^{t-s}} \|\boldsymbol{g}_s\|^2 \|\boldsymbol{G}_s\|^2}{\boldsymbol{\nu}_s \sqrt{\beta_2 \boldsymbol{\nu}_{s-1}}} \right) + 8 \frac{1-\beta_1}{1-\beta_2} \frac{L_1^2}{L_0^2} \left(\sum_{s=1}^t \sqrt[8]{\beta_1^{t-s}} \left(\frac{1}{\sqrt{\beta_2 \boldsymbol{\nu}_{s-1}}} - \frac{1}{\sqrt{\boldsymbol{\nu}_s}} \right) \right).$ Here the last inequality is based on the similar reasoning of Lemma 5. The proof is completed.

В **PROOFS FOR DETERMINISTIC ALGORITHMS**

B.1 **PROOF FOR DETERMINISTIC ADAM**

We will first provide the formal statement of Theorem 1³, and then show the corresponding proof. **Theorem 6** (Theorem 1, restated). Let Assumption 1 hold. Then, $\forall \beta_1, \beta_2$ satisfying $0 \le \beta_1^2 < \beta_2 < 1$, if $T > 16L_1^2 L_0(f(w_0) - f^*)/(1 - \beta_2)$, picking $\eta = \frac{\sqrt{f(w_1) - f^*}\sqrt{1 - \frac{\beta_1^2}{\beta_2}}}{\sqrt{TL_0}(1 - \beta_1)}$, we have $\frac{1}{T} \sum_{t=1}^{I} \|\nabla f(\boldsymbol{w}_t)\| \le \frac{64}{(1-\beta_2)(1-\frac{\beta_1}{\beta_0})\left(1-\frac{\beta_1}{4/\beta_c}\right)^2} \left(\frac{\sqrt{L_0(f(\boldsymbol{w}_1)-f^*)}}{\sqrt{T}}\right).$

Proof. To begin with, according to Lemma 1 and restriction on the value of T, we obtain that

$$\forall t \in \mathbb{N} \& t \ge 1, \|\boldsymbol{w}_{t+1} - \boldsymbol{w}_t\| \le \frac{1}{4L_1}.$$

This is because by Lemma 1, $|w_{t+1} - w_t| \le \underline{\eta(1-\beta_1)/(\sqrt{1-\beta_2}\sqrt{1-\beta_1^2/\beta_2})}$, and by substituting the definition of η , we know $|w_{t+1}-w_t| \leq \sqrt{f(w_0)-f^*}/(\sqrt{1-\beta_2}\sqrt{TL_0})$. Finally, by substituting the requirement for T, we confirm the conclusion holds.

Therefore, the descent lemma can then be applied and thus $\forall t \in \mathbb{N} \& t \geq 1$,

$$f(\boldsymbol{w}_{t+1}) \leq f(\boldsymbol{w}_t) \underbrace{-\eta \left\langle \boldsymbol{G}_t, \frac{\boldsymbol{m}_t}{\lambda + \sqrt{\boldsymbol{\nu}_t}} \right\rangle}_{\text{First Order}} + \underbrace{\eta^2 \frac{L_0 + L_1 \|\boldsymbol{G}_t\|}{2} \frac{\|\boldsymbol{m}_t\|^2}{(\lambda + \sqrt{\boldsymbol{\nu}_t})^2}}_{\text{Second Order}}$$

³In the theorem below and other theorems in this paper afterward, without loss of generality, we analyze the norm version of Adam, i.e., Adam with scalar adaptive learning rate, for a more readable proof. The extension to the coordinate-wise Adam can be easily done, as evidenced by literature such as Xing et al. (2021); Faw et al. (2022; 2023); Wang et al. (2023b)

To begin with, as for the "First Order" term, according to $m_t = \beta_1 m_{t-1} + (1 - \beta_1) G_t$ we have that $-\eta \left\langle \boldsymbol{G}_{t}, \frac{\boldsymbol{m}_{t}}{\lambda + \sqrt{\boldsymbol{\nu}_{t}}} \right\rangle = -\eta \frac{1}{1 - \beta_{1}} \left\langle \boldsymbol{m}_{t}, \frac{\boldsymbol{m}_{t}}{\lambda + \sqrt{\boldsymbol{\nu}_{t}}} \right\rangle + \eta \frac{\beta_{1}}{1 - \beta_{1}} \left\langle \boldsymbol{m}_{t-1}, \frac{\boldsymbol{m}_{t}}{\lambda + \sqrt{\boldsymbol{\nu}_{t}}} \right\rangle$ $\overset{(\star)}{\leq} -\eta \frac{1}{1-\beta_1} \frac{\|\boldsymbol{m}_t\|^2}{\lambda+\sqrt{\boldsymbol{\nu}_t}} + \eta \frac{\beta_1}{(1-\beta_1)\sqrt[4]{\beta_2}} \left\langle \boldsymbol{m}_{t-1}, \frac{\boldsymbol{m}_t}{\sqrt{\lambda+\sqrt{\boldsymbol{\nu}_t}}\sqrt{\lambda+\sqrt{\boldsymbol{\nu}_{t-1}}}} \right\rangle$ $\stackrel{(*)}{\leq} -\eta \frac{1}{1-\beta_1} \frac{\|\boldsymbol{m}_t\|^2}{\lambda+\sqrt{\boldsymbol{\nu}_t}} + \frac{\beta_1}{2(1-\beta_1)\sqrt[4]{\beta_2}} \eta \frac{\|\boldsymbol{m}_t\|^2}{\lambda+\sqrt{\boldsymbol{\nu}_t}} + \frac{\beta_1}{2(1-\beta_1)\sqrt[4]{\beta_2}} \eta \frac{\|\boldsymbol{m}_{t-1}\|^2}{\lambda+\sqrt{\boldsymbol{\nu}_{t-1}}}$ $= -\eta \frac{1 - \frac{\beta_1}{\sqrt[4]{\beta_2}}}{1 - \beta_1} \frac{\|\boldsymbol{m}_t\|^2}{\lambda + \sqrt{\boldsymbol{\nu}_t}} - \frac{\beta_1}{2(1 - \beta_1)\sqrt[4]{\beta_2}} \eta \frac{\|\boldsymbol{m}_t\|^2}{\lambda + \sqrt{\boldsymbol{\nu}_t}} + \frac{\beta_1}{2(1 - \beta_1)\sqrt[4]{\beta_2}} \eta \frac{\|\boldsymbol{m}_{t-1}\|^2}{\lambda + \sqrt{\boldsymbol{\nu}_{t-1}}}.$ where inequality (*) is due to that $\sqrt{\nu_t} \ge \sqrt{\beta_2 \nu_{t-1}}$ and inequality (*) is due to Young's inequality. Meanwhile, as for the "Second Order" term, we have $_{m^2}L_0 + L_1 \|\boldsymbol{G}_t\| \|\boldsymbol{m}_t\|^2 \quad (\bullet)_{T=m^2} \quad (1-\beta_1)^2 \quad |L_1\eta^2| \|\boldsymbol{m}_t\|^2$

$$\frac{\eta}{2} \frac{1}{(\lambda + \sqrt{\nu_t})^2} \leq L_0 \eta \frac{1}{(1 - \beta_2)(1 - \frac{\beta_1^2}{\beta_2})} + \frac{1}{\sqrt{1 - \beta_2}} \frac{1}{\lambda + \sqrt{\nu_t}} \\
\stackrel{(\circ)}{\leq} L_0 \eta^2 \frac{(1 - \beta_1)^2}{(1 - \beta_2)(1 - \frac{\beta_1^2}{\beta_2})} + \frac{\eta}{2} \frac{1 - \frac{\beta_1}{\frac{4}{\beta_2}}}{1 - \beta_1} \frac{\|\boldsymbol{m}_t\|^2}{\lambda + \sqrt{\nu_t}}.$$

Here inequality (\bullet) is due to Lemma 1 and

$$\boldsymbol{\nu}_t \geq (1 - \beta_2) \|\boldsymbol{G}_t\|^2,$$

and inequality (\circ) is due to the requirement over T.

Applying the estimations of both the "First Order" and the "Second Order" terms, we obtain that

$$f(\boldsymbol{w}_{t+1}) - f(\boldsymbol{w}_{t}) \leq -\frac{\eta}{2} \frac{1 - \frac{\beta_{1}}{\sqrt[4]{\beta_{2}}}}{1 - \beta_{1}} \frac{\|\boldsymbol{m}_{t}\|^{2}}{\lambda + \sqrt{\boldsymbol{\nu}_{t}}} - \frac{\beta_{1}}{2(1 - \beta_{1})\sqrt[4]{\beta_{2}}} \eta \frac{\|\boldsymbol{m}_{t}\|^{2}}{\lambda + \sqrt{\boldsymbol{\nu}_{t}}} + \frac{\beta_{1}}{2(1 - \beta_{1})\sqrt[4]{\beta_{2}}} \eta \frac{\|\boldsymbol{m}_{t-1}\|^{2}}{\lambda + \sqrt{\boldsymbol{\nu}_{t-1}}} + L_{0}\eta^{2} \frac{(1 - \beta_{1})^{2}}{(1 - \beta_{2})(1 - \frac{\beta_{1}^{2}}{\beta_{2}})}.$$

Summing the above inequality over $t \in \{1, \dots, T\}$ then gives

$$\sum_{t=1}^{T} \frac{\eta}{2} \frac{1 - \frac{\beta_1}{\sqrt[4]{\beta_2}}}{1 - \beta_1} \frac{\|\boldsymbol{m}_t\|^2}{\lambda + \sqrt{\nu_t}}$$

$$\leq f(\boldsymbol{w}_1) - f(\boldsymbol{w}_{T+1}) - \frac{\beta_1}{2(1 - \beta_1)\sqrt[4]{\beta_2}} \eta \frac{\|\boldsymbol{m}_T\|^2}{\lambda + \sqrt{\nu_T}} + TL_0 \eta^2 \frac{(1 - \beta_1)^2}{(1 - \beta_2)(1 - \frac{\beta_1^2}{\beta_2})} \qquad (2)$$

$$\leq f(\boldsymbol{w}_1) - f(\boldsymbol{w}_{T+1}) + TL_0 \eta^2 \frac{(1 - \beta_1)^2}{(1 - \beta_2)(1 - \frac{\beta_1^2}{\beta_2})}.$$

Furthermore, as $(1 - \beta_1)G_t = m_t - \beta_1 m_{t-1}$, we have that

$$\|\boldsymbol{G}_t\|^2 \leq \frac{1}{(1-\beta_1)^2} \|\boldsymbol{m}_t\|^2 + \frac{1}{(1-\beta_1)^2} \|\boldsymbol{m}_{t-1}\|^2.$$

Applying the above inequality and $\lambda = 0$ to Eq. (2), we obtain that

$$\sum_{t=1}^{T} \frac{\eta}{4} \left(1 - \frac{\beta_1}{\sqrt[4]{\beta_2}} \right) (1 - \beta_1) \frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\boldsymbol{\nu}_t}} \le f(\boldsymbol{w}_1) - f(\boldsymbol{w}_{T+1}) + TL_0 \eta^2 \frac{(1 - \beta_1)^2}{(1 - \beta_2)(1 - \frac{\beta_1^2}{\beta_2})}.$$

Meanwhile, we have

$$\sqrt{\nu_t} - \sqrt{\beta_2 \nu_{t-1}} = \frac{(1 - \beta_2) \|\boldsymbol{G}_t\|^2}{\sqrt{\nu_t} + \sqrt{\beta_2 \nu_{t-1}}} \le (1 - \beta_2) \frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\nu_t}}.$$

Therefore, applying the above inequality and dividing both sides by η , we have

$$\frac{1}{4} \left(1 - \frac{\beta_1}{\sqrt[4]{\beta_2}} \right) (1 - \beta_1) \sum_{t=1}^T (\sqrt{\nu_t} - \sqrt{\beta_2 \nu_{t-1}}) \le \frac{f(w_1) - f(w_{T+1})}{\eta} + TL_0 \eta \frac{(1 - \beta_1)^2}{(1 - \beta_2)(1 - \frac{\beta_1^2}{\beta_2})},$$

which by telescoping further leads to

$$\frac{1}{4} \left(1 - \frac{\beta_1}{\sqrt[4]{\beta_2}} \right) (1 - \beta_1) \sum_{t=1}^T (1 - \beta_2) \sqrt{\nu_t} \le \frac{f(w_1) - f(w_{T+1})}{\eta} + TL_0 \eta \frac{(1 - \beta_1)^2}{(1 - \beta_2)(1 - \frac{\beta_1^2}{\beta_2})}.$$

According to Cauchy-Schwartz's inequality, we then obtain

$$\begin{split} \left(\sum_{t=1}^{T} \|\boldsymbol{G}_{t}\|\right)^{2} &\leq \left(\sum_{t=1}^{T} \sqrt{\nu_{t}}\right) \left(\sum_{t=1}^{T} \frac{\|\boldsymbol{G}_{t}\|^{2}}{\sqrt{\nu_{t}}}\right) \\ &\leq \frac{1}{1-\beta_{2}} \left(\frac{4(f(\boldsymbol{w}_{1})-f(\boldsymbol{w}_{T+1}))}{\eta\left(1-\frac{\beta_{1}}{\sqrt{\beta_{2}}}\right)(1-\beta_{1})} + TL_{0}\eta\frac{(1-\beta_{1})}{(1-\beta_{2})\left(1-\frac{\beta_{1}}{\sqrt{\beta_{2}}}\right)(1-\frac{\beta_{1}^{2}}{\beta_{2}})}\right)^{2} \\ &= \frac{1}{1-\beta_{2}} \left(\frac{4(f(\boldsymbol{w}_{1})-f(\boldsymbol{w}_{T+1}))}{\eta\left(1-\frac{\beta_{1}}{\sqrt{\beta_{2}}}\right)(1-\beta_{1})} + 4TL_{0}\eta\frac{(1-\beta_{1})}{(1-\beta_{2})\left(1-\frac{\beta_{1}}{\sqrt{\beta_{2}}}\right)(1-\frac{\beta_{1}^{2}}{\beta_{2}})}\right)^{2}. \end{split}$$

The proof is completed by applying the value of η .

The proof is completed by applying the value of η .

B.2 PROOF FOR GDM

This section collects the proof of Theorem 2. To begin with, given problem hyperparameters Δ_1 , ε , L_0 , and L_1 . We first construct three 1D functions as follows:

$$f_{1}(x) = \begin{cases} \frac{L_{0}e^{L_{1}x-1}}{L_{1}^{2}} , x \in \left[\frac{1}{L_{1}}, \infty\right), \\ \frac{L_{0}x^{2}}{2} + \frac{L_{0}}{2L_{1}^{2}} , x \in \left[-\frac{1}{L_{1}}, \frac{1}{L_{1}}\right], \\ \frac{L_{0}e^{-L_{1}x-1}}{L_{1}^{2}} , x \in \left(-\infty, -\frac{1}{L_{1}}\right]. \end{cases}$$

$$f_{2}(y) = \begin{cases} \varepsilon(y-1) + \frac{\varepsilon}{2} , y \in [1,\infty), \\ \frac{\varepsilon}{2}y^{2} , y \in [-1,1], \\ -\varepsilon(y+1) + \frac{\varepsilon}{2} , y \in (-\infty, -1]. \end{cases}$$

$$f_{3}(z) = \begin{cases} \varepsilon(z-1) + \frac{\varepsilon}{2L_{1}} + \frac{L_{0}}{2L_{1}^{2}} , z \in \left[\frac{1}{L_{1}}, \infty\right), \\ \frac{\varepsilonL_{1}z^{2}}{2} + \frac{L_{0}}{2L_{1}^{2}} , z \in [0, \frac{1}{L_{1}}], \\ \frac{L_{0}e^{-L_{1}z-1}}{L_{1}^{2}} , z \in \left[-\frac{1}{L_{1}}, 0\right], \\ \frac{L_{0}e^{-L_{1}z-1}}{L_{1}^{2}} , z \in \left(-\infty, -\frac{1}{L_{1}}\right]. \end{cases}$$
(3)

It is easy to verify that these functions satisfy (L_0, L_1) -smooth condition as long as $\varepsilon \leq L_0$. We then respectively the convergence of GDM over these three examples with different learning rate and momentum coefficient.

 $\begin{array}{ll} \textbf{P18} \\ \textbf{P19} \\ \textbf{P19} \\ \textbf{P20} \\ \textbf{P20} \\ \textbf{P21} \\ \textbf{P21} \\ \textbf{P22} \\ \textbf{P22} \\ \textbf{P23} \\ \textbf{P23} \\ \textbf{P24} \\ \textbf{P24} \\ \textbf{L}_1 \\ \textbf{L}_1 \\ \textbf{L}_1 \\ \textbf{L}_1 \\ \textbf{L}_1 \\ \textbf{L}_2 \\ \textbf{L}_2 \\ \textbf{L}_1 \\ \textbf{L}_2 \\ \textbf{L}_2 \\ \textbf{L}_1 \\ \textbf{L}_2 \\ \textbf{L}_2 \\ \textbf{L}_1 \\ \textbf{L}_1 \\ \textbf{L}_2 \\ \textbf{L}_2 \\ \textbf{L}_1 \\ \textbf{L}_1 \\ \textbf{L}_2 \\ \textbf{L}_2 \\ \textbf{L}_1 \\ \textbf{L}_1 \\ \textbf{L}_1 \\ \textbf{L}_2 \\ \textbf{L}_1 \\ \textbf{L}$

Proof. We prove this lemma by proving that $\forall k \ge 1$, $|x_{k+1}| \ge (4 + 8 \log \frac{1}{\varepsilon})|x_k|$ and $\operatorname{Sign}(x_{k+1}) = (-1)^{k+1}$ by induction. When k = 1, according to the update rule of GDM, we have

$$x_2 = x_1 - \eta f_1'(x_1)$$

As
$$\eta \geq \frac{(5+8\log\frac{1}{\varepsilon})(1+\log(\frac{1}{2}+\frac{L_1^2}{L_0}\Delta_1))}{\Delta_1+\frac{L_0}{2L_1^2}} = -\frac{(5+8\log\frac{1}{\varepsilon})x_1}{f_1'(x_1)}$$
, we have

$$x_2 \le -(4+8\log\frac{1}{\varepsilon})x_1$$

936 which leads to the claim.

Now assuming that the claim has been proved for $k \le t - 1$ ($t \ge 2$). Then, for k = t, with induction hypothesis we have

$$x_{t+1} = x_t - \eta \mathbf{m}_t = x_t - \eta \left(\beta^t f_1'(x_1) + (1-\beta) \sum_{s=1}^{t-1} \beta^{t-s} f_1'(x_s) + (1-\beta) f_1'(x_t) \right)$$

Without the loss of generality, we assume t is even. By the induction hypothesis, we obtain that $f'_1(x_t) < 0$ and $f'_1(x_{t-1}) < 0$, and

$$|f_1'(x_1)| \le |f_1'(x_2)| \le \dots \le |f_1'(x_{t-1})|.$$

Therefore, we have

$$\begin{aligned} x_{t+1} \ge & x_t - \eta \left(\beta f_1'(x_{t-1}) + (1-\beta)f_1'(x_t)\right) \\ = & x_t - \frac{L_0}{L_1}\eta \left(\beta e^{L_1 x_{t-1} - 1} - (1-\beta)e^{-L_1 x_t - 1}\right) \\ \ge & x_t - \frac{L_0}{L_1}\eta \left(\beta e^{-\frac{L_1 x_t}{8 \log \frac{1}{\varepsilon} + 4} - 1} - (1-\beta)e^{-L_1 x_t - 1}\right) \end{aligned}$$

Furthermore, according to the definition of x_1 , we have

$$1 - \beta \ge 2e^{-L_1(4\log\frac{1}{\varepsilon} + 2)x_1} \ge 2e^{\frac{L_1x_t}{2}}$$

which leads to

$$x_{t+1} \ge x_t + \frac{L_0}{L_1} \eta e^{-\frac{L_1 x_t}{2} - 1} \ge x_t + \frac{(5 + 8\log\frac{1}{\varepsilon})x_1}{e^{L_1 x_1}} e^{-\frac{L_1 x_t}{2}} \ge x_t + \frac{(5 + 8\log\frac{1}{\varepsilon})x_1}{e^{L_1 x_1}} e^{L_1 x_t (2 + 4\log\frac{1}{\varepsilon})}.$$

Then, as $\frac{e^{\frac{L_1x}{2}}}{x}$ is monotonously increasing for $x \in [\frac{2}{L_1}, \infty)$, and $x_1 \ge \frac{2}{L_1}$, we have

$$x_{t+1} \ge x_t + \frac{(5 + 8\log\frac{1}{\varepsilon})x_1}{e^{L_1 x_1}} e^{L_1 x_t (1 + 2\log\frac{1}{\varepsilon})} \ge x_t - (5 + 8\log\frac{1}{\varepsilon})x_t \ge -(4 + 8\log\frac{1}{\varepsilon})x_t.$$

The proof is completed.

 $\begin{array}{ll} \textbf{PG9} \\ \textbf{970} \\ \textbf{970} \\ \textbf{971} \\ \hline \frac{(5+8\log\frac{1}{\varepsilon})(1+\log(\frac{1}{2}+\frac{L_1^2}{L_0}\Delta_1))}{L_1^2(\Delta_1+\frac{L_0}{2L_1^2})}, \text{ we have that GDM satisfies that } \|\nabla f_2(y_t)\| \geq \varepsilon \text{ if } T \leq \tilde{\Theta}(\frac{L_1^2\Delta_1^2+L_0\Delta_1}{\varepsilon^2}). \end{array}$

973 *Proof.* We have that $m_t = \varepsilon$ before y_t enters the region $(-\infty, 1]$. As the movement of each step 973 before y_t enters the region $(-\infty, 1]$ is $\eta\varepsilon$ and the total length to enter $(-\infty, 1]$ is $y_1 - 1$, the proof is 974 completed.

Lemma 8 (Convergence over f_3). Assume $\Delta_1 \geq \frac{L_0}{L_1^2}e + 4e + \frac{L_0^2}{e^2L_1^2}$, $L_1 \geq 1$, $\varepsilon \leq \frac{1}{2}$, and let $z_1 = -\frac{1 + \log(\frac{1}{2} + \frac{L_1^2}{L_0}\Delta_1)}{L_1}$. Then, we have $f_3(z_1) - f_3^* = \Delta_1$, and if $\eta \geq \frac{(5+8\log\frac{1}{\varepsilon})(1 + \log(\frac{1}{2} + \frac{L_1^2}{L_0}\Delta_1))}{L_1^2(\Delta_1 + \frac{L_0}{2L_1^2})}$ and $\beta \geq 1 - 2\left(\frac{L_1^2}{L_0}e\right)^{-4\log\frac{1}{\varepsilon}-2} (\Delta_1 + \frac{L_0}{2L_1^2})^{-4\log\frac{1}{\varepsilon}-2}$, we have that GDM satisfies that $\forall t \in [1, \Theta(\frac{L_1^2\Delta_1^2}{\varepsilon^3})), |f_3'(x_t)| \geq \varepsilon$.

Proof. To begin with, according to the definition of z_1 , we have $\eta \ge -\frac{(5+8\log \frac{1}{\varepsilon})z_1}{f'_3(x_1)}$ and $1-\beta \ge 2e^{L_1(4\log \frac{1}{\varepsilon}+2)z_1} \ge \frac{1}{2}$. Also, as $\Delta_1 \ge \frac{L_0}{L_1^2}(e-\frac{1}{2})$, we have $z_1 \le -\frac{2}{L_1}$, and thus

$$f_3'(z_1) = -\frac{L_0}{L_1} e^{-L_1 z_1 - 1} \le -L_1 \left(\Delta_1 + \frac{L_0}{2L_1^2} \right) \le -4.$$

We will first prove the following claim by induction: for $k \in [2, \lfloor \frac{1}{1-\beta} \rfloor]$, we have $z_k \ge \frac{1}{L_1}$, and $m_k \le \frac{\beta^{k-1} f'_3(z_1)}{2}$.

As for k = 2, we have

$$z_2 = z_1 - \eta f'_3(z_1) \ge -\left(4 + 8\log\frac{1}{\varepsilon}\right) z_1.$$

According to $\Delta_1 \geq \frac{L_0}{L_1^2}(e-\frac{1}{2})$, we have $z_1 \leq -\frac{2}{L_1}$, and thus $z_2 \geq \frac{1}{L_1}$. Since $m_2 = \beta f'(z_1) + (1-\beta)\varepsilon < \frac{f'_3(z_1)}{2}$, the claim is proved for k = 2.

Now assuming that we have prove the claim for $k \le t - 1$. According to the induction hypothesis, we have

$$f_3'(z_2) = \cdots = f_3'(z_{t-1}) = \varepsilon,$$

1003 and thus

$$\boldsymbol{m}_t = \beta^{t-1} f_3'(z_1) + (1 - \beta^{t-1}) \varepsilon \stackrel{(\star)}{\leq} \beta^{t-1} f_3'(z_1) - \frac{\beta^{t-1} f_3'(z_1)}{2} \le \frac{\beta^{t-1} f_3'(z_1)}{2}.$$

Here inequality (*) is due to $\beta^{\lfloor \frac{1}{1-\beta} \rfloor} \ge \frac{1}{4}$ as $\beta \ge \frac{1}{2}$. Therefore, as $z_t = z_{t-1} - \eta m_t \ge z_{t-1} \ge \frac{1}{L_1}$, we prove the claim.

It should be noticed that $\forall t \in [1, \lfloor \frac{1}{1-\beta} \rfloor], ||f'_3(z_t)| \geq \varepsilon$. Furthermore, according to the claim, $z_{\lfloor \frac{1}{1-\beta} \rfloor+1}$ can now be bounded as

$$z_{\lfloor \frac{1}{1-\beta} \rfloor+1} = z_1 - \eta \sum_{k=1}^{\lfloor \frac{1}{1-\beta} \rfloor} m_t \ge \frac{\eta}{5+8\log \frac{1}{\varepsilon}} f'_3(z_1) - \eta \sum_{k=1}^{\lfloor \frac{1}{1-\beta} \rfloor} \frac{\beta^{k-1} f'_3(z_1)}{2} \ge \frac{\eta}{5+8\log \frac{1}{\varepsilon}} f'_3(z_1) - \eta \frac{1-\frac{1}{e}}{(1-\beta)} \frac{f'_3(z_1)}{2} \ge \frac{\eta}{5+8\log \frac{1}{\varepsilon}} f'_3(z_1) - \eta \frac{1-\frac{1}{e}}{(1-\beta)} \frac{f'_3(z_1)}{2} \ge \frac{1}{L_1} - \eta \frac{1-\frac{1}{e}}{(1-\beta)} \frac{f'_3(z_1)}{4} \ge \frac{1}{L_1} - \eta \left(1-\frac{1}{e}\right) \frac{f'_3(z_1)}{8} \left(\frac{L_1^2}{L_0}e\right)^{4\log \frac{1}{\varepsilon}+2} \left(\Delta_1 + \frac{L_0}{2L_1^2}\right)^{4\log \frac{1}{\varepsilon}+2} \ge \frac{1}{L_1} + \frac{\eta}{16} \frac{L_1^2 \Delta_1^2 + L_0 \Delta_1}{\varepsilon^2}.$$

1022 As $f'_3(z) = \varepsilon$ for all $z \ge \frac{1}{L_1}$, the iterates needs additional $\frac{\frac{\eta}{16} \frac{L_1^2 \Delta_1^2}{\varepsilon^2}}{\eta \varepsilon} = \frac{1}{16} \frac{L_1^2 \Delta_1^2}{\varepsilon^3}$ steps to make $f'_3(z_t) < \varepsilon$. The proof is completed.

Lemma 9. Let $f_1, f_2, f_3 : \mathcal{R} \to \mathcal{R}$ satisfies (L_0, L_1) -smooth condition. Then $f_1(x) + f_2(y) + f_3(z)$ satisfies (L_0, L_1) -smooth condition.

1027 Proof. If we consider a point (x_1, y_1, z_1) within a ball centered at (x_2, y_2, z_2) with radius $1/L_1$, it 1028 follows that x_1 is within a ball centered at x_2 with the same radius. Thus, we have:

$$|\nabla f_1(x_1) - \nabla f_1(x_2)| \le (L_0 + L_1 |\nabla f_1(x_1)|) |x_1 - x_2| \le (L_0 + L_1 |\nabla f(x_1, y_1, z_1)|) |x_1 - x_2|.$$

1032 The last inequality holds because $\nabla f_1(x_1)$ is one coordinate of $\nabla f(x_1, y_1, z_1)$. Similarly, we can derive:

 $|\nabla f_2(y_1) - \nabla f_2(y_2)| \le (L_0 + L_1 |\nabla f(x_1, y_1, z_1)|) |y_1 - y_2|,$

$$|\nabla f_3(z_1) - \nabla f_3(z_2)| \le (L_0 + L_1 |\nabla f(x_1, y_1, z_1)|) |z_1 - z_2|.$$

Taking the squared sum of these inequalities confirms that f is indeed (L_0, L_1) -smooth.

Theorem 7 (Theorem 2, restated). Assume that $\Delta_1 \ge 4\frac{L_0}{L_1}e + 16e + 4\frac{L_0^2}{e^2L_1^2}$, $L_1 \ge 1$ and $\varepsilon \le 1$, then there exists objective function f satisfying (L_0, L_1) -smooth condition and $f(w_1) - f^* = \Delta_1$, such that for any learning rate $\eta > 0$ and $\beta \in [0, 1]$, the minimum step T of GDM to achieve final error ε satisfies $(L^2 \Delta^2 + L_1 \Delta_1)$

$$T = \tilde{\Omega}\left(\frac{L_1^2 \Delta_1^2 + L_0 \Delta_1}{\varepsilon^2}\right)$$

The proof is completed.

B.3 PROOF FOR DETERMINISTIC ADAGRAD

¹⁰⁵⁷ To begin with, we recall the following result from Wang et al. (2023b):

Proposition 2. For every learning rate $\eta \ge \Theta(\frac{1}{L_1})$ and Δ_1 , there exist a lower-bounded objective function g_1 obeying Assumption 1 and a corresponding initialization point w_0 with $g_1(w_1) - g_1^* = \Delta_1$, such that AdaGrad with learning rate η and initialized at w_0 diverges over g_1 .

1062 We then define g_2 as the f_2 in the proof of Theorem 2, i.e.,

$$g_2(y) = \begin{cases} \varepsilon(y-1) + \frac{\varepsilon}{2} & , y \in [1,\infty), \\ \frac{\varepsilon}{2}y^2 & , y \in [-1,1], \\ -\varepsilon(y+1) + \frac{\varepsilon}{2} & , y \in (-\infty,-1]. \end{cases}$$
(6)

We then have the following lemma characterizing the convergence of AdaGrad over g_2 .

Lemma 10 (Convergence over g_2). Assume that $\Delta_1 \geq \frac{\varepsilon}{2} + \frac{L_1}{L_0}$, and let $y_1 \triangleq \frac{\Delta_1}{\varepsilon} + \frac{1}{2}$. Then, if $\eta \leq \Theta(\frac{1}{L_1})$, we have that AdaGrad satisfies that $\|\nabla g_2(y_t)\| \geq \varepsilon$ if $T \leq \tilde{\Theta}(\frac{L_1^2 \Delta_1^2}{\varepsilon^2})$.

Proof. We have that $g_t = \varepsilon$ before y_t enters the region $(-\infty, 1]$. Therefore, the sum of movement of 1075 each step before y_t enters the region $(-\infty, 1]$ is

1077
1078
$$\eta \sum_{s=1}^{t} \frac{\varepsilon}{\sqrt{s\varepsilon}} = \eta \Theta(\sqrt{t}).$$

Solving $\eta \Theta(\sqrt{t}) = \frac{\Delta_1}{\varepsilon} + \frac{1}{2} - 1$ gives $t = \frac{L_1^2 \Delta_1^2}{\varepsilon^2}$, and the proof is completed.

¹⁰⁸⁰ We then have the following lower bound for deterministic AdaGrad.

Theorem 8. Assume that $\Delta_1 \geq \frac{\varepsilon}{2} + \frac{L_1}{L_0}$. Then, there exists objective function f satisfying (L_0, L_1) smooth condition and $f(w_1) - f^* = \Delta_1$, such that for any learning rate $\eta > 0$ and $\beta \in [0, 1]$, the minimum step T of AdaGrad to achieve final error ε satisfies

$$T = \Omega(\frac{L_1^2 \Delta_1^2}{\varepsilon^2}).$$

Proof. The proof is completed by letting $f(x, y) = g_1(x) + g_2(y)$ following the same routine as Theorem 7.

C PROOF FOR STOCHASTIC ALGORITHMS

C.1 PROOF FOR ADAM

⁷ To begin with, we restate the theorem as follows:

Theorem 9 (Theorem 3, restated). Let Assumptions 1 and 2 hold. Then, $\forall \beta_1 \ge 0$ and $\lambda = 0$, if $\varepsilon \leq \frac{1}{\operatorname{poly}(f(\boldsymbol{w}_1) - f^*, L_0, L_1, \sigma_0, \sigma_1)}, \text{ with } \eta = \frac{\sqrt{f(\boldsymbol{w}_1) - f^*}}{\sqrt{L_0 + L_1}\sqrt{T\sigma_0\sigma_1^2}} \text{ and momentum hyperparameter } \beta_2 = 1$ $1 - \eta^{2} \left(\frac{1024\sigma_{1}^{2}(L_{1}+L_{0})(1-\beta_{1})}{\sqrt{1 - \frac{\beta_{1}^{2}}{\beta_{2}}(1 - \frac{\beta_{1}}{\sqrt{\beta_{2}}})}} \right)^{2}, \text{ we have if } T \geq \Theta \left(\frac{(L_{0}+L_{1})\sigma_{0}^{3}\sigma_{1}^{2}(f(\boldsymbol{w}_{1})-f^{*})}{\varepsilon^{4}} \right), \text{ then Algorithm 1}$ satisfies $\frac{1}{T}\mathbb{E}\sum_{t=1}^{T} \|\nabla f(\boldsymbol{w}_t)\| \leq \varepsilon.$ *Proof.* Let the approximate iterative sequence be defined as $u_t \triangleq \frac{w_t - \frac{p_1}{\sqrt{\beta_2}} w_{t-1}}{\frac{1 - \frac{\beta_1}{\sqrt{\beta_2}}}{1 - \frac{\beta_1}{\sqrt{\beta_2}}}}$ and the surro-

gate second-order momentum be defined as $\tilde{\nu}_t \triangleq \beta_2 \nu_{t-1} + (1 - \beta_2) \sigma_0^2$. Then, as $\frac{\eta}{\sqrt{1 - \beta_2}} = \sqrt{1 - \frac{\beta_1^2}{2}} (1 - \frac{\beta_1}{2})$

15
$$\frac{\sqrt{\beta_2}}{1024\sigma_1^2(L_1+L_0)(1-\beta_1)}$$
, we have

$$\|\boldsymbol{u}_t - \boldsymbol{w}_t\| = \frac{\frac{\beta_1}{\sqrt{\beta_2}}}{1 - \frac{\beta_1}{\sqrt{\beta_2}}} \|\boldsymbol{w}_t - \boldsymbol{w}_{t-1}\| \stackrel{(*)}{\leq} \eta \frac{\frac{\beta_1}{\sqrt{\beta_2}}}{1 - \frac{\beta_1}{\sqrt{\beta_2}}} \frac{1 - \beta_1}{\sqrt{1 - \beta_2}\sqrt{1 - \frac{\beta_1^2}{\beta_2}}} \le \frac{1}{4L_1},$$

1121 and

$$\|\boldsymbol{u}_{t+1} - \boldsymbol{w}_t\| = \frac{1}{1 - \frac{\beta_1}{\sqrt{\beta_2}}} \|\boldsymbol{w}_{t+1} - \boldsymbol{w}_t\| \stackrel{(*)}{\leq} \eta \frac{1}{1 - \frac{\beta_1}{\sqrt{\beta_2}}} \frac{1 - \beta_1}{\sqrt{1 - \beta_2}\sqrt{1 - \frac{\beta_1^2}{\beta_2}}} \leq \frac{1}{4L_1}.$$

1127 Therefore, if choosing $w^1 = w_t$, $w^2 = u_{t+1}$, and $w^3 = u_t$ in Lemma 2, we see the conditions of 1128 Lemma 2 is satisfied, which after taking expectation gives

1129
1130
$$\mathbb{E}^{|\mathcal{F}_{t}}f(\boldsymbol{u}_{t+1}) \leq f(\boldsymbol{u}_{t}) + \mathbb{E}^{|\mathcal{F}_{t}}\langle \nabla f(\boldsymbol{w}_{t}), \boldsymbol{u}_{t+1} - \boldsymbol{u}_{t} \rangle + \frac{1}{2}(L_{0} + L_{1} \|\nabla f(\boldsymbol{w}_{t})\|)\mathbb{E}^{|\mathcal{F}_{t}|}(\|\boldsymbol{u}_{t+1} - \boldsymbol{w}_{t}\| + \|\boldsymbol{u}_{t} - \boldsymbol{w}_{t}\|)\|\boldsymbol{u}_{t+1} - \boldsymbol{u}_{t}\|.$$
1131

1132 We call $\langle \nabla f(\boldsymbol{w}_t), \boldsymbol{u}_{t+1} - \boldsymbol{u}_t \rangle$ the first-order term and $\frac{1}{2}(L_0 + L_1 \|\nabla f(\boldsymbol{w}_t)\|)(\|\boldsymbol{u}_{t+1} - \boldsymbol{w}_t\| + \|\boldsymbol{u}_t - \boldsymbol{u}_t\|)$

 $w_t \| \| \| u_{t+1} - u_t \|$ the second-order term, as they respectively correspond to the first-order and second-order Taylor's expansion. We then respectively bound these two terms as follows.

Analysis for the first-order term. Before we start, denote $\tilde{\nu}_t \triangleq \beta_2 \nu_{t-1} + (1 - \beta_2) \sigma_0^2$ $u_{t+1} - u_t = \frac{w_{t+1} - w_t}{1 - \frac{\beta_1}{\sqrt{\beta_2}}} - \frac{\beta_1}{\sqrt{\beta_2}} \frac{w_t - w_{t-1}}{1 - \frac{\beta_1}{\sqrt{\beta_2}}}$ $=-\frac{\eta}{1-\frac{\beta_1}{\sqrt{\nu_t}}}\frac{1}{\sqrt{\nu_t}}\boldsymbol{m}_t+\beta_1\frac{\eta}{1-\frac{\beta_1}{\sqrt{\beta_2}}}\frac{1}{\sqrt{\beta_2\nu_{t-1}}}\boldsymbol{m}_{t-1}$ $= -rac{\eta}{1-rac{eta_1}{\sqrt{m{arepsilon}_t}}}rac{1}{\sqrt{m{
u}_t}}m{m}_t + eta_1rac{\eta}{1-rac{eta_1}{\sqrt{m{m{arepsilon}_t}}}}rac{1}{\sqrt{m{
u}_t}}m{m}_{t-1} - rac{\eta}{1-rac{eta_1}{\sqrt{m{m{m{m{m{\mu}}}_t}}}}\left(rac{1}{\sqrt{m{
u}_t}}-rac{1}{\sqrt{m{
u}_t}}
ight)m{m}_t$ $+ \beta_1 \frac{\eta}{1 - \frac{\beta_1}{\sqrt{\alpha}}} \left(\frac{1}{\sqrt{\beta_2 \nu_{t-1}}} - \frac{1}{\sqrt{\widetilde{\nu}_t}} \right) m_{t-1}$ $=-\eta \frac{1-\beta_1}{1-\frac{\beta_1}{\sqrt{\tilde{\boldsymbol{\nu}}_t}}} \frac{1}{\sqrt{\tilde{\boldsymbol{\nu}}_t}} \boldsymbol{g}_t - \frac{\eta}{1-\frac{\beta_1}{\sqrt{\tilde{\boldsymbol{\mu}}_t}}} \left(\frac{1}{\sqrt{\boldsymbol{\nu}_t}} - \frac{1}{\sqrt{\tilde{\boldsymbol{\nu}}_t}}\right) \boldsymbol{m}_t + \beta_1 \frac{\eta}{1-\frac{\beta_1}{\sqrt{\tilde{\boldsymbol{\mu}}_t}}} \left(\frac{1}{\sqrt{\beta_2 \boldsymbol{\nu}_{t-1}}} - \frac{1}{\sqrt{\tilde{\boldsymbol{\nu}}_t}}\right) \boldsymbol{m}_{t-1}.$ According to the above decomposition, we have the first-order term can also be decomposed into

 $\mathbb{E}^{|\mathcal{F}_t|}[\langle \nabla f(\boldsymbol{w}_t), \boldsymbol{u}_{t+1} - \boldsymbol{u}_t \rangle]$ $= \frac{1 - \beta_1}{1 - \frac{\beta_1}{\sqrt{\mathbf{z}_t}}} \mathbb{E}^{|\mathcal{F}_t} \left[\left\langle \boldsymbol{G}_t, -\eta \frac{1}{\sqrt{\tilde{\boldsymbol{\nu}}_t}} \boldsymbol{g}_t \right\rangle \right] + \mathbb{E}^{|\mathcal{F}_t} \left| \left\langle \boldsymbol{G}_t, -\frac{\eta}{1 - \frac{\beta_1}{\sqrt{\mathbf{z}_t}}} \left(\frac{1}{\sqrt{\boldsymbol{\nu}_t}} - \frac{1}{\sqrt{\tilde{\boldsymbol{\nu}}_t}} \right) \boldsymbol{m}_t \right\rangle \right|$ $+ \mathbb{E}^{|\mathcal{F}_t|} \left| \left\langle \boldsymbol{G}_t, eta_1 rac{\eta}{1 - rac{eta_1}{1 - rac{eta_1}{1 - rac{eta_1}{2}}} \left(rac{1}{\sqrt{eta_2 oldsymbol{
u}_{t-1}}} - rac{1}{\sqrt{\widetilde{oldsymbol{
u}}_t}}
ight) oldsymbol{m}_{t-1}
ight
angle
ight|.$ (7)

1159
1160 As
$$\mathbb{E}^{|\mathcal{F}_t} \left[\left\langle \boldsymbol{G}_t, -\eta \frac{1}{\sqrt{\tilde{\nu}_t}} \boldsymbol{g}_t \right\rangle \right] = -\eta \frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\tilde{\nu}_t}}$$
, we have
1161
1162 $\frac{1-\beta_1}{1-\frac{\beta_1}{\sqrt{\beta_2}}} \mathbb{E}^{|\mathcal{F}_t} \left[\left\langle \boldsymbol{G}_t, -\eta \frac{1}{\sqrt{\tilde{\nu}_t}} \boldsymbol{g}_t \right\rangle \right] \leq -\frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\tilde{\nu}_t}}.$
1164

We then respectively bound the rest of the two terms in Eq. (7). To begin with,

$$\begin{aligned} & 1166 \\ & 1167 \\ & \mathbb{E}^{|\mathcal{F}_{t}} \left[\left\langle \boldsymbol{G}_{t}, -\frac{\eta}{1-\frac{\beta_{1}}{\sqrt{\beta_{2}}}} \left(\frac{1}{\sqrt{\nu_{t}}} - \frac{1}{\sqrt{\tilde{\nu}_{t}}} \right) \boldsymbol{m}_{t} \right\rangle \right] \\ & 1169 \\ & 1169 \\ & 1170 \\ & = \mathbb{E}^{|\mathcal{F}_{t}} \left[\left\langle \boldsymbol{G}_{t}, -\frac{\eta}{1-\frac{\beta_{1}}{\sqrt{\beta_{2}}}} \left(\frac{(1-\beta_{2})(\sigma_{0}^{2} - \|\boldsymbol{g}_{t}\|^{2})}{\sqrt{\nu_{t}}\sqrt{\tilde{\nu}_{t}}(\sqrt{\nu_{t}} + \sqrt{\tilde{\nu}_{t}})} \right) \boldsymbol{m}_{t} \right\rangle \right] \\ & 1171 \\ & 1172 \\ & \leq \frac{\eta}{1-\frac{\beta_{1}}{\sqrt{\beta_{2}}}} \mathbb{E}^{|\mathcal{F}_{t}} \left[\left\| \boldsymbol{G}_{t} \right\| \left(\frac{(1-\beta_{2})(\sigma_{0}^{2} + \|\boldsymbol{g}_{t}\|^{2})}{\sqrt{\nu_{t}}\sqrt{\tilde{\nu}_{t}}(\sqrt{\nu_{t}} + \sqrt{\tilde{\nu}_{t}})} \right) \|\boldsymbol{m}_{t}\| \right] \\ & 1174 \\ & = \frac{\eta}{1-\frac{\beta_{1}}{\sqrt{\beta_{2}}}} \mathbb{E}^{|\mathcal{F}_{t}} \left[\left\| \boldsymbol{G}_{t} \right\| \left(\frac{(1-\beta_{2})\|\boldsymbol{g}_{t}\|^{2}}{\sqrt{\nu_{t}}\sqrt{\tilde{\nu}_{t}}(\sqrt{\nu_{t}} + \sqrt{\tilde{\nu}_{t}})} \right) \|\boldsymbol{m}_{t}\| \right] + \frac{\eta}{1-\frac{\beta_{1}}{\sqrt{\beta_{2}}}} \mathbb{E}^{|\mathcal{F}_{t}} \left[\left\| \boldsymbol{G}_{t} \right\| \left(\frac{(1-\beta_{2})\sigma_{0}^{2}}{\sqrt{\nu_{t}}\sqrt{\tilde{\nu}_{t}}(\sqrt{\nu_{t}} + \sqrt{\tilde{\nu}_{t}})} \right) \|\boldsymbol{m}_{t}\| \right] \\ & (8) \end{aligned}$$

The first term in the right-hand-side of Eq. (8) can be bounded as

where inequality (*) uses Lemma 1, inequality (
$$\circ$$
) is due to Holder's in-
equality, and inequality (\bullet) is due to Assumption 2. Applying mean-
value inequality respectively to $\frac{\eta(1-\beta_1)\sqrt{1-\beta_2}}{\left(\sqrt{1-\beta_1}\right)^3} \mathbb{E}^{|\mathcal{F}_t|} \frac{||\mathbf{G}_t||}{\sqrt{\nu_t}} \sigma_0 \sqrt{\mathbb{E}^{|\mathcal{F}_t|} \frac{||\mathbf{g}_t||^2}{(\sqrt{\nu_t}+\sqrt{\nu_t})^2}}$ and
 $\frac{\eta(1-\beta_1)\sqrt{1-\beta_2}}{\left(\sqrt{1-\beta_2}\right)^3} \mathbb{E}^{|\mathcal{F}_t|} \frac{||\mathbf{G}_t||}{\sqrt{\nu_t}} \sigma_1 ||\mathbf{G}_t|| \sqrt{\mathbb{E}^{|\mathcal{F}_t|} \frac{||\mathbf{g}_t||^2}{(\sqrt{\nu_t}+\sqrt{\nu_t})^2}}}$ and due to $\beta_1 \leq \beta_2$, we obtain that the
right-hand-side of the above inequality can be bounded by
 $\frac{1}{16} \eta \frac{1-\beta_1}{1-\frac{\beta_1}{\sqrt{\beta_2}}} \sqrt{1-\beta_2} \sigma_0 \frac{||\mathbf{G}_t||^2}{\tilde{\nu}_t} + \frac{4\eta\sqrt{1-\beta_2}\sigma_0}{\left(1-\frac{\beta_1}{\sqrt{\beta_2}}\right)^2} \mathbb{E}^{|\mathcal{F}_t|} \frac{||\mathbf{g}_t||^2}{(\sqrt{\nu_t}+\sqrt{\tilde{\nu}_t})^2}$
 $+ \frac{1}{16} \eta \frac{1-\beta_1}{1-\frac{\beta_1}{\sqrt{\beta_2}}} \frac{||\mathbf{G}_t||^2}{\sqrt{\tilde{\nu}_t}} + 4\eta \frac{(1-\beta_2)(1-\beta_1)}{(1-\frac{\beta_1}{\sqrt{\beta_2}})^2} \sigma_1^2 \frac{||\mathbf{G}_t||^2}{\sqrt{\tilde{\nu}_t}} \mathbb{E}^{|\mathcal{F}_t|} \frac{||\mathbf{g}_t||^2}{(\sqrt{\nu_t}+\sqrt{\tilde{\nu}_t})^2}$
 $\leq \frac{1}{8} \eta \frac{||\mathbf{G}_t||^2}{\sqrt{\tilde{\nu}_t}} + \frac{4\eta\sqrt{1-\beta_2}\sigma_0}{\left(1-\frac{\beta_1}{\sqrt{\beta_2}}\right)^2} \mathbb{E}^{|\mathcal{F}_t|} \frac{||\mathbf{g}_t||^2}{\nu_t} + \frac{1}{8} \eta \frac{||\mathbf{G}_t||^2}{\sqrt{\tilde{\nu}_t}} + 16\eta \frac{(1-\beta_2)}{(1-\beta_1)} \sigma_1^2 \frac{||\mathbf{G}_t||^2}{\sqrt{\tilde{\nu}_t}} \mathbb{E}^{|\mathcal{F}_t|} \frac{||\mathbf{g}_t||^2}{(\sqrt{\nu_t}+\sqrt{\tilde{\nu}_t})^2}.$
(9)
Here the inequality is due to $\tilde{\nu}_t = (1-\beta_2)\sigma_0^2 + \beta_2\nu_{t-1} \ge (1-\beta_2)\sigma_0^2$. Meanwhile, we have

$$\begin{split} & \left(\frac{1}{\sqrt{\beta_2 \widetilde{\boldsymbol{\nu}}_t}} - \frac{1}{\sqrt{\widetilde{\boldsymbol{\nu}}_{t+1}}}\right) \|\boldsymbol{G}_t\|^2 \\ &= \frac{\|\boldsymbol{G}_t\|^2 ((1-\beta_2)^2 \sigma_0^2 + \beta_2 (1-\beta_2) \|\boldsymbol{g}_t\|^2)}{\sqrt{\beta_2 \widetilde{\boldsymbol{\nu}}_t} \sqrt{\widetilde{\boldsymbol{\nu}}_{t+1}} (\sqrt{\beta_2 \widetilde{\boldsymbol{\nu}}_t} + \sqrt{\widetilde{\boldsymbol{\nu}}_{t+1}})} \ge \frac{\|\boldsymbol{G}_t\|^2 \beta_2 (1-\beta_2) \|\boldsymbol{g}_t\|^2}{\sqrt{\beta_2 \widetilde{\boldsymbol{\nu}}_t} \sqrt{\widetilde{\boldsymbol{\nu}}_{t+1}} (\sqrt{\beta_2 \widetilde{\boldsymbol{\nu}}_t} + \sqrt{\widetilde{\boldsymbol{\nu}}_{t+1}})} \\ &\ge \frac{1}{4} \frac{\|\boldsymbol{G}_t\|^2 (1-\beta_2) \|\boldsymbol{g}_t\|^2}{\sqrt{\widetilde{\boldsymbol{\nu}}_t} (\sqrt{\boldsymbol{\nu}_t} + \sqrt{\widetilde{\boldsymbol{\nu}}_t})^2}, \end{split}$$

where in the last inequality, we use $\sqrt{\beta_2} \ge \frac{1}{2}$. Applying the above inequality back to Eq. (9), we obtain that

$$\frac{\eta}{1-\beta_{1}}\mathbb{E}^{|\mathcal{F}_{t}}\left[\left\|\boldsymbol{G}_{t}\right\|\left(\frac{(1-\beta_{2})\boldsymbol{g}_{t}^{2}}{\sqrt{\boldsymbol{\nu}_{t}}\sqrt{\boldsymbol{\tilde{\nu}_{t}}}(\sqrt{\boldsymbol{\nu}_{t}}+\sqrt{\boldsymbol{\tilde{\nu}_{t}}})\right)\left\|\boldsymbol{m}_{t}\right\|\right] \\
\leq \frac{1}{4}\eta\frac{\|\boldsymbol{G}_{t}\|^{2}}{\sqrt{\boldsymbol{\tilde{\nu}_{t}}}} + \frac{4\eta\sqrt{1-\beta_{2}}\sigma_{0}}{\left(1-\frac{\beta_{1}^{2}}{\beta_{2}}\right)^{2}}\mathbb{E}^{|\mathcal{F}_{t}}\frac{\|\boldsymbol{g}_{t}\|^{2}}{\boldsymbol{\nu}_{t}} + \eta\frac{64}{(1-\beta_{1})}\sigma_{1}^{2}\mathbb{E}^{|\mathcal{F}_{t}}\left(\frac{1}{\sqrt{\beta_{2}}\boldsymbol{\tilde{\nu}_{t}}}-\frac{1}{\sqrt{\boldsymbol{\tilde{\nu}_{t+1}}}}\right)\left\|\boldsymbol{G}_{t}\right\|^{2}.$$
(10)

 Furthermore, due to Assumption 1, we have (we define $G_0 \triangleq G_1$)

$$\begin{aligned} \|\boldsymbol{G}_{t+1}\|^{2} &\leq \|\boldsymbol{G}_{t}\|^{2} + 2\|\boldsymbol{G}_{t}\|\|\boldsymbol{G}_{t+1} - \boldsymbol{G}_{t}\| + \|\boldsymbol{G}_{t+1} - \boldsymbol{G}_{t}\|^{2} \\ &\leq \|\boldsymbol{G}_{t}\|^{2} + 2(L_{0} + L_{1}\|\boldsymbol{G}_{t}\|)\|\boldsymbol{G}_{t}\|\|\boldsymbol{w}_{t+1} - \boldsymbol{w}_{t}\| + 2(L_{0}^{2} + L_{1}^{2}\|\boldsymbol{G}_{t}\|^{2})\|\boldsymbol{w}_{t+1} - \boldsymbol{w}_{t}\|^{2}, \end{aligned}$$

which by
$$\begin{aligned} &\frac{\eta}{\sqrt{1-\beta_2}} = \frac{\sqrt{1-\frac{\beta_1^2}{\beta_2}(1-\frac{\beta_1}{\sqrt{\beta_2}})^2}}{1024\sigma_1^2(L_1+L_0)(1-\beta_1)} \text{ further leads to} \\ &\frac{1}{\sqrt{\beta_2 \widetilde{\boldsymbol{\nu}}_{t+1}}} \|\boldsymbol{G}_t\|^2 \\ \geq &\frac{1}{\sqrt{\beta_2 \widetilde{\boldsymbol{\nu}}_{t+1}}} \left(\|\boldsymbol{G}_{t+1}\|^2 - 2(L_0 + L_1\|\boldsymbol{G}_t\|)\|\boldsymbol{G}_t\|\|\boldsymbol{w}_{t+1} - \boldsymbol{w}_t\| - 2(L_0^2 + L_1^2\|\boldsymbol{G}_t\|^2)\|\boldsymbol{w}_{t+1} - \boldsymbol{w}_t\|^2\right) \\ \geq &\left(\frac{1}{\sqrt{\beta_2 \widetilde{\boldsymbol{\nu}}_{t+1}}} \|\boldsymbol{G}_{t+1}\|^2 - \frac{2L_0}{\sigma_0}\frac{(1-\beta_1)}{64\sigma_1^2}\|\boldsymbol{w}_{t+1} - \boldsymbol{w}_t\|^2 - \frac{3}{8}\frac{(1-\beta_1)}{64\sigma_1^2}\frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\widetilde{\boldsymbol{\nu}}_t}}\right). \end{aligned}$$

Applying the above inequality back to Eq. (10) leads to that

$$\frac{\eta}{1-\frac{\beta_1}{\sqrt{\beta_2}}} \mathbb{E}^{|\mathcal{F}_t} \left[\|\boldsymbol{G}_t\| \left(\frac{(1-\beta_2)\boldsymbol{g}_t^2}{\sqrt{\boldsymbol{\nu}_t}\sqrt{\boldsymbol{\tilde{\nu}}_t}(\sqrt{\boldsymbol{\nu}_t}+\sqrt{\boldsymbol{\tilde{\nu}}_t})} \right) \|\boldsymbol{m}_t\| \right] \\ \leq \frac{5}{\eta} \frac{\|\boldsymbol{G}_t\|^2}{2\eta} + \frac{4\eta\sqrt{1-\beta_2}\sigma_0}{\eta} \mathbb{E}^{|\mathcal{F}_t} \frac{\|\boldsymbol{g}_t\|^2}{2\eta} + \eta \frac{64}{\eta} \sigma_t^2 \mathbb{E}^{|\mathcal{F}_t} \left(\frac{\|\boldsymbol{G}_t\|^2}{2\eta} + \eta \frac{64}{\eta} \right) \mathbb{E}^{|\mathcal{F}_t|} \left(\frac{\|\boldsymbol{G}_t\|^2}{2\eta} + \eta \frac{64}{\eta} \right) \mathbb{E}^$$

$$\leq \frac{5}{8} \eta \frac{\|\boldsymbol{G}_{t}\|^{2}}{\sqrt{\tilde{\boldsymbol{\nu}}_{t}}} + \frac{4\eta\sqrt{1-\beta_{2}}\sigma_{0}}{(1-\beta_{1})^{2}} \mathbb{E}^{|\mathcal{F}_{t}} \frac{\|\boldsymbol{g}_{t}\|^{2}}{\boldsymbol{\nu}_{t}} + \eta \frac{64}{(1-\beta_{1})}\sigma_{1}^{2} \mathbb{E}^{|\mathcal{F}_{t}} \left(\frac{\|\boldsymbol{G}_{t}\|^{2}}{\sqrt{\beta_{2}}\tilde{\boldsymbol{\nu}}_{t}} - \frac{\|\boldsymbol{G}_{t+1}\|^{2}}{\sqrt{\tilde{\boldsymbol{\nu}}_{t+1}}}\right) + 2\frac{L_{0}}{\sigma_{0}} \mathbb{E}^{|\mathcal{F}_{t}} \|\boldsymbol{w}_{t+1} - \boldsymbol{w}_{t}\|^{2}.$$

$$(11)$$

As for the second term in the right-hand-side of Eq. (8), we have

$$= \frac{1 - \frac{\beta_1}{\sqrt{\beta_2}}}{\leq \frac{1}{8}\eta \frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\boldsymbol{\tilde{\nu}}_t}} + \frac{8\eta\sqrt{1 - \beta_2}\sigma_0}{(1 - \beta_1)^2} \mathbb{E}^{|\mathcal{F}_t} \left[\left(\frac{\|\boldsymbol{m}_t\|^2}{\boldsymbol{\nu}_t} \right) \right].$$
(12)

1261 In the last inequality we use again $\beta_2 \ge \beta_1$. With Inequalities (11) and (12), we conclude that the 1262 first-order term can be bounded by

$$\begin{aligned}
\mathbf{E}^{|\mathcal{F}_{t}} & \left[\langle \nabla f(\boldsymbol{w}_{t}), \boldsymbol{u}_{t+1} - \boldsymbol{u}_{t} \rangle \right] \leq -\frac{1}{4} \eta \mathbb{E} \frac{\|\boldsymbol{G}_{t}\|^{2}}{\sqrt{\tilde{\boldsymbol{\nu}}_{t}}} + \frac{4\eta \sqrt{1 - \beta_{2}} \sigma_{0}}{\left(1 - \beta_{1}\right)^{2}} \mathbb{E}^{|\mathcal{F}_{t}} \frac{\|\boldsymbol{g}_{t}\|^{2}}{\boldsymbol{\nu}_{t}} + \eta \frac{64}{\left(1 - \beta_{1}\right)} \sigma_{1}^{2} \mathbb{E}^{|\mathcal{F}_{t}} \left(\frac{\|\boldsymbol{G}_{t}\|^{2}}{\sqrt{\beta_{2}} \tilde{\boldsymbol{\nu}}_{t}} - \frac{\|\boldsymbol{G}_{t+1}\|^{2}}{\sqrt{\tilde{\boldsymbol{\nu}}_{t+1}}} \right) \\ + 2 \frac{L_{0}}{\sigma_{0}} \mathbb{E}^{|\mathcal{F}_{t}} \|\boldsymbol{w}_{t+1} - \boldsymbol{w}_{t}\|^{2} + \frac{8\eta \sqrt{1 - \beta_{2}} \sigma_{0}}{\left(1 - \beta_{1}\right)^{2}} \mathbb{E}^{|\mathcal{F}_{t}} \left[\left(\frac{\|\boldsymbol{m}_{t}\|^{2}}{\boldsymbol{\nu}_{t}} \right) \right]. \end{aligned}$$

$$(13)$$

1270 Analysis for the second-order term. To recall, the second order term is $\frac{1}{2}(L_0 + L_1 \|\nabla f(\boldsymbol{w}_t)\|)(\|\boldsymbol{u}_{t+1} - \boldsymbol{w}_t\| + \|\boldsymbol{u}_t - \boldsymbol{w}_t\|)\|\boldsymbol{u}_{t+1} - \boldsymbol{u}_t\|$. Before we start, we have the following expansion for $\boldsymbol{u}_{t+1} - \boldsymbol{u}_t$: (here the operations are all coordinate-wisely)

$$\boldsymbol{u}_{t+1} - \boldsymbol{u}_{t} = \frac{\boldsymbol{w}_{t+1} - \boldsymbol{w}_{t} - \frac{\beta_{1}}{\sqrt{\beta_{2}}}(\boldsymbol{w}_{t} - \boldsymbol{w}_{t-1})}{1 - \frac{\beta_{1}}{\sqrt{\beta_{2}}}} \\ = \frac{-\eta \frac{\boldsymbol{m}_{t}}{\sqrt{\nu_{t}}} + \eta \frac{\beta_{1}}{\sqrt{\beta_{2}}} \frac{\boldsymbol{m}_{t-1}}{\sqrt{\nu_{t-1}}}}{1 - \frac{\beta_{1}}{\sqrt{\beta_{2}}}} = \frac{-\eta \frac{\boldsymbol{m}_{t}}{\sqrt{\nu_{t}}} + \eta \beta_{1} \frac{\boldsymbol{m}_{t-1}}{\sqrt{\nu_{t}}} - \eta \beta_{1} \frac{\boldsymbol{m}_{t-1}}{\sqrt{\nu_{t}}} + \eta \frac{\beta_{1}}{\sqrt{\beta_{2}}} \frac{\boldsymbol{m}_{t-1}}{\sqrt{\nu_{t-1}}}}{1 - \frac{\beta_{1}}{\sqrt{\beta_{2}}}} \\ = \frac{-\eta \frac{(1 - \beta_{1})\boldsymbol{g}_{t}}{\sqrt{\nu_{t}}} + \eta \frac{\beta_{1}(1 - \beta_{2})\|\boldsymbol{g}_{t}\|^{2}}{\sqrt{\beta_{2}}} \frac{\boldsymbol{m}_{t-1}}{\sqrt{\nu_{t-1}}\sqrt{\nu_{t}}(\sqrt{\nu_{t}} + \sqrt{\beta_{2}\nu_{t-1}})}}{1 - \frac{\beta_{1}}{\sqrt{\beta_{2}}}}$$
(14)

Then firstly, we have

$$\begin{aligned}
1284 & \frac{1}{2}L_{0}(\|\boldsymbol{u}_{t+1} - \boldsymbol{w}_{t}\| + \|\boldsymbol{u}_{t} - \boldsymbol{w}_{t}\|)\|\boldsymbol{u}_{t+1} - \boldsymbol{u}_{t}\| \\
1285 & \leq \frac{1}{2}L_{0}\left(\|\boldsymbol{u}_{t+1} - \boldsymbol{w}_{t}\| + \|\boldsymbol{u}_{t} - \boldsymbol{w}_{t}\|^{2} + \frac{1}{2}\|\boldsymbol{u}_{t} - \boldsymbol{w}_{t}\|^{2}\right) \\
1287 & = \frac{1}{2}L_{0}\left(\left\|\frac{-\eta\frac{(1-\beta_{1})\boldsymbol{g}_{t}}{\sqrt{\boldsymbol{\nu}_{t}}} + \eta\frac{\beta_{1}(1-\beta_{2})\|\boldsymbol{g}_{t}\|^{2}}{\sqrt{\beta_{2}}}\frac{\boldsymbol{m}_{t-1}}{\sqrt{\boldsymbol{\nu}_{t-1}}\sqrt{\boldsymbol{\nu}_{t}}(\sqrt{\boldsymbol{\nu}_{t}} + \sqrt{\beta_{2}}\boldsymbol{\nu}_{t-1})}\right\|^{2} + \frac{1}{2}\left\|\frac{\beta_{1}}{\sqrt{\beta_{2}}}\left(\boldsymbol{w}_{t} - \boldsymbol{w}_{t-1}\right)\right\|^{2} + \frac{1}{2}\left\|\frac{1}{1 - \frac{\beta_{1}}{\sqrt{\beta_{2}}}}\left(\boldsymbol{w}_{t+1} - \boldsymbol{w}_{t}\right)\right\|^{2}\right) \\
1289 & = \frac{1}{2}L_{0}\left(\left\|\frac{-\eta\frac{(1-\beta_{1})\boldsymbol{g}_{t}}{\sqrt{\boldsymbol{\nu}_{t}}} + \eta\frac{\beta_{1}(1-\beta_{1})}{\sqrt{\beta_{2}}}\right)^{2}\left\|\boldsymbol{g}_{t}\right\|^{2} + \frac{1}{2}\left(\frac{\beta_{1}}{\sqrt{\beta_{2}}}\right)^{2}\left\|\boldsymbol{w}_{t} - \boldsymbol{w}_{t-1}\right\|^{2} + \frac{1}{2}\left(\frac{1}{1 - \frac{\beta_{1}}{\sqrt{\beta_{2}}}}\right)^{2}\left\|\boldsymbol{w}_{t+1} - \boldsymbol{w}_{t}\right\|^{2}\right) \\
1290 & \leq \frac{L_{0}\eta^{2}}{2}\left(\left(\frac{1-\beta_{1}}{1 - \frac{\beta_{1}}{\sqrt{\beta_{2}}}} + \frac{\beta_{1}(1-\beta_{1})}{(\sqrt{\beta_{2}} - \beta_{1})\sqrt{1 - \frac{\beta_{1}^{2}}{\beta_{2}}}}\right)^{2}\left\|\boldsymbol{g}_{t}\right\|^{2} + \frac{1}{2}\left(\frac{\beta_{1}}{\sqrt{\beta_{2}}}\right)^{2}\left\|\boldsymbol{m}_{t-1}\right\|^{2} + \frac{1}{2}\left(\frac{1}{1 - \frac{\beta_{1}}{\sqrt{\beta_{2}}}}\right)^{2}\left\|\boldsymbol{w}_{t}\right\|^{2}\right) \\
1293 & (\bullet) \frac{L_{0}\eta^{2}}{2}\left(2\left(\frac{1-\beta_{1}}{1 - \frac{\beta_{1}}{\sqrt{\beta_{2}}}} + \frac{\beta_{1}(1-\beta_{1})}{(\sqrt{\beta_{2}} - \beta_{1})\sqrt{1 - \frac{\beta_{1}^{2}}{\beta_{2}^{2}}}}\right)^{2}\left\|\boldsymbol{g}_{t}\right\|^{2} + \left(\frac{\beta_{1}}{\sqrt{\beta_{2}}}\right)^{2}\left\|\boldsymbol{m}_{t-1}\right\|^{2}\right). \\
1294 & (\bullet) \frac{L_{0}\eta^{2}}{2}\left(2\left(\frac{1-\beta_{1}}{1 - \frac{\beta_{1}}{\sqrt{\beta_{2}}}} + \frac{\beta_{1}(1-\beta_{1})}{(\sqrt{\beta_{2}} - \beta_{1})\sqrt{1 - \frac{\beta_{1}^{2}}{\beta_{2}^{2}}}}\right)^{2}\left\|\boldsymbol{g}_{t}\right\|^{2} + \left(\frac{\beta_{1}}{\sqrt{\beta_{2}}}\right)^{2}\left\|\boldsymbol{m}_{t-1}\right\|^{2}\right). \\
1295 & (\bullet) \frac{L_{0}\eta^{2}}{2}\left(2\left(\frac{1-\beta_{1}}{1 - \frac{\beta_{1}}{\sqrt{\beta_{2}}}} + \frac{\beta_{1}(1-\beta_{1})}{(\sqrt{\beta_{2}} - \beta_{1})\sqrt{1 - \frac{\beta_{1}^{2}}{\beta_{2}^{2}}}}\right)^{2}\left\|\boldsymbol{g}_{t}\right\|^{2} + \left(\frac{\beta_{1}}{\sqrt{\beta_{2}}}\right)^{2}\left\|\boldsymbol{m}_{t-1}\right\|^{2}\right). \\
1295 & (\bullet) \frac{L_{0}\eta^{2}}{2}\left(\frac{1-\beta_{1}}{1 - \frac{\beta_{1}}{\sqrt{\beta_{2}}}} + \frac{\beta_{1}(1-\beta_{1})}{(\sqrt{\beta_{2}} - \beta_{1})\sqrt{1 - \frac{\beta_{1}^{2}}{\beta_{2}^{2}}}}\right)^{2}\left\|\boldsymbol{g}_{t}\right\|^{2} + \frac{\beta_{1}}{\sqrt{\beta_{2}}}\left(\frac{\beta_{1}}{\sqrt{\boldsymbol{\mu}_{1}}}\right)^{2}\left\|\boldsymbol{g}_{t}\right\|^{2} + \frac{\beta_{1}}{\sqrt{\beta_{1}}}\left(\frac{\beta_{1}}{\sqrt{\beta_{2}}}\right)^{2}\left\|\boldsymbol{g}_{t}\right\|^{2} + \frac{\beta_{1}}{\sqrt{\beta_{2}}}\left(\frac{\beta_{1}}{\sqrt{\beta_{2}}}\right)^{2}\left\|\boldsymbol{g}_{t}\right\|^{$$

Secondly, we have $\frac{1}{2}L_1 \|\nabla f(\boldsymbol{w}_t)\| (\|\boldsymbol{u}_{t+1} - \boldsymbol{w}_t\| + \|\boldsymbol{u}_t - \boldsymbol{w}_t\|) \|\boldsymbol{u}_{t+1} - \boldsymbol{u}_t\|$ $\leq \frac{1}{2}L_1 \|\nabla f(\boldsymbol{w}_t)\| (2\|\boldsymbol{u}_{t+1} - \boldsymbol{w}_t\| + \|\boldsymbol{u}_{t+1} - \boldsymbol{u}_t\|) \left(\frac{\left\| \eta \frac{(1-\beta_1)\boldsymbol{g}_t}{\sqrt{\boldsymbol{\nu}_t}} \right\|}{1 - \frac{\beta_1}{\sqrt{\beta_2}}} + \frac{\eta \frac{\beta_1(1-\beta_2)\|\boldsymbol{g}_t\|^2}{\sqrt{\beta_2}} \frac{\|\boldsymbol{m}_{t-1}\|}{\sqrt{\boldsymbol{\nu}_{t-1}}\sqrt{\boldsymbol{\nu}_t}(\sqrt{\boldsymbol{\nu}_t} + \sqrt{\beta_2\boldsymbol{\nu}_{t-1}})}{1 - \frac{\beta_1}{\sqrt{\beta_2}}} \right)$ $\overset{(*)}{\leq} \frac{1}{2} L_1 \|\nabla f(\boldsymbol{w}_t)\| (2\|\boldsymbol{u}_{t+1} - \boldsymbol{w}_t\| + \|\boldsymbol{u}_{t+1} - \boldsymbol{u}_t\|) \left(\frac{\left\| \eta \frac{(1-\beta_1)\boldsymbol{g}_t}{\sqrt{\nu_t}} \right\|}{1 - \frac{\beta_1}{\sqrt{\beta_2}}} + \frac{\eta \frac{\beta_1(1-\beta_1)}{\sqrt{\beta_2}} \frac{\|\boldsymbol{g}_t\|}{\sqrt{\nu_t}}}{(1 - \frac{\beta_1}{\sqrt{\beta_2}})\sqrt{1 - \frac{\beta_1^2}{2}}} \right)$ $=\frac{L_1}{2}\eta\left(\frac{1-\beta_1}{1-\frac{\beta_1}{\sqrt{\beta_2}}}+\frac{\beta_1(1-\beta_1)}{(\sqrt{\beta_2}-\beta_1)\sqrt{1-\frac{\beta_1^2}{2}}}\right)\|\nabla f(\boldsymbol{w}_t)\|(2\|\boldsymbol{u}_{t+1}-\boldsymbol{w}_t\|+\|\boldsymbol{u}_t-\boldsymbol{u}_{t+1}\|)\frac{\|\boldsymbol{g}_t\|}{\sqrt{\boldsymbol{\nu}_t}}$ $\stackrel{(\circ)}{=} \frac{L_1}{2} \eta \left(\frac{1 - \beta_1}{1 - \frac{\beta_1}{\sqrt{\beta_2}}} + \frac{\beta_1 (1 - \beta_1)}{(\sqrt{\beta_2} - \beta_1)\sqrt{1 - \frac{\beta_1}{2}}} \right) \|\boldsymbol{G}_t\| \left(\|\boldsymbol{u}_{t+1} - \boldsymbol{u}_t\| + 2\frac{1}{1 - \frac{\beta_1}{\sqrt{\beta_2}}} \eta \left\| \frac{\boldsymbol{m}_t}{\sqrt{\boldsymbol{\nu}_t}} \right\| \right) \frac{\|\boldsymbol{g}_t\|}{\sqrt{\boldsymbol{\nu}_t}}$ where inequality (*) is due to that $\frac{\|\boldsymbol{m}_{t-1}\|}{\sqrt{\boldsymbol{\nu}_{t-1}}} \leq \frac{1-\beta_1}{\sqrt{1-\beta_2}\sqrt{1-\frac{\beta_1^2}{\sigma}}}, \frac{\|\boldsymbol{g}_t\|}{\sqrt{\boldsymbol{\nu}_t}} \leq \frac{1}{\sqrt{1-\beta_2}},$ and equation (\circ) is due to $u_t - w_t = \frac{\frac{\beta_1}{\sqrt{\beta_2}}}{1 - \frac{\beta_1}{\sqrt{\beta_2}}} (w_t - w_{t-1})$ and $u_{t+1} - w_t = \frac{1}{1 - \frac{\beta_1}{\sqrt{\beta_2}}} (w_{t+1} - w_t)$. As for the term $\|G_t\| \frac{\|m_t\|}{\sqrt{\nu_t}} \frac{\|g_t\|}{\sqrt{\nu_t}}$, we first add additional denominator for it. Specifically, we have $\|\boldsymbol{G}_t\| \frac{\|\boldsymbol{m}_t\|}{\sqrt{\boldsymbol{\nu}_t}} \frac{\|\boldsymbol{g}_t\|}{\sqrt{\boldsymbol{\nu}_t}} = \frac{\|\boldsymbol{G}_t\| \|\boldsymbol{m}_t\| \|\boldsymbol{g}_t\|}{\boldsymbol{\nu}_t + (1-\beta_2)\sigma_0^2} + \frac{\|\boldsymbol{G}_t\| \|\boldsymbol{m}_t\| \|\boldsymbol{g}_t\| (1-\beta_2)\sigma_0^2}{(\boldsymbol{\nu}_t + (1-\beta_2)\sigma_0^2)\boldsymbol{\nu}_t}$ $\leq \frac{\|\boldsymbol{G}_t\|\|\boldsymbol{m}_t\|\|\boldsymbol{g}_t\|}{\boldsymbol{\nu}_t + (1-\beta_2)\sigma_0^2} + \frac{\|\boldsymbol{G}_t\|\|\boldsymbol{m}_t\|\sigma_0}{\sqrt{\boldsymbol{\nu}_t + (1-\beta_2)\sigma_0^2}\sqrt{\boldsymbol{\nu}_t}}$ $\leq \frac{\|\boldsymbol{G}_t\|\|\boldsymbol{m}_t\|\|\boldsymbol{g}_t\|}{\boldsymbol{\nu}_t + (1-\beta_2)\sigma_0^2} + \frac{1}{2}\frac{\|\boldsymbol{G}_t\|^2\sigma_0}{\boldsymbol{\nu}_t + (1-\beta_2)\sigma_0^2} + \frac{1}{2}\sigma_0\frac{\|\boldsymbol{m}_t\|^2}{\boldsymbol{\nu}_t}$ $\leq \frac{\|\boldsymbol{G}_t\|\|\boldsymbol{m}_t\|\|\boldsymbol{g}_t\|}{\boldsymbol{\nu}_t + (1-\beta_2)\sigma_0^2} + \frac{1}{2\sqrt{1-\beta_2}}\frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\boldsymbol{\nu}_t + (1-\beta_2)\sigma_0^2}} + \frac{1}{2}\sigma_0\frac{\|\boldsymbol{m}_t\|^2}{\boldsymbol{\nu}_t}.$ We analyze the first term in the right-hand-side of above inequality more carefully. Specifically, this term with expectation can be bounded as $\mathbb{E}^{|\mathcal{F}_t} \frac{\|\boldsymbol{G}_t\| \|\boldsymbol{m}_t\| \|\boldsymbol{g}_t\|}{\boldsymbol{\nu}_t + (1-\beta_2)\sigma_0^2}$ $\leq \mathbb{E}^{|\mathcal{F}_{t}} \frac{\|\boldsymbol{G}_{t}\| \|\boldsymbol{m}_{t}\| \|\boldsymbol{g}_{t}\|}{\sqrt{\boldsymbol{\nu}_{t} + (1 - \beta_{2})\sigma_{0}^{2}}\sqrt{\beta_{2}\boldsymbol{\nu}_{t-1} + (1 - \beta_{2})\sigma_{0}^{2}}}$ $\leq \frac{\|\boldsymbol{G}_t\|}{\sqrt{\beta_2 \boldsymbol{\nu}_{t-1} + (1-\beta_2)\sigma_2^2}} \sqrt{\|\boldsymbol{g}_t\|^2} \sqrt{\mathbb{E}^{|\mathcal{F}_t|} \frac{\|\boldsymbol{m}_t\|^2}{\boldsymbol{\nu}_t + (1-\beta_2)\sigma_0^2}}$ $\stackrel{(\star)}{\leq} \frac{\|\boldsymbol{G}_t\|}{\sqrt{\beta_2 \boldsymbol{\nu}_{t-1} + (1-\beta_2)\sigma_0^2}} \sqrt{\sigma_1^2 \|\boldsymbol{G}_t\|^2 + \sigma_0^2} \sqrt{\mathbb{E}^{|\mathcal{F}_t} \frac{\|\boldsymbol{m}_t\|^2}{\boldsymbol{\nu}_t + (1-\beta_2)\sigma_0^2}}$ $\leq \frac{\|\boldsymbol{G}_t\|}{\sqrt{\beta_2 \boldsymbol{\nu}_{t-1} + (1-\beta_2)\sigma_0^2}} (\sigma_1 \|\boldsymbol{G}_t\| + \sigma_0) \sqrt{\mathbb{E}^{|\mathcal{F}_t|} \frac{\|\boldsymbol{m}_t\|^2}{\boldsymbol{\nu}_t + (1-\beta_2)\sigma_0^2}}$ $\leq \frac{1-\beta_1}{\sqrt{1-\beta_2}\sqrt{1-\frac{\beta_1^2}{\beta_2}}}\sigma_1 \frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\beta_2\boldsymbol{\nu}_{t-1}} + (1-\beta_2)\sigma_0^2} + \frac{1}{2\sqrt{1-\beta_2}}\frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\beta_2\boldsymbol{\nu}_{t-1}} + (1-\beta_2)\sigma_0^2} + \frac{\sigma_0}{2}\mathbb{E}^{|\mathcal{F}_t}\frac{\|\boldsymbol{m}_t\|^2}{\boldsymbol{\nu}_t + (1-\beta_2)\sigma_0^2}$

where Eq. (\star) is due to Holder's inequality.

Meanwhile, due to Eq. (14), we have that the term $\|G_t\| \|u_{t+1} - u_t\| \frac{\|g_t\|}{\sqrt{\nu_t}}$ can be be bounded as

$$\|\boldsymbol{G}_t\|\|\boldsymbol{u}_{t+1} - \boldsymbol{u}_t\|\frac{\|\boldsymbol{g}_t\|}{\sqrt{\boldsymbol{\nu}_t}} \le \eta \left(\frac{1 - \beta_1}{1 - \frac{\beta_1}{\sqrt{\beta_2}}} + \frac{\beta_1(1 - \beta_1)}{(\sqrt{\beta_2} - \beta_1)\sqrt{1 - \frac{\beta_1^2}{\beta_2}}}\right)\|\boldsymbol{G}_t\|\frac{\|\boldsymbol{g}_t\|}{\sqrt{\boldsymbol{\nu}_t}}\frac{\|\boldsymbol{g}_t\|}{\sqrt{\boldsymbol{\nu}_t}} \frac{\|\boldsymbol{g}_t\|}{\sqrt{\boldsymbol{\nu}_t}}.$$

Then, following the similar reasoning above, we have $|G_t|| ||u_{t+1} - u_t|| \frac{||g_t||}{\sqrt{\nu_t}}$ can be bounded as

$$\begin{split} & \mathbb{E}^{|\mathcal{F}_{t}} \| \boldsymbol{G}_{t} \| \frac{\| \boldsymbol{g}_{t} \|}{\sqrt{\nu_{t}}} \frac{\| \boldsymbol{g}_{t} \|}{\sqrt{\nu_{t}}} \\ & \mathbb{E}^{|\mathcal{F}_{t}} \| \boldsymbol{G}_{t} \| \frac{\| \boldsymbol{g}_{t} \|}{\sqrt{\nu_{t}}} \frac{\| \boldsymbol{g}_{t} \|}{\sqrt{\nu_{t}}} \\ & \mathbb{E}^{|\mathcal{F}_{t}} \| \boldsymbol{G}_{t} \|^{2} \\ & \mathbb{E}^{|\mathcal{F}_{t}} \| \boldsymbol{G}_{t} \|^{2} \\ & \mathbb{E}^{|\mathcal{F}_{t}} \frac{\| \boldsymbol{G}_{t} \|^{2}}{\sqrt{\beta_{2} \nu_{t-1} + (1 - \beta_{2}) \sigma_{0}^{2}}} + \frac{1}{2\sqrt{1 - \beta_{2}}} \frac{\| \boldsymbol{G}_{t} \|^{2}}{\sqrt{\beta_{2} \nu_{t-1} + (1 - \beta_{2}) \sigma_{0}^{2}}} + \sigma_{0} \mathbb{E}^{|\mathcal{F}_{t}} \frac{\| \boldsymbol{g}_{t} \|^{2}}{\nu_{t} + (1 - \beta_{2}) \sigma_{0}^{2}} \\ & \mathbb{E}^{|\mathcal{F}_{t}} \frac{\| \boldsymbol{G}_{t} \|^{2}}{\sqrt{\nu_{t} + (1 - \beta_{2}) \sigma_{0}^{2}}} + \frac{1}{2} \sigma_{0} \mathbb{E}^{|\mathcal{F}_{t}} \frac{\| \boldsymbol{g}_{t} \|^{2}}{\nu_{t}}. \end{split}$$

Putting all the estimations together, we have that the second-order term can be bounded by

> Here in the second inequality we use $\beta_2 \geq \beta_1$, and in the last inequality we use $\frac{\eta}{\sqrt{1-\beta_2}}$ $\frac{\sqrt{1\!-\!\frac{\beta_1^2}{\beta_2}}(1\!-\!\frac{\beta_1}{\sqrt{\beta_2}})^2}{1024\sigma_1^2(L_1\!+\!L_0)(1\!-\!\beta_1)}.$

Applying the estimations of the first-order term (Eq. (13)) and the second-order term (Eq. (15)) back into the descent lemma, we derive that

$$\mathbb{E}^{|\mathcal{F}_{t}} f(\boldsymbol{u}_{t+1}) \leq f(\boldsymbol{u}_{t}) - \frac{1}{8} \eta \frac{\|\boldsymbol{G}_{t}\|^{2}}{\sqrt{\tilde{\boldsymbol{\nu}}_{t}}} + \frac{4\eta\sqrt{1-\beta_{2}}\sigma_{0}}{(1-\beta_{1})^{2}} \mathbb{E}^{|\mathcal{F}_{t}} \frac{\|\boldsymbol{g}_{t}\|^{2}}{\boldsymbol{\nu}_{t}} + \eta \frac{64}{(1-\beta_{1})}\sigma_{1}^{2} \mathbb{E}^{|\mathcal{F}_{t}} \left(\frac{\|\boldsymbol{G}_{t}\|^{2}}{\sqrt{\beta_{2}}\tilde{\boldsymbol{\nu}}_{t}} - \frac{\|\boldsymbol{G}_{t+1}\|^{2}}{\sqrt{\tilde{\boldsymbol{\nu}}_{t+1}}}\right)$$

1400
1401
$$+ 2\frac{L_0}{\sigma_0} \mathbb{E}^{|\mathcal{F}_t|} \|\boldsymbol{w}_{t+1} - \boldsymbol{w}_t\|^2 + \frac{8\eta\sqrt{1-\beta_2}\sigma_0}{(1-\beta_1)^2} \mathbb{E}^{|\mathcal{F}_t|} \left[\left(\frac{\|\boldsymbol{m}_t\|^2}{\boldsymbol{\nu}_t} \right) \right]$$

1402
1403
$$+4\frac{L_1\eta^2\sigma_0}{(1-\beta_1)^{\frac{3}{2}}}\mathbb{E}^{|\mathcal{F}_t}\frac{\|\boldsymbol{g}_t\|^2}{\boldsymbol{\nu}_t}+2L_0\eta^2\left(8\frac{1}{1-\beta_1}\mathbb{E}^{|\mathcal{F}_t}\left\|\frac{\boldsymbol{g}_t}{\sqrt{\boldsymbol{\nu}_t}}\right\|^2+\left(\frac{1}{1-\beta_1}\right)^2\left\|\frac{\boldsymbol{m}_{t-1}}{\sqrt{\boldsymbol{\nu}_{t-1}}}\right\|^2\right).$$

Taking expectation to the above inequality and summing it over $t \in [1, T]$ then gives

Since $\beta_2 \geq \frac{1}{2}$ and $1 - \beta_2 \leq \frac{1 - \beta_1}{1024\sigma_1^2}$, we have

$$\eta \frac{64}{(1-\beta_1)} \sigma_1^2 \left(\frac{1}{\sqrt{\beta_2}} - 1\right) \sum_{t=1}^T \mathbb{E} \frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\widetilde{\boldsymbol{\nu}}_t}} \le \frac{1}{16} \eta \sum_{t=1}^T \mathbb{E} \frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\widetilde{\boldsymbol{\nu}}_t}}.$$

By further applying Lemma 3 and $\beta_2 \ge \beta_1$, we obtain

Here last inequality we apply $\frac{\eta}{\sqrt{1-\beta_2}} = \frac{\sqrt{1-\frac{\beta_1}{\beta_2}}(1-\frac{\beta_1}{\sqrt{\beta_2}})^2}{1024\sigma_1^2(L_1+L_0)(1-\beta_1)}.$

Below we transfer the above bound to the bound of $\sum_{t=1}^{T} \|G_t\|$ by two rounds of divide-and-conquer. In the first round, we will bound $\mathbb{E} \ln \nu_T$. To start with, we have that

$$\frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\widetilde{\boldsymbol{\nu}}_t}}\mathbb{1}_{\|G_t\|\geq\frac{\sigma_0}{\sigma_1}}\geq\frac{\frac{1}{2\sigma_1^2}\mathbb{E}^{|\mathcal{F}_t}\|\boldsymbol{g}_t\|^2}{\sqrt{\widetilde{\boldsymbol{\nu}}_t}}\mathbb{1}_{\|G_t\|\geq\frac{\sigma_0}{\sigma_1}}\\-\frac{\frac{1}{2\sigma_1^2}\mathbb{E}^{|\mathcal{F}_t}\|\boldsymbol{g}_t\|^2}{\mathbb{I}_{\|G_t\|\geq\frac{\sigma_0}{\sigma_1}}}=\mathbb{I}$$

$$= \frac{1}{\sqrt{\beta_2 \boldsymbol{\nu}_{t-1} + (1-\beta_2)\sigma_0^2}} \mathbb{I}_{\|G_t\| \ge \frac{\sigma_0}{\sigma_1}}$$

$$1 - \tau - \beta_1^2$$

 $\geq \frac{1}{2\sigma_1^2} \mathbb{E}^{|\mathcal{F}_t} \frac{\beta_2^{T-t} \| \boldsymbol{g}_t \|^2}{\sqrt{\boldsymbol{\nu}_T + (1-\beta_2)\sigma_0^2}} \mathbb{1}_{\|G_t\| \geq \frac{\sigma_0}{\sigma_1}},$

where the last inequality is due to that

$$\beta_2 \boldsymbol{\nu}_{t-1} + (1-\beta_2)\sigma_0^2 \le \beta_2^{t-T} \boldsymbol{\nu}_T + (1-\beta_2)\sigma_0^2 \le (\boldsymbol{\nu}_T + (1-\beta_2)\sigma_0^2)\beta_2^{2(t-T)}.$$
 (18)

Furthermore, we have $\frac{\sigma_0^2 + \frac{\beta_2^T \nu_0}{1 - \beta_2}}{\sqrt{\nu_T + (1 - \beta_2)\sigma_0^2}} + \sum_{t=1}^T \mathbb{E} \frac{\beta_2^{T-t} \|\boldsymbol{g}_t\|^2}{\sqrt{\nu_T + (1 - \beta_2)\sigma_0^2}} \mathbb{1}_{\|\boldsymbol{G}_t\| < \frac{\sigma_0}{\sigma_1}}$ $\leq \frac{\sigma_0^2 + \frac{\beta_2^T \nu_0}{1-\beta_2}}{\sqrt{\nu_0 \beta_2^T + \sum_{s=1}^T \beta_2^{T-s} \|g_s\|^2 \mathbb{1}_{\|G_s\| < \frac{\sigma_0}{\sigma_1}} + (1-\beta_2)\sigma_0^2}} + \sum_{t=1}^T \mathbb{E} \frac{\beta_2^{T-s} \|g_t\|^2}{\sqrt{\nu_0 \beta_2^T + \sum_{s=1}^T \beta_2^{T-s} \|g_s\|^2 \mathbb{1}_{\|G_s\| < \frac{\sigma_0}{\sigma_1}} + (1-\beta_2)\sigma_0^2}} \mathbb{1}_{\|G_t\| < \frac{\sigma_0}{\sigma_1}} \mathbb{1}_{\|G_s\| < \frac{\sigma_0}{\sigma_1}} + (1-\beta_2)\sigma_0^2} = \frac{\sigma_0^T (1-\beta_2)\sigma_0^2}{\sqrt{\nu_0 \beta_2^T + \sum_{s=1}^T \beta_2^{T-s} \|g_s\|^2 \mathbb{1}_{\|G_s\| < \frac{\sigma_0}{\sigma_1}} + (1-\beta_2)\sigma_0^2}} \mathbb{1}_{\|G_s\| < \frac{\sigma_0}{\sigma_1}} \mathbb{1}_{\|G_s\| < \frac{\sigma_0}{\sigma_1}} + (1-\beta_2)\sigma_0^2} = \frac{\sigma_0^T (1-\beta_2)\sigma_0^2}{\sqrt{\nu_0 \beta_2^T + \sum_{s=1}^T \beta_2^{T-s} \|g_s\|^2 \mathbb{1}_{\|G_s\| < \frac{\sigma_0}{\sigma_1}} + (1-\beta_2)\sigma_0^2}} \mathbb{1}_{\|G_s\| < \frac{\sigma_0}{\sigma_1}} \mathbb{1}_{\|G_s\| < \frac{\sigma_0}{\sigma_1}} + (1-\beta_2)\sigma_0^2} = \frac{\sigma_0^T (1-\beta_2)\sigma_0^2}{\sqrt{\nu_0 \beta_2^T + \sum_{s=1}^T \beta_2^{T-s} \|g_s\|^2 \mathbb{1}_{\|G_s\| < \frac{\sigma_0}{\sigma_1}} + (1-\beta_2)\sigma_0^2}} \mathbb{1}_{\|G_s\| < \frac{\sigma_0}{\sigma_1}} \mathbb{1}_{\|G_s\| < \frac{\sigma_0}{\sigma_1}} + (1-\beta_2)\sigma_0^2} \mathbb{1}_{\|G_s\| < \frac{\sigma_0}{\sigma_1}} \mathbb{1$ $=\frac{1}{1-\beta_2}\mathbb{E}\sqrt{\nu_0\beta_2^T+\sum_{s=1}^T\beta_2^{T-s}\|g_s\|^2\mathbb{1}_{\|\boldsymbol{G}_s\|<\frac{\sigma_0}{\sigma_1}}+(1-\beta_2)\sigma_0^2}\leq\frac{1}{1-\beta_2}\sqrt{\beta_2^T\nu_0+2\sigma_0^2}.$ (19) Conclusively, we obtain $\mathbb{E}\sqrt{\boldsymbol{\nu}_T + (1-\beta_2)\sigma_0^2}$ $=(1-\beta_2)\left(\frac{\sigma_0^2+\frac{\beta_2^2\,\boldsymbol{\nu}_0}{1-\beta_2}}{\sqrt{\boldsymbol{\nu}_T+(1-\beta_2)\sigma_0^2}}+\sum_{t=1}^T \mathbb{E}\frac{\beta_2^{T-t}\|\boldsymbol{g}_t\|^2}{\sqrt{\boldsymbol{\nu}_T+(1-\beta_2)\sigma_0^2}}\mathbb{1}_{\|\boldsymbol{G}_t\|<\frac{\sigma_0}{\sigma_1}}\right)$ + $\sum_{t=1}^{T} \mathbb{E} \frac{\beta_2^{T-t} \|\boldsymbol{g}_t\|^2}{\sqrt{\boldsymbol{\nu}_T + (1-\beta_2)\sigma_2^2}} \mathbb{1}_{\|\boldsymbol{G}_t\| \geq \frac{\sigma_0}{\sigma_1}} \right)$ $\leq \sqrt{\beta_2^T \boldsymbol{\nu}_0 + 2\sigma_0^2} + 2(1 - \beta_2)\sigma_1^2 \mathbb{E} \sum_{t=1}^T \frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\boldsymbol{\widetilde{\boldsymbol{\nu}}_t}}} \mathbb{1}_{\|\boldsymbol{G}_t\| \geq \frac{\sigma_0}{\sigma_1}}$ $\leq \sqrt{\beta_2^T \boldsymbol{\nu}_0 + 2\sigma_0^2} + 2(1-\beta_2)\sigma_1^2 \mathbb{E} \sum_{\boldsymbol{\nu}}^T \frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\boldsymbol{\mu}_t}}.$ Substituting $\mathbb{E} \sum_{t=1}^{T} \frac{\|G_t\|^2}{\sqrt{\mu_t}}$ according to Eq. (17), we obtain that $\mathbb{E}\sqrt{\boldsymbol{\nu}_T + (1-\beta_2)\sigma_0^2}$ $\leq \sqrt{\beta_2^T \boldsymbol{\nu}_0 + 2\sigma_0^2} + \frac{2(1-\beta_2)\sigma_1^2}{\eta} \mathbb{E} \sum_{t=1}^T \frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\boldsymbol{\widetilde{\nu}}_t}}$ $\leq \sqrt{\beta_2^T \boldsymbol{\nu}_0 + 2\sigma_0^2} + \frac{2(1-\beta_2)\sigma_1^2}{n} \left(f(\boldsymbol{w}_1) - f^* + \eta \frac{64}{(1-\beta_1)} \sigma_1^2 \frac{\|\boldsymbol{G}_1\|^2}{\sqrt{\beta_2 \widetilde{\boldsymbol{\nu}}_1}} \right)$ $+\frac{1}{1-\beta_2}\left(\frac{147456\eta^2(L_0+L_1)\sigma_1^2\sigma_0}{(1-\beta_1)^{\frac{5}{2}}}+4\frac{L_1\eta^2\sigma_0}{(1-\beta_1)^{\frac{3}{2}}}+\frac{24L_0\eta^2}{1-\beta_1}+8\frac{L_0}{\sigma_0}\eta^2\right)\left(\mathbb{E}\ln\boldsymbol{\nu}_T-T\ln\beta_2\right)\right)$ $\leq \sqrt{\beta_2^T \boldsymbol{\nu}_0 + 2\sigma_0^2} + \sigma_0 + \frac{1}{4} \mathbb{E} \ln \boldsymbol{\nu}_T$ $\leq \sqrt{\beta_2^T \boldsymbol{\nu}_0 + 2\sigma_0^2} + \sigma_0 + \frac{1}{2} \mathbb{E} \sqrt{\boldsymbol{\nu}_T + (1 - \beta_2)\sigma_0^2}$ where the third inequality is due to $T \geq \frac{36 * 2048^4 (L_0 + L_1)^3 \sigma_1^{12} (f(\boldsymbol{w}_1) - f^*)}{(1 - \beta_1)^6 \sigma_0^2} + \frac{768 * 2048^2 (f(\boldsymbol{w}_1) - f^*) \sigma_1^8 (8L_1^2 (f(\boldsymbol{w}_1) - f^*)^2 + 4L_0 (f(\boldsymbol{w}_1) - f^*))}{(1 - \beta_1)^4 \sigma_0^2}$ $+\frac{24^2*147456(L_0+L_1)\sigma_1^8(f(\boldsymbol{w}_1)-f^*)\sigma_0^2}{(1-\beta_2)^5}+\frac{128^2(L_0+L_1)(f(\boldsymbol{w}_1)-f^*)\sigma_1^4}{\sigma^2}$ $+\frac{24^2*147456*2048^2(L_0+L_1)^3\sigma_1^{16}(f(\boldsymbol{w}_1)-f^*)^3}{(1-\beta_2)^{11}}+\frac{128^2*2048^2(L_0+L_1)^3(f(\boldsymbol{w}_1)-f^*)^3\sigma^{12}}{\sigma_0^4(1-\beta_1)^6}$ and the last inequality is due to $\ln x \leq x$. Solving the above inequality with respect to $\mathbb{E}\sqrt{oldsymbol{
u}_T+(1-eta_2)\sigma_0^2}$ and applying $oldsymbol{
u}_0=\sigma_0^2$ then gives

$$\sqrt{\boldsymbol{\nu}_T} \le \mathbb{E}\sqrt{\boldsymbol{\nu}_T + (1 - \beta_2)\sigma_0^2} \le 6\sigma_0.$$
⁽²⁰⁾

 \mathbb{E}

¹⁵¹² Therefore, Eq. (17) can be rewritten as

We then execute the second round of divide-and-conquer. To begin with, we have that

$$\sum_{t=1}^{T} \mathbb{E}\left[\frac{\|\boldsymbol{G}_{t}\|^{2}}{\sqrt{\widetilde{\boldsymbol{\nu}}_{t}}}\mathbb{1}_{\|\boldsymbol{G}_{t}\| \geq \frac{\sigma_{0}}{\sigma_{1}}}\right] \leq \sum_{t=1}^{T} \mathbb{E}\left[\frac{\|\boldsymbol{G}_{t}\|^{2}}{\sqrt{\widetilde{\boldsymbol{\nu}}_{t}}}\right].$$
(22)

On the other hand, we have that

$$\frac{\|\boldsymbol{G}_{t}\|^{2}}{\sqrt{\tilde{\boldsymbol{\nu}}_{t}}} \mathbb{1}_{\|\boldsymbol{G}_{t}\| \geq \frac{\sigma_{0}}{\sigma_{1}}} \geq \frac{\frac{2}{3}\|\boldsymbol{G}_{t}\|^{2} + \frac{1}{3}\frac{\sigma_{0}^{2}}{\sigma_{1}^{2}}}{\sqrt{\tilde{\boldsymbol{\nu}}_{t}}} \mathbb{1}_{\|\boldsymbol{G}_{t}\| \geq \frac{\sigma_{0}}{\sigma_{1}}} \geq \frac{\frac{\beta_{2}}{3\sigma_{1}^{2}}\mathbb{E}^{|\mathcal{F}_{t}}\|\boldsymbol{g}_{t}\|^{2} + \frac{1-\beta_{2}}{3}\frac{\sigma_{0}^{2}}{\sigma_{1}^{2}}}{\sqrt{\tilde{\boldsymbol{\nu}}_{t}}} \mathbb{1}_{\|\boldsymbol{G}_{t}\| \geq \frac{\sigma_{0}}{\sigma_{1}}} \geq \frac{\beta_{2}}{2} \mathbb{E}^{|\mathcal{F}_{t}}\frac{\frac{\beta_{2}}{3\sigma_{1}^{2}}}{\sqrt{\tilde{\boldsymbol{\nu}}_{t}}} \mathbb{I}_{\|\boldsymbol{G}_{t}\| \geq \frac{\sigma_{0}}{\sigma_{1}}} \mathbb{I}_{\|\boldsymbol{G}_{t}\| \geq \frac{\sigma_{0}}{\sigma_{1}}} \geq \frac{1}{2} \mathbb{E}^{|\mathcal{F}_{t}}\frac{\frac{\beta_{2}}{3\sigma_{1}^{2}}}{\sqrt{\tilde{\boldsymbol{\nu}}_{t+1}}} \mathbb{I}_{\|\boldsymbol{G}_{t}\| \geq \frac{\sigma_{0}}{\sigma_{1}}} \mathbb{I}_{\|\boldsymbol{G}_{t}\| \geq \frac{\sigma_{0}}{\sigma_{1}}} \cdot$$

As a conclusion,

$$\sum_{t=1}^{T} \mathbb{E}\left[\frac{\|\boldsymbol{G}_{t}\|^{2}}{\sqrt{\tilde{\boldsymbol{\nu}}_{t}}}\mathbb{1}_{\|\boldsymbol{G}_{t}\|\geq\frac{\sigma_{0}}{\sigma_{1}}}\right] \geq \frac{1}{2}\sum_{t=1}^{T} \mathbb{E}\left[\frac{\frac{\beta_{2}}{3\sigma_{1}^{2}}\|\boldsymbol{g}_{t}\|^{2} + \frac{1-\beta_{2}}{3\sigma_{1}^{2}}\sigma_{0}^{2}}{\sqrt{\tilde{\boldsymbol{\nu}}_{t+1}} + \sqrt{\beta_{2}\tilde{\boldsymbol{\nu}}_{t}}}\mathbb{1}_{\|\boldsymbol{G}_{t}\|\geq\frac{\sigma_{0}}{\sigma_{1}}}\right] \\ \geq \frac{1}{6(1-\beta_{2})\sigma_{1}^{2}}\sum_{t=1}^{T} \mathbb{E}\left[\left(\sqrt{\tilde{\boldsymbol{\nu}}_{t+1}} - \sqrt{\beta_{2}\tilde{\boldsymbol{\nu}}_{t}}\right)\mathbb{1}_{\|\boldsymbol{G}_{t}\|\geq\frac{\sigma_{0}}{\sigma_{1}}}\right].$$

1546 Meanwhile, for convenience, we define $\{\bar{\boldsymbol{\nu}}_t\}_{t=0}^{\infty}$ as $\bar{\boldsymbol{\nu}}_0 = \boldsymbol{\nu}_0$, $\bar{\boldsymbol{\nu}}_t = \beta_2 \bar{\boldsymbol{\nu}}_{t-1} + (1-\beta_2)|g_t|^2 \mathbb{1}_{\|\boldsymbol{G}_t\| < \frac{\sigma_0^2}{\sigma_1^2}}$. 1547 One can easily observe that $\bar{\boldsymbol{\nu}}_t \leq \boldsymbol{\nu}_t$, and thus

$$\begin{split} & \sum_{t=1}^{T} \mathbb{E} \left[\left(\sqrt{\tilde{\nu}_{t+1}} - \sqrt{\beta_2 \tilde{\nu}_t} \right) \mathbb{1}_{\|G_t\| < \frac{\sigma_t^2}{\sigma_1^2}} \right] \\ & = \sum_{t=1}^{T} \mathbb{E} \left(\sqrt{\beta_2^2 \nu_{t-1}} + \beta_2 (1 - \beta_2) \|g_t\|^2 + (1 - \beta_2) \sigma_0^2} - \sqrt{\beta_2 (\beta_2 \nu_{t-1} + (1 - \beta_2) \sigma_0^2)} \right) \mathbb{1}_{\|G_t\| < \frac{\sigma_t^2}{\sigma_1^2}} \\ & = \sum_{t=1}^{T} \mathbb{E} \left(\sqrt{\beta_2^2 \bar{\nu}_{t-1}} + \beta_2 (1 - \beta_2) \|g_t\|^2 + (1 - \beta_2) \sigma_0^2} - \sqrt{\beta_2 (\beta_2 \bar{\nu}_{t-1} + (1 - \beta_2) \sigma_0^2)} \right) \mathbb{1}_{\|G_t\| < \frac{\sigma_t^2}{\sigma_1^2}} \\ & = \sum_{t=1}^{T} \mathbb{E} \left(\sqrt{\beta_2^2 \bar{\nu}_{t-1}} + \beta_2 (1 - \beta_2) \|g_t\|^2 \mathbb{1}_{\|G_t\| < \frac{\sigma_t^2}{\sigma_1^2}} + (1 - \beta_2) \sigma_0^2} - \sqrt{\beta_2 (\beta_2 \bar{\nu}_{t-1} + (1 - \beta_2) \sigma_0^2)} \right) \mathbb{1}_{\|G_t\| < \frac{\sigma_t^2}{\sigma_1^2}} \\ & = \sum_{t=1}^{T} \mathbb{E} \left(\sqrt{\beta_2^2 \bar{\nu}_{t-1}} + \beta_2 (1 - \beta_2) \|g_t\|^2 \mathbb{1}_{\|G_t\| < \frac{\sigma_t^2}{\sigma_1^2}} + (1 - \beta_2) \sigma_0^2} - \sqrt{\beta_2 (\beta_2 \bar{\nu}_{t-1} + (1 - \beta_2) \sigma_0^2)} \right) \\ & = \sum_{t=1}^{T} \mathbb{E} \left(\sqrt{\beta_2 \bar{\nu}_t} + (1 - \beta_2) \sigma_0^2 + \sqrt{\beta_2 (\beta_2 \bar{\nu}_{t-1} + (1 - \beta_2) \sigma_0^2)} - \sqrt{\beta_2 (\beta_2 \bar{\nu}_{t-1} + (1 - \beta_2) \sigma_0^2)} \right) \\ & = \mathbb{E} \sqrt{\beta_2 \bar{\nu}_t} + (1 - \beta_2) \sigma_0^2 + (1 - \sqrt{\beta_2}) \sum_{t=1}^{T-1} \mathbb{E} \sqrt{\beta_2 \bar{\nu}_t} + (1 - \beta_2) \sigma_0^2} - \mathbb{E} \sqrt{\beta_2 (\beta_2 \bar{\nu}_0 + (1 - \beta_2) \sigma_0^2)} \\ & = \mathbb{E} \sqrt{\beta_2 \bar{\nu}_t} + (1 - \beta_2) \sigma_0^2 + (1 - \sqrt{\beta_2}) \sum_{t=1}^{T-1} \mathbb{E} \sqrt{\beta_2 \bar{\nu}_t} + (1 - \beta_2) \sigma_0^2} - \mathbb{E} \sqrt{\beta_2 (\beta_2 \bar{\nu}_0 + (1 - \beta_2) \sigma_0^2)} \\ & = \mathbb{E} \sqrt{\beta_2 \bar{\nu}_t} + (1 - \beta_2) \sigma_0^2 + (1 - \sqrt{\beta_2}) \sum_{t=1}^{T-1} \mathbb{E} \sqrt{\beta_2 \bar{\nu}_t} + (1 - \beta_2) \sigma_0^2} - \mathbb{E} \sqrt{\beta_2 (\beta_2 \bar{\nu}_0 + (1 - \beta_2) \sigma_0^2)} \\ & = \mathbb{E} \sqrt{\beta_2 \bar{\nu}_t} + (1 - \beta_2) \sigma_0^2 + (1 - \sqrt{\beta_2}) \sum_{t=1}^{T-1} \mathbb{E} \sqrt{\beta_2 \bar{\nu}_t} + (1 - \beta_2) \sigma_0^2} \\ & = \mathbb{E} \sqrt{\beta_2 \bar{\nu}_t} + (1 - \beta_2) \sigma_0^2 + (1 - \sqrt{\beta_2}) \sum_{t=1}^{T-1} \mathbb{E} \sqrt{\beta_2 \bar{\nu}_t} + (1 - \beta_2) \sigma_0^2} \\ & = \mathbb{E} \sqrt{\beta_2 \bar{\nu}_t} + (1 - \beta_2) \sigma_0^2 + (1 - \sqrt{\beta_2}) \sum_{t=1}^{T-1} \mathbb{E} \sqrt{\beta_2 \bar{\nu}_t} + (1 - \beta_2) \sigma_0^2} \\ & = \mathbb{E} \sqrt{\beta_2 \bar{\nu}_t} + (1 - \beta_2) \sigma_0^2 + (1 - \sqrt{\beta_2}) \sum_{t=1}^{T-1} \mathbb{E} \sqrt{\beta_2 \bar{\nu}_t} + (1 - \beta_2) \sigma_0^2} \\ & = \mathbb{E} \sqrt{\beta_2 \bar{\nu}_t} + (1 - \beta_2) \sigma_0^2 + (1 - \sqrt{\beta_2}) \sum_{t=1}^{T-1} \mathbb{E} \sqrt{\beta_2 \bar{\nu}_t} \\ & = \mathbb{E} \sqrt{\beta_2 \bar{\nu}_t} + (1 - \beta_2) \sigma_0^2 + (1 - \sqrt{\beta_2}) \sum_{t=1}^{T-1} \mathbb{E} \sqrt{\beta_2 \bar{\nu}_t} + (1 - \beta_2) \sigma_0^2} \\$$

All in all, summing the above two inequalities together, we obtain that

$$\begin{split} & \mathbb{E}\sqrt{\tilde{\nu}_{t+1}} + (1 - \sqrt{\beta_2}) \sum_{t=2}^{T} \mathbb{E}\sqrt{\tilde{\nu}_t} - \sqrt{\beta_2} \tilde{\nu}_1 \\ & 1569 \\ & 1569 \\ \\ & 1570 \\ & 1571 \\ & = \sum_{t=1}^{T} \mathbb{E}\left(\sqrt{\tilde{\nu}_t} - \sqrt{\beta_2} \tilde{\nu}_{t-1}\right) \\ & 1572 \\ & 1573 \\ & \leq \sum_{t=1}^{T} \mathbb{E}\left(\sqrt{\tilde{\nu}_t} - \sqrt{\beta_2} \tilde{\nu}_{t-1}\right) \mathbb{1}_{\|G_t\| \ge \frac{\sigma_0}{\sigma_1}} + \sum_{t=1}^{T} \mathbb{E}\left(\sqrt{\tilde{\nu}_t} - \sqrt{\beta_2} \tilde{\nu}_{t-1}\right) \mathbb{1}_{\|G_t\| < \frac{\sigma_0^2}{\sigma_1^2}} \\ & 1575 \\ & \leq \frac{3(1 - \beta_2)\sigma_1^2}{\sqrt{\beta_2}} \sum_{t=1}^{T} \mathbb{E}\left[\frac{\|G_t\|^2}{\sqrt{\tilde{\nu}_t}}\right] + \mathbb{E}\sqrt{\beta_2 \bar{\nu}_t + (1 - \beta_2)\sigma_0^2} + (1 - \sqrt{\beta_2}) \sum_{t=1}^{T-1} \mathbb{E}\sqrt{\beta_2 \bar{\nu}_t + (1 - \beta_2)\sigma_0^2} - \sqrt{\beta_2(\beta_2 \bar{\nu}_0 + (1 - \beta_2)\sigma_0^2)} . \end{split}$$

$$\mathbb{E}\sqrt{\beta_2 \bar{\boldsymbol{\nu}}_t + (1 - \beta_2)\sigma_0^2} \le \sqrt{\beta_2 \mathbb{E} \bar{\boldsymbol{\nu}}_t + (1 - \beta_2)\sigma_0^2} \le \sqrt{\sigma_0^2 + \boldsymbol{\nu}_0} \le \sqrt{2}\sigma_0,$$

combining with $\sqrt{\beta_2 \widetilde{\boldsymbol{\nu}}_1} = \sqrt{\beta_2 (\beta_2 \overline{\boldsymbol{\nu}}_0 + (1 - \beta_2) \sigma_0^2)}$ and $\mathbb{E}\sqrt{\widetilde{\boldsymbol{\nu}}_{t+1}} = \mathbb{E}\sqrt{\beta_2 \boldsymbol{\nu}_t + (1 - \beta_2) \sigma_0^2} \geq 0$ $\mathbb{E}\sqrt{\beta_2 \bar{\boldsymbol{\nu}}_t + (1-\beta_2)\sigma_0^2}$, we obtain

$$(1 - \sqrt{\beta_2}) \sum_{t=1}^{T} \mathbb{E}\sqrt{\tilde{\nu}_t} \leq \frac{3(1 - \beta_2)\sigma_1^2}{\sqrt{\beta_2}} \sum_{t=2}^{T} \mathbb{E}\left[\frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\tilde{\nu}_t}}\right] + (1 - \sqrt{\beta_2}) \sum_{t=1}^{T} \mathbb{E}\sqrt{\beta_2 \bar{\nu}_t + (1 - \beta_2)\sigma_0^2} \\ \leq \frac{3(1 - \beta_2)\sigma_1^2}{\sqrt{\beta_2}} \sum_{t=1}^{T} \mathbb{E}\left[\frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\tilde{\nu}_t}}\right] + \sqrt{2}(1 - \sqrt{\beta_2})T\sigma_0..$$

Dividing both sides of the above equation by $1 - \sqrt{\beta_2}$ then gives

$$\sum_{t=1}^{T} \mathbb{E}\sqrt{\widetilde{\nu}_{t}} \leq \frac{3(1-\beta_{2})\sigma_{1}^{2}}{\sqrt{\beta_{2}}} \sum_{t=2}^{T} \mathbb{E}\left[\frac{\|\boldsymbol{G}_{t}\|^{2}}{\sqrt{\widetilde{\nu}_{t}}}\right] + (1-\sqrt{\beta_{2}}) \sum_{t=1}^{T} \mathbb{E}\sqrt{\beta_{2}\bar{\nu}_{t}} + (1-\beta_{2})\sigma_{0}^{2}$$
$$\leq 12\sigma_{1}^{2} \sum_{t=1}^{T} \mathbb{E}\left[\frac{\|\boldsymbol{G}_{t}\|^{2}}{\sqrt{\widetilde{\nu}_{t}}}\right] + \sqrt{2}T\sigma_{0}.$$
(23)

By applying Eq. (21) and the constraint of T, we obtain that

$$\sum_{t=1}^{T} \mathbb{E}\sqrt{\widetilde{\nu}_{t}} \leq \sqrt{2}T\sigma_{0} + 12\frac{\sigma_{1}^{2}}{\eta} \left(f(\boldsymbol{w}_{1}) - f^{*} + \eta \frac{64}{(1-\beta_{1})}\sigma_{1}^{2}\frac{\|\boldsymbol{G}_{1}\|^{2}}{\sqrt{\beta_{2}\widetilde{\nu}_{1}}} + \frac{1}{1-\beta_{2}} \left(\frac{147456\eta^{2}(L_{0}+L_{1})\sigma_{1}^{2}\sigma_{0}}{(1-\beta_{1})^{\frac{5}{2}}} + 4\frac{L_{1}\eta^{2}\sigma_{0}}{(1-\beta_{1})^{\frac{3}{2}}} + \frac{24L_{0}\eta^{2}}{1-\beta_{1}} + 8\frac{L_{0}}{\sigma_{0}}\eta^{2} \right) (2\ln 6\sigma_{0} - T\ln \beta_{2}) \right)$$

$$\leq 4T\sigma_{0}.$$

Combining the above inequality and Eq. (21) and applying Cauchy's inequality, we obtain that

$$\begin{aligned} & \begin{bmatrix} 1609 \\ 1610 \\ 1611 \\ 1612 \\ 1612 \\ 1613 \\ 1614 \end{aligned} \\ & \left(\mathbb{E} \sum_{t=1}^{T} \|\nabla f(\boldsymbol{w}_{t})\| \right)^{2} \leq \left(\sum_{t=1}^{T} \mathbb{E} \sqrt{\widetilde{\nu}_{t}} \right) \left(\sum_{t=1}^{T} \mathbb{E} \left[\frac{\|\boldsymbol{G}_{t}\|^{2}}{\sqrt{\widetilde{\nu}_{t}}} \right] \right) \\ & \leq 4T\sigma_{0} \times \frac{1}{\eta} \left(f(\boldsymbol{w}_{1}) - f^{*} + \eta \frac{64}{(1-\beta_{1})} \sigma_{1}^{2} \frac{\|\boldsymbol{G}_{1}\|^{2}}{\sqrt{\beta_{2}\widetilde{\nu}_{1}}} \right) \end{aligned}$$

$$\leq 4T\sigma_0 \times \frac{1}{\eta} \left(f(\boldsymbol{w}_1) - f^* + \eta \frac{\sigma_1}{(1-\beta_1)} \sigma_1^2 \frac{\eta \sigma_1}{\sqrt{\beta_2 \widetilde{\boldsymbol{\nu}}_1}} \right)$$

$$+ \frac{1}{1 - \beta_2} \left(\frac{147456\eta^2 (L_0 + L_1)\sigma_1^2 \sigma_0}{(1 - \beta_1)^{\frac{5}{2}}} + 4 \frac{L_1 \eta^2 \sigma_0}{(1 - \beta_1)^{\frac{3}{2}}} + \frac{24L_0 \eta^2}{1 - \beta_1} + 8 \frac{L_0}{\sigma_0} \eta^2 \right) (2 \ln 6\sigma_0 - T \ln \beta_2) \right)$$

By $\eta = \frac{\sqrt{f(w_1) - f^*}}{\sqrt{L_0 + L_1}\sqrt{T\sigma_0\sigma_1^2}}$ and the constraint of *T*, the proof is completed.

1620 C.2 PROOF FOR SGDM

Theorem 10 (Informal). Fix $L_0 \ge 0$, $L_1 > 0$, and $\Delta_1 \ge 0$, there exists objective function f satisfying (L_0, L_1)-smooth condition and $f(w_1) - f^* = \Delta_1$, and a noise oracle $\mathcal{O}(w, z)$ generating stochastic gradient by $g_t = \nabla f(w_t) + \mathcal{O}(w_t, z_t)$ and satisfying Assumption 2 (z_t is i.i.d. sampled from some underlying distribution), such that for any learning rate $\eta > 0$ and $\beta \in [0, 1]$, for all T > 0,

$$\min_{t \in [T]} \mathbb{E} \|\nabla f(\boldsymbol{w}_t)\| = \|\nabla f(\boldsymbol{w}_1)\| \ge L_1 \Delta_1.$$

Proof. Define the objective function f as f_1 used in the proof of Theorem 2 as

$$f_{1}(x) = \begin{cases} \frac{L_{0}e^{L_{1}x-1}}{L_{1}^{2}} & , x \in \left[\frac{1}{L_{1}}, \infty\right), \\ \frac{L_{0}x^{2}}{2} + \frac{L_{0}}{2L_{1}^{2}} & , x \in \left[-\frac{1}{L_{1}}, \frac{1}{L_{1}}\right], \\ \frac{L_{0}e^{-L_{1}x-1}}{L_{1}^{2}} & , x \in \left(-\infty, -\frac{1}{L_{1}}\right]. \end{cases}$$
(24)

1637 1638 1639 1640 1641 It is easy to verify that f_1 obeys Assumption 1. Then, we set w_1 as the solution of $f_1(x) - \frac{L_0}{2L_1^2} = \Delta_1$, thus $f(w_1) - f^* = \Delta_1$ is satisfied. We then construct the noise oracle as $O_f(w, z) = z$, where $z \sim e^{-\frac{\sqrt{|z|}}{6\sqrt{\sigma_0^2/960}}}$. One can easily verify that $Var(z) = \sigma_0^2$ and Assumption 2 is meet.

Now, we prove the following claim: starting any point w_t and with any previous momentum m_{t-1} , one step of SGDM

$$\mathbb{E}[\|\nabla f(\boldsymbol{w}_{t+1})\| | \boldsymbol{w}_t] = \infty$$

1645 Specifically, we have one step of SGDM gives $(1 - 0)\nabla$

$$\boldsymbol{w}_{t+1} = \boldsymbol{w}_t - \eta(1-\beta)\nabla f(\boldsymbol{w}_t) - \eta\beta\boldsymbol{m}_{t-1} - \eta(1-\beta)\boldsymbol{z}_t.$$

1648 Therefore, we have

$$\begin{aligned}
\mathbf{E}[|\nabla f(w_{t+1})||\boldsymbol{w}_{t}] \geq \mathbb{E}\left[|\nabla f(w_{t+1})|\mathbb{1}_{w_{t+1}\geq \max\{\frac{1}{L_{1}},\boldsymbol{w}_{t}-\eta(1-\beta)\nabla f(\boldsymbol{w}_{t})-\eta\beta\boldsymbol{m}_{t-1}\}}\right] \\
\approx \mathbb{E}\left[|\nabla f(w_{t+1})|\mathbb{1}_{z_{t}\leq \min\{\frac{\boldsymbol{w}_{t}-\frac{1}{L_{1}}}{\eta(1-\beta)}-\nabla f(\boldsymbol{w}_{t})-\frac{\beta}{1-\beta}\boldsymbol{m}_{t-1},0\}}\right] \\
\approx \frac{1}{2}\int_{-\infty}^{\min\{\frac{\boldsymbol{w}_{t}-\frac{1}{L_{1}}}{\eta(1-\beta)}-\nabla f(\boldsymbol{w}_{t})-\frac{\beta}{1-\beta}\boldsymbol{m}_{t-1},0\}}\frac{L_{0}}{L_{1}}e^{L_{1}(\boldsymbol{w}_{t}-\eta(1-\beta)\nabla f(\boldsymbol{w}_{t})-\eta\beta\boldsymbol{m}_{t-1}-\eta(1-\beta)z)-1}e^{-\frac{\sqrt{-z}}{\sqrt[6]{\sigma_{0}^{2}/960}}}dz.
\end{aligned}$$

1650

1662

1664

1626 1627 1628

1629 1630

1633

1635

1644

1647

1657 1658 Since $\lim_{z\to-\infty} e^{L_1(w_t-\eta(1-\beta)\nabla f(w_t)-\eta\beta m_{t-1}-\eta(1-\beta)z)-1}e^{-\frac{\sqrt{-z}}{\sqrt[6]{\sigma_0^2/960}}} = \infty$ regardless of η, β , and m_{t-1} , we have $\mathbb{E}[|\nabla f(w_{t+1})||w_t] = \infty$ based on the above inequalities. This means that an update 1660 from any point over this example will always lead to the divergence on expected gradient norm, thus we have $\forall t > 1$, $\lim_{z\to-\infty} e^{L_1(w_t-\eta(1-\beta)\nabla f(w_t)-\eta\beta m_{t-1}-\eta(1-\beta)z)-1}e^{-\frac{\sqrt{-z}}{\sqrt[6]{\sigma_0^2/960}}} = \infty$ regardless of η, β , and m_{t-1} , we have $\mathbb{E}[|\nabla f(w_{t+1})||w_t] = \infty$ based on the above inequalities. This means that an update $\lim_{z\to-\infty} e^{L_1(w_t-\eta(1-\beta)\nabla f(w_t)-\eta\beta m_{t-1}-\eta(1-\beta)z)-1}e^{-\frac{\sqrt{-z}}{\sqrt[6]{\sigma_0^2/960}}} = \infty$ regardless of η, β , and $\lim_{z\to-\infty} e^{L_1(w_t-\eta(1-\beta)\nabla f(w_t)-\eta\beta m_{t-1}-\eta(1-\beta)z)-1}e^{-\frac{\sqrt{-z}}{\sqrt[6]{\sigma_0^2/960}}} = \infty$ regardless of η, β , and $\lim_{z\to-\infty} e^{L_1(w_t-\eta(1-\beta)\nabla f(w_t)-\eta\beta m_{t-1}-\eta(1-\beta)z)-1}e^{-\frac{\sqrt{-z}}{\sqrt[6]{\sigma_0^2/960}}} = \infty$ regardless of η, β , and $\lim_{z\to-\infty} e^{L_1(w_t-\eta(1-\beta)\nabla f(w_t)-\eta\beta m_{t-1}-\eta(1-\beta)z)-1}e^{-\frac{\sqrt{-z}}{\sqrt[6]{\sigma_0^2/960}}} = \infty$ regardless of η, β , and $\lim_{z\to-\infty} e^{L_1(w_t-\eta(1-\beta)\nabla f(w_t)-\eta\beta m_{t-1}-\eta(1-\beta)z)-1}e^{-\frac{\sqrt{-z}}{\sqrt[6]{\sigma_0^2/960}}} = \infty$ regardless of η, β .

$$\min_{t \in [T]} \mathbb{E} |\nabla f(w_t)| = |\nabla f(w_1)|.$$

1663 The proof is completed.

1665 D PROOFS FOR SECTION 5

⁶⁷ D.1 Proof for Theorem 5

 $\begin{array}{ll} \text{1669} & \text{Theorem 11} (\text{Theorem 5, restated}). \ Let \ Assumption \ 1 \ hold. \ Then, \ \forall \beta_1 \geq 0, \ if \ \varepsilon \leq \\ \hline 1670 & \frac{1}{\operatorname{Poly}(L_0,L_1,\sigma_0,\sigma_1,\frac{1}{1-\beta_1},f(\boldsymbol{w}_1)-f^*)}, \ with \ \eta = (1-\beta_1)\frac{\sqrt{L_0(f(\boldsymbol{w}_1)-f^*)}}{\sqrt{T}} \ and \ momentum \ hyperparam-\\ \hline 1672 & eter \ \beta_2 = 1 - \eta^2 \frac{(256\sigma_1^2 L_1)^2}{1-\beta_1}, \ we \ have \ that \ if \ T \geq \Theta(\frac{L_0\sigma_0^2(f(\boldsymbol{w}_1)-f^*)}{\varepsilon^4}) \\ & \mathbb{E} \ \min_{t \in [1,T]} \|\nabla f(\boldsymbol{w}_t)\| \leq \varepsilon. \end{array}$

Proof. Recall that $u_t \triangleq \frac{w_t - \frac{\beta_1}{\sqrt{\beta_2}} w_{t-1}}{1 - \frac{\beta_1}{\sqrt{\beta_2}}}$. and the surrogate second-order momentum be defined as $\widetilde{\boldsymbol{\nu}}_t \triangleq \beta_2 \boldsymbol{\nu}_{t-1} + (1-\beta_2)\sigma_0^2$. Due to $\frac{\eta}{\sqrt{1-\beta_2}} \leq \frac{\sqrt{1-\beta_1}}{8L_1}$ and following the similar routine as Theorem 3, one can easily verify that

$$\|\boldsymbol{u}_t - \boldsymbol{w}_t\| \le \frac{1}{4L_1}, \|\boldsymbol{u}_{t+1} - \boldsymbol{w}_t\| \le \frac{1}{4L_1}.$$

Therefore, if Lemma 2 can be applied with $w^1 = w_t$, $w^2 = u_{t+1}$, and $w^3 = u_t$, we see the conditions of Lemma 2 is satisfied, which after taking expectation gives

$$\mathbb{E}^{|\mathcal{F}_t} f(\boldsymbol{u}_{t+1}) \\ \leq f(\boldsymbol{u}_t) + \mathbb{E}^{|\mathcal{F}_t} \langle \nabla f(\boldsymbol{w}_t), \boldsymbol{u}_{t+1} - \boldsymbol{u}_t \rangle + \frac{1}{2} (L_0 + L_1 \|\nabla f(\boldsymbol{w}_t)\|) \mathbb{E}^{|\mathcal{F}_t} (\|\boldsymbol{u}_{t+1} - \boldsymbol{w}_t\| + \|\boldsymbol{u}_t - \boldsymbol{w}_t\|) \|\boldsymbol{u}_{t+1} - \boldsymbol{u}_t\|.$$

We call $\langle \nabla f(\boldsymbol{w}_t), \boldsymbol{u}_{t+1} - \boldsymbol{u}_t \rangle$ the first-order term and $\frac{1}{2}(L_0 + L_1 \|\nabla f(\boldsymbol{w}_t)\|)(\|\boldsymbol{u}_{t+1} - \boldsymbol{w}_t\| + \|\boldsymbol{u}_t - \boldsymbol{u}_t\|)$ $|w_t||)||u_{t+1} - u_t||$ the second-order term, as they respectively correspond to the first-order and second-order Taylor's expansion. We then respectively bound these two terms as follows.

Analysis for the first-order term. Similar to bounding the first-order term in the proof of Theorem 3, we have the following decomposition :

$$\begin{aligned} \boldsymbol{u}_{t+1} - \boldsymbol{u}_t &= \frac{\boldsymbol{w}_{t+1} - \boldsymbol{w}_t}{1 - \frac{\beta_1}{\sqrt{\beta_2}}} - \frac{\beta_1}{\sqrt{\beta_2}} \frac{\boldsymbol{w}_t - \boldsymbol{w}_{t-1}}{1 - \frac{\beta_1}{\sqrt{\beta_2}}} \\ &= -\eta \frac{1 - \beta_1}{1 - \frac{\beta_1}{\sqrt{\beta_2}}} \frac{1}{\sqrt{\beta_2 \boldsymbol{\nu}_{t-1}}} \boldsymbol{g}_t - \frac{\eta}{1 - \frac{\beta_1}{\sqrt{\beta_2}}} \left(\frac{1}{\sqrt{\boldsymbol{\nu}_t}} - \frac{1}{\sqrt{\beta_2 \boldsymbol{\nu}_{t-1}}} \right) \boldsymbol{m}_t \end{aligned}$$

 1 77

According to the above decomposition, we have the first-order term can also be decomposed into

$$\begin{aligned} & \mathbb{E}^{|\mathcal{F}_t} \left[\langle \nabla f(\boldsymbol{w}_t), \boldsymbol{u}_{t+1} - \boldsymbol{u}_t \rangle \right] \\ & \text{1705} \\ & \text{1706} \\ & \text{1707} \end{aligned} = \frac{1 - \beta_1}{1 - \frac{\beta_1}{\sqrt{\beta_2}}} \mathbb{E}^{|\mathcal{F}_t} \left[\left\langle \boldsymbol{G}_t, -\eta \frac{1}{\sqrt{\beta_2 \boldsymbol{\nu}_{t-1}}} \boldsymbol{g}_t \right\rangle \right] + \mathbb{E}^{|\mathcal{F}_t} \left[\left\langle \boldsymbol{G}_t, -\frac{\eta}{1 - \frac{\beta_1}{\sqrt{\beta_2}}} \left(\frac{1}{\boldsymbol{\nu}_t} - \frac{1}{\sqrt{\beta_2 \boldsymbol{\nu}_{t-1}}} \right) \boldsymbol{m}_t \right\rangle \right]. \\ & \text{1708} \end{aligned}$$

1709
1710 As
$$\mathbb{E}^{|\mathcal{F}_t}\left[\left\langle \boldsymbol{G}_t, -\eta \frac{1}{\sqrt{\beta_2 \boldsymbol{\nu}_{t-1}}} \boldsymbol{g}_t \right\rangle\right] = -\eta \frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\beta_2 \boldsymbol{\nu}_{t-1}}}$$
, we have
1711

$$\frac{1-\beta_1}{1-\frac{\beta_1}{\sqrt{\beta_2}}} \mathbb{E}^{|\mathcal{F}_t} \left[\left\langle \boldsymbol{G}_t, -\eta \frac{1}{\sqrt{\beta_2 \boldsymbol{\nu}_{t-1}}} \boldsymbol{g}_t \right\rangle \right] \leq -\frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\beta_2 \boldsymbol{\nu}_{t-1}}}$$

1715
1716 We then bound
$$\mathbb{E}^{|\mathcal{F}_t} \left[\left\langle G_t, -\frac{\eta}{1-\frac{\beta_1}{\sqrt{\beta_2}}} \left(\frac{1}{\nu_t} - \frac{1}{\sqrt{\beta_2}\nu_{t-1}} \right) m_t \right\rangle \right]$$
 as follows
1717
1718

1719
1720
$$\mathbb{E}^{|\mathcal{F}_t} \left[\left\langle \boldsymbol{G}_t, -\frac{\eta}{1 - \frac{\beta_1}{\sqrt{\beta_2}}} \left(\frac{1}{\sqrt{\boldsymbol{\nu}_t}} - \frac{1}{\sqrt{\beta_2}\boldsymbol{\nu}_{t-1}} \right) \boldsymbol{m}_t \right\rangle \right]$$
1721
$$\left[\left\langle \boldsymbol{G}_t, -\frac{\eta}{1 - \frac{\beta_1}{\sqrt{\beta_2}}} \left(\frac{1}{\sqrt{\boldsymbol{\nu}_t}} - \frac{1}{\sqrt{\beta_2}\boldsymbol{\nu}_{t-1}} \right) \boldsymbol{m}_t \right\rangle \right]$$

$$= \mathbb{E}^{|\mathcal{F}_t} \left| \left\langle \boldsymbol{G}_t, -\frac{\eta}{1 - \frac{\beta_1}{2\pi}} \left(\frac{(1 - \beta_2) \|\boldsymbol{g}_t\|^2}{\sqrt{\boldsymbol{\nu}_t} \sqrt{\beta_2 \boldsymbol{\nu}_{t-1}} (\sqrt{\boldsymbol{\nu}_t} + \sqrt{\beta_2 \boldsymbol{\nu}_{t-1}})} \right) \boldsymbol{m}_t \right\rangle$$

$$= \begin{bmatrix} \eta \\ \mathbb{R}^{|\mathcal{F}_t|} \end{bmatrix} \begin{bmatrix} \|\mathbf{C}_t\| \begin{pmatrix} (1-\beta_2) \|\mathbf{g}_t\|^2 \\ \|\mathbf{g}_t\| \end{bmatrix} \end{bmatrix}$$

1725
$$\leq \frac{\beta_1}{1 - \frac{\beta_1}{\sqrt{\beta_2}}} \mathbb{E} \left[\|\mathbf{G}_t\| \left(\frac{1}{\sqrt{\nu_t}\sqrt{\beta_2\nu_{t-1}}} (\sqrt{\nu_t} + \sqrt{\beta_2\nu_{t-1}}) \right) \|\mathbf{u}_t\| \right]$$

1720
1727
$$= \frac{\eta}{1 - \frac{\beta_1}{\sqrt{\beta_2}}} \mathbb{E}^{|\mathcal{F}_t} \left[\| \boldsymbol{G}_t \| \left(\frac{(1 - \beta_2) \| \boldsymbol{g}_t \|^2}{\sqrt{\boldsymbol{\nu}_t} \sqrt{\beta_2 \boldsymbol{\nu}_{t-1}} (\sqrt{\boldsymbol{\nu}_t} + \sqrt{\beta_2 \boldsymbol{\nu}_{t-1}})} \right) \| \boldsymbol{m}_t \| \right].$$

Due to Lemma 1, the right-hand-side of the above inequality can be further bounded as $-\frac{\eta}{1-\frac{\beta_1}{\sqrt{\beta_2}}}\mathbb{E}^{|\mathcal{F}_t}\left[\|\boldsymbol{G}_t\|\left(\frac{(1-\beta_2)\|\boldsymbol{g}_t\|^2}{\sqrt{\boldsymbol{\nu}_t}\sqrt{\beta_2\boldsymbol{\nu}_{t-1}}(\sqrt{\boldsymbol{\nu}_t}+\sqrt{\beta_2\boldsymbol{\nu}_{t-1}})}\right)\|\boldsymbol{m}_t\|\right] \leq \frac{\eta(1-\beta_1)}{\left(\sqrt{1-\frac{\beta_1}{/2}}\right)^3}\mathbb{E}^{|\mathcal{F}_t}\left[\|\boldsymbol{G}_t\|\left(\frac{\sqrt{1-\beta_2}\|\boldsymbol{g}_t\|^2}{\sqrt{\beta_2\boldsymbol{\nu}_{t-1}}(\sqrt{\boldsymbol{\nu}_t}+\sqrt{\beta_2\boldsymbol{\nu}_{t-1}})}\right)\right]$ $\stackrel{(\circ)}{\leq} \frac{\eta(1-\beta_1)}{\left(\sqrt{1-\frac{\beta_1}{\sqrt{\beta_2}}}\right)^3} \frac{\|\boldsymbol{G}_t\|}{\sqrt{\beta_2\boldsymbol{\nu}_{t-1}}} \sqrt{\mathbb{E}^{|\mathcal{F}_t|} \|\boldsymbol{g}_t\|^2} \sqrt{\mathbb{E}^{|\mathcal{F}_t|} \frac{\|\boldsymbol{g}_t\|^2}{(\sqrt{\boldsymbol{\nu}_t}+\sqrt{\beta_2\boldsymbol{\nu}_{t-1}})^2}} \stackrel{(\bullet)}{\leq} \frac{\eta(1-\beta_1)\sqrt{1-\beta_2}}{\left(\sqrt{1-\frac{\beta_1}{\sqrt{\beta_2}}}\right)^3} \frac{\|\boldsymbol{G}_t\|}{\sqrt{\beta_2\boldsymbol{\nu}_{t-1}}} \sqrt{\sigma_0^2 + \sigma_1^2 \|\boldsymbol{G}_t\|^2} \sqrt{\mathbb{E}^{|\mathcal{F}_t|} \frac{\|\boldsymbol{g}_t\|^2}{(\sqrt{\boldsymbol{\nu}_t}+\sqrt{\beta_2\boldsymbol{\nu}_{t-1}})^2}}$ $\leq \frac{\eta(1-\beta_1)\sqrt{1-\beta_2}}{\left(\sqrt{1-\frac{\beta_1}{\sqrt{\beta_2}}}\right)^3} \frac{\|\boldsymbol{G}_t\|}{\sqrt{\beta_2\boldsymbol{\nu}_{t-1}}} (\sigma_0+\sigma_1\|\boldsymbol{G}_t\|) \sqrt{\mathbb{E}^{|\mathcal{F}_t|} \frac{\|\boldsymbol{g}_t\|^2}{(\sqrt{\boldsymbol{\nu}_t}+\sqrt{\beta_2\boldsymbol{\nu}_{t-1}})^2}}$ where inequality (\circ) is due to Holder's inequality, and inequality (\bullet) is due to Assumption 2. Apply-ing mean-value inequality respectively to $\frac{\eta(1-\beta_1)\sqrt{1-\beta_2}}{\left(\sqrt{1-\frac{\beta_1}{\sqrt{\beta_2}}}\right)^3} \mathbb{E}^{|\mathcal{F}_t|} \frac{\|\mathcal{G}_t\|}{\sqrt{\beta_2\nu_{t-1}}} \sigma_0 \sqrt{\mathbb{E}^{|\mathcal{F}_t|} \frac{\|g_t\|^2}{(\sqrt{\nu_t}+\sqrt{\beta_2\nu_{t-1}})^2}} \text{ and } \frac{\eta(1-\beta_1)\sqrt{1-\beta_2}}{\left(\sqrt{1-\frac{\beta_1}{\sqrt{\beta_2}}}\right)^3} \mathbb{E}^{|\mathcal{F}_t|} \frac{\|G_t\|}{\sqrt{\beta_2\nu_{t-1}}} \sigma_1 \|G_t\| \sqrt{\mathbb{E}^{|\mathcal{F}_t|} \frac{\|g_t\|^2}{(\sqrt{\nu_t}+\sqrt{\beta_2\nu_{t-1}})^2}}} \text{ and due to } \beta_1 \leq \beta_2, \text{ we obtain that the }$ right-hand-side of the above inequality can be bounded by $\frac{1}{16}\eta \frac{1-\beta_1}{1-\frac{\beta_1}{\sqrt{\beta_2}}} \frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\beta_2 \boldsymbol{\nu}_{t-1}}} + \frac{4\eta(1-\beta_2)\sigma_0^2}{\left(1-\frac{\beta_1}{\sqrt{\beta_2}}\right)^2 \sqrt{\beta_2 \boldsymbol{\nu}_{t-1}}} \mathbb{E}^{|\mathcal{F}_t} \frac{\|\boldsymbol{g}_t\|^2}{(\sqrt{\boldsymbol{\nu}_t}+\sqrt{\beta_2 \boldsymbol{\nu}_{t-1}})^2}$ $+\frac{1}{16}\eta \frac{1-\beta_1}{1-\frac{\beta_1}{\sqrt{\beta_2}}} \frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\beta_2\boldsymbol{\nu}_{t-1}}} + 4\eta \frac{(1-\beta_2)(1-\beta_1)}{(1-\frac{\beta_1}{\sqrt{\beta_2}})^2} \sigma_1^2 \frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\beta_2\boldsymbol{\nu}_{t-1}}} \mathbb{E}^{|\mathcal{F}_t} \frac{\|\boldsymbol{g}_t\|^2}{(\sqrt{\boldsymbol{\nu}_t}+\sqrt{\beta_2\boldsymbol{\nu}_{t-1}})^2}$ $\leq \frac{1}{8} \eta \frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\beta_2 \boldsymbol{\nu}_{t-1}}} + \frac{8\eta (1-\beta_2) \sigma_0^2}{\left(1-\beta_1\right)^2} \mathbb{E}^{|\mathcal{F}_t|} \frac{\|\boldsymbol{g}_t\|^2}{\sqrt{\beta_2 \boldsymbol{\nu}_{t-1}} (\sqrt{\boldsymbol{\nu}_t} + \sqrt{\beta_2 \boldsymbol{\nu}_{t-1}})^2}$

$$+\frac{1}{8}\eta \frac{\|\mathbf{G}_{t}\|}{\sqrt{\beta_{2}\boldsymbol{\nu}_{t-1}}} + 16\eta \frac{(1-\beta_{2})}{(1-\beta_{1})}\sigma_{1}^{2} \frac{\|\mathbf{G}_{t}\|}{\sqrt{\beta_{2}\boldsymbol{\nu}_{t-1}}} \mathbb{E}^{|\mathcal{F}_{t}|} \frac{\|\mathbf{g}_{t}\|}{(\sqrt{\boldsymbol{\nu}_{t}} + \sqrt{\beta_{2}\boldsymbol{\nu}_{t-1}})^{2}}$$

 $\left(\frac{1}{\sqrt{eta_2 oldsymbol{
u}_{t-1}}} - \frac{1}{\sqrt{oldsymbol{
u}_t}}
ight) \|oldsymbol{G}_t\|^2$

Meanwhile, we have

and

$$-\frac{1}{\sqrt{\beta_2 \boldsymbol{\nu}_{t-1}}} - \frac{1}{\sqrt{\boldsymbol{\nu}_t}} = \frac{(1-\beta_2) \|\boldsymbol{g}_t\|^2}{\sqrt{\beta_2 \boldsymbol{\nu}_{t-1}} \sqrt{\boldsymbol{\nu}_t} (\sqrt{\beta_2 \boldsymbol{\nu}_{t-1}} + \sqrt{\boldsymbol{\nu}_t})} \ge \frac{(1-\beta_2) \|\boldsymbol{g}_t\|^2}{\sqrt{\beta_2 \boldsymbol{\nu}_{t-1}} (\sqrt{\beta_2 \boldsymbol{\nu}_{t-1}} + \sqrt{\boldsymbol{\nu}_t})^2}.$$

 $=\frac{\|\boldsymbol{G}_t\|^2((1-\beta_2)\|\boldsymbol{g}_t\|^2)}{\sqrt{\beta_2\boldsymbol{\nu}_{t-1}}\sqrt{\boldsymbol{\nu}_t}(\sqrt{\beta_2\boldsymbol{\nu}_{t-1}}+\sqrt{\boldsymbol{\nu}_t})} \geq \frac{\|\boldsymbol{G}_t\|^2((1-\beta_2)\|\boldsymbol{g}_t\|^2)}{\sqrt{\beta_2\boldsymbol{\nu}_{t-1}}(\sqrt{\beta_2\boldsymbol{\nu}_{t-1}}+\sqrt{\boldsymbol{\nu}_t})^2}$

Combing the above two inequalities, we further obtain

$$\begin{array}{l} & \begin{array}{l} 1771 \\ 1772 \\ 1773 \\ 1773 \\ 1774 \\ 1775 \\ 1776 \\ 1776 \\ 1776 \end{array} & \begin{array}{l} \frac{\eta}{1-\beta_1} \mathbb{E}^{|\mathcal{F}_t} \left[\|\boldsymbol{G}_t\| \left(\frac{(1-\beta_2)\boldsymbol{g}_t^2}{\sqrt{\boldsymbol{\nu}_t}\sqrt{\beta_2\boldsymbol{\nu}_{t-1}}(\sqrt{\boldsymbol{\nu}_t}+\sqrt{\beta_2\boldsymbol{\nu}_{t-1}})} \right) \|\boldsymbol{m}_t\| \right] \\ 1774 \\ 1775 \\ 1776 \\ 1776 \end{array} & \begin{array}{l} \leq \frac{1}{4} \eta \frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\beta_2\boldsymbol{\nu}_{t-1}}} + \frac{8\eta\sigma_0^2}{(1-\beta_1)^2} \mathbb{E}^{|\mathcal{F}_t} \left(\frac{1}{\sqrt{\beta_2\boldsymbol{\nu}_{t-1}}} - \frac{1}{\sqrt{\boldsymbol{\nu}_t}} \right) + \eta \frac{16}{(1-\beta_1)} \sigma_1^2 \mathbb{E}^{|\mathcal{F}_t} \left(\frac{1}{\sqrt{\beta_2\boldsymbol{\nu}_{t-1}}} - \frac{1}{\sqrt{\boldsymbol{\nu}_t}} \right) \|\boldsymbol{G}_t\|^2. \end{array}$$

Furthermore, due to Assumption 1, we have (we define $G_0 \triangleq G_1$)

780
781
$$\|\boldsymbol{G}_{t+1}\|^{2} \leq \|\boldsymbol{G}_{t}\|^{2} + 2\|\boldsymbol{G}_{t}\|\|\boldsymbol{G}_{t+1} - \boldsymbol{G}_{t}\| + \|\boldsymbol{G}_{t+1} - \boldsymbol{G}_{t}\|^{2} \\ \leq \|\boldsymbol{G}_{t}\|^{2} + 2(L_{0} + L_{1}\|\boldsymbol{G}_{t}\|)\|\boldsymbol{G}_{t}\|\|\boldsymbol{w}_{t+1} - \boldsymbol{w}_{t}\| + 2(L_{0}^{2} + L_{1}^{2}\|\boldsymbol{G}_{t}\|^{2})\|\boldsymbol{w}_{t+1} - \boldsymbol{w}_{t}\|^{2}$$

which further leads to

$$\frac{1}{\sqrt{\nu_{t}}} \|G_{t}\|^{2} = \frac{1}{\sqrt{\nu_{t}}} \|G_{t+1}\|^{2} - 2(L_{0} + L_{1}\|G_{t}\|) \|G_{t}\| \|w_{t+1} - w_{t}\| - 2(L_{0}^{2} + L_{1}^{2}\|G_{t}\|^{2}) \|w_{t+1} - w_{t}\|^{2}) = \frac{1}{\sqrt{\nu_{t}}} \left(\|G_{t+1}\|^{2} - 2(L_{0} + L_{1}\|G_{t}\|) \|G_{t}\| \|w_{t+1} - w_{t}\| - 2(L_{0}^{2} + L_{1}^{2}\|G_{t}\|^{2}) \|w_{t+1} - w_{t}\|^{2} \right) = \frac{1}{\sqrt{\nu_{t}}} \|G_{t+1}\|^{2} - \frac{1 - \beta_{1}}{128\sigma_{1}^{2}} \frac{\|G_{t}\|^{2}}{\sqrt{\nu_{t}}} - \frac{128L_{0}^{2}\sigma_{1}^{2}}{(1 - \beta_{1})\sqrt{\nu_{t}}} \|w_{t+1} - w_{t}\|^{2} - 2L_{1} \frac{\|G_{t}\|^{2}}{\sqrt{\nu_{t}}} \|w_{t+1} - w_{t}\| - 2\frac{L_{0}^{2}}{\sqrt{\nu_{t}}} \|w_{t+1} - w_{t}\|^{2} - \frac{2L_{1}^{2}\|G_{t}\|^{2}\|w_{t+1} - w_{t}\|^{2}}{\sqrt{\nu_{t}}} = \frac{128L_{0}^{2}\sigma_{1}^{2}}{(1 - \beta_{1})\sqrt{\nu_{t}}} \|w_{t+1} - w_{t}\|^{2} - \frac{1 - \beta_{1}}{128\sigma_{1}^{2}} \frac{\|G_{t}\|^{2}}{\sqrt{\nu_{t}}} + \frac{1 - \beta_{1}}{128\sigma_$$

where the second inequality is due to Young's inequality, and the last inequality is due to $\|w_{t+1} - w_{t+1}\|$ $\| w_t \| \leq rac{\eta(1-eta_1)}{\sqrt{1-eta_2}\sqrt{1-rac{eta_1^2}{eta_2^2}}} \leq rac{\eta\sqrt{1-eta_1}}{\sqrt{1-eta_2}} \leq rac{1-eta_1}{256\sigma_1^2 L_1}.$

All in all, we conclude that the first-order term can be bounded by

$$\begin{split} & \mathbb{E}^{|\mathcal{F}_{t}}\left[\langle \nabla f(\boldsymbol{w}_{t}), \boldsymbol{u}_{t+1} - \boldsymbol{u}_{t} \rangle\right] \leq -\frac{3}{8} \eta \frac{\|\boldsymbol{G}_{t}\|^{2}}{\sqrt{\beta_{2}\boldsymbol{\nu}_{t-1}}} + \frac{8\eta\sigma_{0}^{2}}{(1-\beta_{1})^{2}} \mathbb{E}^{|\mathcal{F}_{t}}\left(\frac{1}{\sqrt{\beta_{2}\boldsymbol{\nu}_{t-1}}} - \frac{1}{\sqrt{\boldsymbol{\nu}_{t}}}\right) + \eta \frac{16}{(1-\beta_{1})}\sigma_{1}^{2} \mathbb{E}^{|\mathcal{F}_{t}}\left(\frac{\|\boldsymbol{G}_{t}\|^{2}}{\sqrt{\beta_{2}\boldsymbol{\nu}_{t-1}}} - \frac{\|\boldsymbol{G}_{t+1}\|^{2}}{\sqrt{\boldsymbol{\nu}_{t}}}\right) \\ & + \frac{32768L_{0}^{2}\sigma_{1}^{4}\eta^{3}}{(1-\beta_{1})(1-\beta_{2})} \mathbb{E}^{|\mathcal{F}_{t}}\left(\sum_{s=1}^{t}\frac{\beta_{1}^{t-s}}{\sqrt{\beta_{2}^{3}(t-s)}}\left(\frac{1}{\sqrt{\beta_{2}\boldsymbol{\nu}_{s-1}}} - \frac{1}{\sqrt{\boldsymbol{\nu}_{s}}}\right)\right). \end{split}$$

Analysis for the second-order term. To recall, the second order term is $\frac{1}{2}(L_0 +$ $L_1 \|\nabla f(\boldsymbol{w}_t)\|)(\|\boldsymbol{u}_{t+1} - \boldsymbol{w}_t\| + \|\boldsymbol{u}_t - \boldsymbol{w}_t\|)\|\boldsymbol{u}_{t+1} - \boldsymbol{u}_t\|$. Before we start, we have the follow-ing expansion for $u_{t+1} - u_t$: (here the operations are all coordinate-wisely)

$$\boldsymbol{u}_{t+1} - \boldsymbol{u}_t = \frac{\boldsymbol{w}_{t+1} - \boldsymbol{w}_t - \frac{\beta_1}{\sqrt{\beta_2}} (\boldsymbol{w}_t - \boldsymbol{w}_{t-1})}{1 - \frac{\beta_1}{\sqrt{\beta_2}}} \\ - \eta \frac{\boldsymbol{m}_t}{\sqrt{\boldsymbol{\nu}_t}} + \eta \frac{\beta_1}{\sqrt{\beta_2}} \frac{\boldsymbol{m}_{t-1}}{\sqrt{\boldsymbol{\nu}_{t-1}}} - \eta \frac{\boldsymbol{m}_t}{\sqrt{\boldsymbol{\nu}_t}} + \eta \beta_1 \frac{\boldsymbol{m}_{t-1}}{\sqrt{\boldsymbol{\nu}_t}} - \eta \beta_1 \frac{\boldsymbol{m}_{t-1}}{\sqrt{\boldsymbol{\nu}_t}} + \eta \frac{\beta_1}{\sqrt{\beta_2}} \frac{\boldsymbol{m}_{t-1}}{\sqrt{\boldsymbol{\nu}_{t-1}}}$$

$$= \frac{1 - \frac{\beta_1}{\sqrt{\beta_2}}}{1 - \frac{\beta_1}{\sqrt{\beta_2}}} = \frac{1 - \frac{\beta_1}{\sqrt{\beta_2}}}{1 - \frac{\beta_1}{\sqrt{\beta_2}}}$$

1833
1834
$$= \frac{-\eta \frac{(1-\beta_1)g_t}{\sqrt{\nu_t}} + \eta \frac{\beta_1(1-\beta_2) \|g_t\|^2}{\sqrt{\beta_2}} \frac{m_{t-1}}{\sqrt{\nu_{t-1}}\sqrt{\nu_t}(\sqrt{\nu_t} + \sqrt{\beta_2\nu_{t-1}})}}{\frac{m_{t-1}}{\sqrt{\nu_t}}}$$

$$\begin{split} & \text{Then firstly, we have} & \frac{1}{2} I_{2}(|u_{n+1}-u_{n}|| + ||u_{r}-u_{n}||)||u_{n+1}-u_{n}||^{2}}{2} \\ & \leq \frac{1}{2} I_{2}\left(\left\| ||u_{n+1}-u_{n}||^{2} + \frac{1}{2} ||u_{n+1}-u_{n}||^{2} + \frac{1}{2} ||u_{n}-u_{n}||^{2} \right) \\ & = \frac{1}{2} I_{2}\left(\left\| \left\| -\frac{q^{(1-q)}(u_{n})-q_{n}}{1-\frac{1}{\sqrt{2}}} + \frac{1}{2} ||u_{n}-u_{n}||^{2} + \frac{1}{2} ||u_{n}-u_{n}||^{2} \right) \\ & = \frac{1}{2} I_{2}\left(\left\| \left(\frac{1-g_{n}}{1-\frac{1}{\sqrt{2}}} + \frac{g_{n}(1-g_{n})}{\sqrt{2}} - \frac{g_{n}(1-g_{n})}{\sqrt{2}} \right) \right\| \frac{g_{n}}{\sqrt{2\pi}} + \frac{1}{\sqrt{2\pi}} \left(\frac{1}{\sqrt{2}} + \frac{g_{n}(1-g_{n})}{\sqrt{2}} \right) \\ & \leq \frac{Lov^{2}}{2} \left(\left(\frac{1-g_{n}}{1-\frac{1}{\sqrt{2}}} + \frac{g_{n}(1-g_{n})}{\sqrt{2}} - \frac{1}{g_{n}} \right) \right) \left\| \frac{g_{n}}{\sqrt{2\pi}} \right|^{2} + \frac{1}{2} \left(\frac{1}{\sqrt{2}} + \frac{g_{n}}{\sqrt{2}} \right) \\ & = \frac{Lov^{2}}{2} \left(2 \left(\frac{1-g_{n}}{1-\frac{1}{\sqrt{2}}} + \frac{g_{n}(1-g_{n})}{\sqrt{2}} - \frac{1}{g_{n}} \right) \right) \\ & \leq \frac{Lov^{2}}{2} \left(2 \left(\frac{1-g_{n}}{1-\frac{1}{\sqrt{2}}} + \frac{g_{n}(1-g_{n})}{\sqrt{2}} \right) \left\| \frac{g_{n}}{\sqrt{2}} \right\|^{2} + \left(\frac{1}{\sqrt{2}} + \frac{g_{n}}{\sqrt{2}} \right)^{2} \left\| \frac{g_{n}}{\sqrt{2}} \right\|^{2} \right) \\ & \leq \frac{Lov^{2}}{2} \left(2 \left(\frac{1-g_{n}}{1-\frac{1}{\sqrt{2}}} + \frac{g_{n}(1-g_{n})}{\sqrt{2}} \right) \left\| \frac{g_{n}}{\sqrt{2}} \right\|^{2} + \left(\frac{1}{\sqrt{2}} + \frac{g_{n}}{\sqrt{2}} \right)^{2} \right\| \frac{g_{n}}{\sqrt{2}} \right\|^{2} \right) \\ & \leq \frac{Lov^{2}}{2} \left(2 \left(\frac{1-g_{n}}{1-\frac{1}{\sqrt{2}}} + \frac{g_{n}(1-g_{n})}{\sqrt{2}} \right) \left\| \frac{g_{n}}{\sqrt{2}} \right\|^{2} + \left(\frac{1}{\sqrt{2}} + \frac{g_{n}}{\sqrt{2}} \right)^{2} \right\| \frac{g_{n}}{\sqrt{2}} \right\|^{2} \right) \\ & \leq \frac{Lov^{2}}{2} \left(2 \left(\frac{1}{1-g_{n}} + \frac{g_{n}}{\sqrt{2}} \right) \left\| \frac{g_{n}}{\sqrt{2}} \right\| \frac{g_{n}}{\sqrt{2}} \right\|^{2} + \left(\frac{1}{\sqrt{2}} + \frac{g_{n}}{\sqrt{2}} \right)^{2} \right\| \frac{g_{n}}{\sqrt{2}} \right\|^{2} \right) \\ & \leq \frac{Lov^{2}}{2} \left(2 \left(\frac{1-g_{n}}{\sqrt{2}} + \frac{g_{n}}{\sqrt{2}} \right) \left\| \frac{g_{n}}{\sqrt{2}} \right\|^{2} + \left(\frac{1-g_{n}}{\sqrt{2}} + \frac{g_{n}}{\sqrt{2}} \right) \left\| \frac{g_{n}}{\sqrt{2}} \right\|^{2} \right) \\ & \leq \frac{Lov^{2}}{2} \left(2 \left(\frac{1-g_{n}}{\sqrt{2}} + \frac{g_{n}}{\sqrt{2}} \right) \left\| \frac{g_{n}}{\sqrt{2}} \right\|^{2} \\ & = \frac{1}{2} \int \frac{g_{n}}{\sqrt{2}} \left\| \frac{g_{n}}{\sqrt{2}} \right\| \frac{g_{n}}{\sqrt{2}} \right\|^{2} \\ & = \frac{1}{2} \int \frac{g_{n}}{\sqrt{2}} \left\| \frac{g_{n}}{\sqrt{2}} \right\|^{2} \\ & = \frac{1}{2} \int \frac{g_{n}}{\sqrt{2}} \left\| \frac{g_{n}}{\sqrt{2}} \right\|^{2} \\ & \leq \frac{1}{2} \int \frac{g_{n}}{\sqrt{2}} \left\| \frac{g_{n}}{\sqrt{2}} \right\|^{2} \\ & = \frac{1}{2} \int \frac{g_{n}}}{\sqrt{2}} \left\| \frac{g_{n}}{\sqrt{2}} \right\|^$$

By applying Lemma 5, we further obtain

$$\begin{split} & \mathbb{E}^{|\mathcal{F}_{t}} \frac{\|\boldsymbol{G}_{t}\| \|\boldsymbol{m}_{t}\| \|\boldsymbol{g}_{t}\|}{\boldsymbol{\nu}_{t}} \\ & \leq \frac{\sqrt{(1-\beta_{1})^{3}}}{256\eta L_{1}} \frac{\|\boldsymbol{G}_{t}\|^{2}}{\sqrt{\beta_{2}\boldsymbol{\nu}_{t-1}}} + \frac{64\eta\sigma_{0}^{2}L_{1}}{\sqrt{(1-\beta_{1})^{3}}\sqrt{\beta_{2}\boldsymbol{\nu}_{t-1}}} \mathbb{E}^{|\mathcal{F}_{t}} \frac{\|\boldsymbol{m}_{t}\|^{2}}{\boldsymbol{\nu}_{t}} + \frac{\sqrt{(1-\beta_{1})^{3}}}{256\eta L_{1}} \frac{\|\boldsymbol{G}_{t}\|^{2}}{\sqrt{\beta_{2}\boldsymbol{\nu}_{t-1}}} \\ & + \frac{64\eta\sigma_{1}^{2}L_{1}}{\sqrt{(1-\beta_{1})^{3}}} \mathbb{E}^{|\mathcal{F}_{t}} \left(4(1-\beta_{1}) \left(\sum_{s=1}^{t} \frac{\sqrt[8]{\beta_{1}^{t-s}}}{\sqrt{\beta_{2}\boldsymbol{\nu}_{s-1}}\sqrt{\boldsymbol{\nu}_{s}^{2}}} \right) + 8\frac{1-\beta_{1}}{1-\beta_{2}} \frac{L_{1}^{2}}{L_{0}^{2}} \left(\sum_{s=1}^{t} \sqrt[8]{\beta_{1}^{t-s}} \left(\frac{1}{\sqrt{\beta_{2}\boldsymbol{\nu}_{s-1}}} - \frac{1}{\sqrt{\boldsymbol{\nu}_{s}}} \right) \right) \right) \end{split}$$

which further indicates that

$$\begin{split} & \frac{8L_1\eta^2}{(1-\beta_1)^{\frac{3}{2}}} \mathbb{E}^{|\mathcal{F}_t} \| \mathbf{G}_t \| \frac{\| \mathbf{g}_t \|}{\sqrt{\nu_t}} \frac{\| \mathbf{m}_t \|}{\sqrt{\nu_t}} \\ & \leq \frac{1}{16} \eta \frac{\| \mathbf{G}_t \|^2}{\sqrt{\beta_2 \nu_{t-1}}} + \frac{8L_1\eta^2}{(1-\beta_1)^{\frac{3}{2}}} \frac{64\eta \sigma_0^2 L_1}{\sqrt{(1-\beta_1)^3}\sqrt{\beta_2 \nu_{t-1}}} \mathbb{E}^{|\mathcal{F}_t} \frac{\| \mathbf{m}_t \|^2}{\nu_t} \\ & + \frac{64\eta \sigma_1^2 L_1}{\sqrt{(1-\beta_1)^3}} \mathbb{E}^{|\mathcal{F}_t} \frac{32L_1\eta^2}{(1-\beta_1)^{\frac{1}{2}}} \left(\sum_{s=1}^t \frac{\sqrt[8]{\beta_1^{t-s}} \| \mathbf{g}_s \|^2 \| \mathbf{G}_s \|^2}{\sqrt{\beta_2 \nu_{s-1}} \sqrt{\nu_s^2}} \right) \\ & + \frac{64L_1\eta^2}{(1-\beta_1)^{\frac{3}{2}}} \frac{64\eta \sigma_1^2 L_1}{\sqrt{(1-\beta_1)^3}} \frac{1-\beta_1}{1-\beta_2} \frac{L_1^2}{L_0^2} \left(\sum_{s=1}^t \sqrt[8]{\beta_1^{t-s}} \left(\frac{1}{\sqrt{\beta_2 \nu_{s-1}}} - \frac{1}{\sqrt{\nu_s}} \right) \right) \\ & \leq \frac{1}{16} \eta \frac{\| \mathbf{G}_t \|^2}{\sqrt{\beta_2 \nu_{t-1}}} + \frac{8L_1\eta^2}{(1-\beta_1)^{\frac{3}{2}}} \frac{64\eta \sigma_0^2 L_1}{\sqrt{(1-\beta_1)^3}} \mathbb{E}^{|\mathcal{F}_t} 4(1-\beta_1) \left(\sum_{s=1}^t \sqrt[4]{\beta_1^{t-s}} \frac{2}{1-\beta_2} \left(\frac{1}{\sqrt{\beta_2 \nu_{s-1}}} - \frac{1}{\sqrt{\nu_s}} \right) \right) \\ & + \frac{64\eta \sigma_1^2 L_1}{\sqrt{(1-\beta_1)^3}} \mathbb{E}^{|\mathcal{F}_t} \frac{32L_1\eta^2}{(1-\beta_1)^{\frac{1}{2}}} \left(\sum_{s=1}^t \frac{\sqrt[8]{\beta_1^{t-s}} \| \mathbf{g}_s \|^2 \| \mathbf{G}_s \|^2}{\sqrt{\beta_2 \nu_{s-1}} \sqrt{\nu_s^2}} \right) \\ & + \frac{64\eta \sigma_1^2 L_1}{\sqrt{(1-\beta_1)^3}} \mathbb{E}^{|\mathcal{F}_t} \frac{32L_1\eta^2}{(1-\beta_1)^{\frac{1}{2}}} \left(\sum_{s=1}^t \frac{\sqrt[8]{\beta_1^{t-s}} \| \mathbf{g}_s \|^2 \| \mathbf{G}_s \|^2}{\sqrt{\beta_2 \nu_{s-1}} \sqrt{\nu_s^2}} \right) \\ & + \frac{64L_1\eta^2}{(1-\beta_1)^{\frac{3}{2}}} \frac{64\eta \sigma_1^2 L_1}{\sqrt{(1-\beta_1)^3}} \frac{1-\beta_1}{1-\beta_2} \frac{L_1^2}{L_0^2} \left(\sum_{s=1}^t \sqrt[8]{\beta_1^{t-s}}} \left(\frac{1}{\sqrt{\beta_2 \nu_{s-1}}} - \frac{1}{\sqrt{\nu_s}} \right) \right). \\ & = \frac{1}{100} \left(\frac{1}{\sqrt{(1-\beta_1)^3}} \frac{1-\beta_1}{\sqrt{(1-\beta_1)^3}} \frac{1-\beta_1}{1-\beta_2} \frac{L_1^2}{L_0^2} \left(\sum_{s=1}^t \sqrt[8]{\beta_1^{t-s}}} \left(\frac{1}{\sqrt{\beta_2 \nu_{s-1}}} - \frac{1}{\sqrt{\nu_s}} \right) \right) \\ & = \frac{1}{100} \left(\frac{1}{\sqrt{(1-\beta_1)^3}} \frac{1-\beta_1}{\sqrt{(1-\beta_1)^3}} \frac{1-\beta_1}{1-\beta_2} \frac{L_1^2}{L_0^2} \left(\sum_{s=1}^t \sqrt[8]{\beta_1^{t-s}} \left(\frac{1}{\sqrt{\beta_2 \nu_{s-1}}} - \frac{1}{\sqrt{\nu_s}} \right) \right). \\ & = \frac{1}{100} \left(\frac{1}{\sqrt{(1-\beta_1)^3}} \frac{1-\beta_1}{\sqrt{(1-\beta_1)^3}} \frac{1-\beta_1}{1-\beta_2} \frac{L_1^2}{L_0^2} \left(\sum_{s=1}^t \sqrt[8]{\beta_1^{t-s}} \left(\frac{1}{\sqrt{\beta_2 \nu_{s-1}}} - \frac{1}{\sqrt{\nu_s}} \right) \right) \\ & = \frac{1}{100} \left(\frac{1}{\sqrt{(1-\beta_1)^3}} \frac{1-\beta_1}{\sqrt{(1-\beta_1)^3}} \frac{1-\beta_1}{1-\beta_2} \frac{L_1^2}{L_0^2} \left(\sum_{s=1}^t \sqrt{\beta_1^{t-s}} \left(\frac{1}{\sqrt{\beta_2 \nu_{s-1}}} - \frac{1}{\sqrt{\nu_s}} \right) \right) \\ & = \frac{1}{100} \left(\frac{1}{\sqrt{1-\beta_1$$

Here the last inequality is due to Lemma 4.

Following similar reasoning, we have $\|\boldsymbol{u}_{t+1} - \boldsymbol{u}_t\| \le \frac{4\eta}{\sqrt{1-\beta_2}} \frac{\|\boldsymbol{g}_t\|}{\sqrt{\nu_t}}$, and

$$\begin{split} & \mathbb{E}^{|\mathbf{F}_{t}} \frac{\|\mathbf{G}_{t}\| \|\mathbf{g}_{t}\| \|\mathbf{g}_{t}\|}{\boldsymbol{\nu}_{t}} \\ & \mathbb{E}^{|\mathbf{F}_{t}} \frac{\|\mathbf{G}_{t}\| \|\mathbf{g}_{t}\| \|\mathbf{g}_{t}\|}{\boldsymbol{\nu}_{t}} \\ & \leq \frac{\sqrt{(1-\beta_{1})^{3}}}{256\eta L_{1}} \frac{\|\mathbf{G}_{t}\|^{2}}{\sqrt{\beta_{2}\boldsymbol{\nu}_{t-1}}} + \frac{64\eta\sigma_{0}^{2}L_{1}}{\sqrt{(1-\beta_{1})^{3}}\sqrt{\beta_{2}\boldsymbol{\nu}_{t-1}}} \mathbb{E}^{|\mathbf{F}_{t}} \frac{\|\mathbf{g}_{t}\|^{2}}{\boldsymbol{\nu}_{t}} + \frac{\sqrt{(1-\beta_{1})^{3}}}{256\eta L_{1}} \frac{\|\mathbf{G}_{t}\|^{2}}{\sqrt{\beta_{2}\boldsymbol{\nu}_{t-1}}} \\ & + \frac{64\eta\sigma_{1}^{2}L_{1}\|\mathbf{G}_{t}\|^{2}}{\sqrt{(1-\beta_{1})^{3}}\sqrt{\beta_{2}\boldsymbol{\nu}_{t-1}}} \mathbb{E}^{|\mathbf{F}_{t}} \frac{\|\mathbf{g}_{t}\|^{2}}{\boldsymbol{\nu}_{t}} \\ \\ & = \frac{\sqrt{(1-\beta_{1})^{3}}}{256\eta L_{1}} \frac{\|\mathbf{G}_{t}\|^{2}}{\sqrt{\beta_{2}\boldsymbol{\nu}_{t-1}}} + \frac{128\eta\sigma_{1}^{2}L_{1}}{\sqrt{(1-\beta_{1})^{3}}(1-\beta_{2})} \mathbb{E}^{|\mathbf{F}_{t}} \left(\frac{1}{\sqrt{\beta_{2}\boldsymbol{\nu}_{t-1}}} - \frac{1}{\sqrt{\boldsymbol{\nu}_{t}}}\right) + \frac{\sqrt{(1-\beta_{1})^{3}}}{256\eta L_{1}} \frac{\|\mathbf{G}_{t}\|^{2}}{\sqrt{\beta_{2}\boldsymbol{\nu}_{t-1}}} \\ & + \frac{128\eta\sigma_{1}^{2}L_{1}\|\mathbf{G}_{t}\|^{2}}{\sqrt{(1-\beta_{1})^{3}}(1-\beta_{2})} \mathbb{E}^{|\mathbf{F}_{t}} \left(\frac{1}{\sqrt{\beta_{2}\boldsymbol{\nu}_{t-1}}} - \frac{1}{\sqrt{\boldsymbol{\nu}_{t}}}\right). \end{split}$$

Then, following the similar routine as Eq. (25) and due to $\frac{\eta}{\sqrt{1-\beta_2}} \leq \frac{1-\beta_1}{128L_1\sigma_1}$, we have $\frac{2L_1\eta}{(1-\beta_t)^{\frac{1}{2}}} \mathbb{E}^{|\mathcal{F}_t|} \|\boldsymbol{G}_t\| \|\boldsymbol{u}_{t+1} - \boldsymbol{u}_t\| \frac{\|\boldsymbol{g}_t\|}{\sqrt{\boldsymbol{\nu}_t}} \le \frac{8L_1\eta^2}{(1-\beta_t)} \mathbb{E}^{|\mathcal{F}_t|} \|\boldsymbol{G}_t\| \frac{\|\boldsymbol{g}_t\|^2}{\boldsymbol{\nu}_t}$ $\leq \frac{1}{16} \frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\beta_2 \boldsymbol{\nu}_{t-1}}} + \frac{8L_1 \eta^2}{(1-\beta_1)^{\frac{3}{2}}} \frac{128\eta \sigma_1^2 L_1}{\sqrt{(1-\beta_1)^3}(1-\beta_2)} \mathbb{E}^{|\mathcal{F}_t} \left(\frac{1}{\sqrt{\beta_2 \boldsymbol{\nu}_{t-1}}} - \frac{1}{\sqrt{\boldsymbol{\nu}_t}}\right)$ $+\frac{8L_1\eta^2}{(1-\beta_1)^{\frac{3}{2}}}\frac{128\eta\sigma_1^2L_1}{\sqrt{(1-\beta_1)^3}(1-\beta_2)}\mathbb{E}^{|\mathcal{F}_t}\left(\frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\beta_2\boldsymbol{\nu}_{t-1}}}-\frac{\|\boldsymbol{G}_{t+1}\|^2}{\sqrt{\boldsymbol{\nu}_t}}\right)$ $+\frac{1}{16}\eta \frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\beta_2 \boldsymbol{\nu}_{t-1}}} + \eta^3 \frac{64L_0^2}{(1-\beta_1)^2} \mathbb{E}^{|\mathcal{F}_t} 4(1-\beta_1) \left(\sum_{i=1}^t \sqrt[4]{\beta_1^{t-s}} \frac{2}{1-\beta_2} \left(\frac{1}{\sqrt{\beta_2 \boldsymbol{\nu}_{s-1}}} - \frac{1}{\sqrt{\boldsymbol{\nu}_s}}\right)\right).$

Putting all the estimations together, we have that the second-order term can be bounded by (note here due to the complexity of coefficients, we use $Poly(L_0, L_1, \sigma_0, \sigma_1, \frac{1}{1-\beta_1}, f(\boldsymbol{w}_1)$ $f^*)(L_0, L_1, \sigma_0, \sigma_1, \frac{1}{1-\beta_1})$ to denote the polynomial function of $L_0, L_1, \sigma_0, \sigma_1, \frac{1}{1-\beta_1})$

$$\begin{aligned}
& \mathbb{E}^{|\mathcal{F}_{t}} \frac{1}{2} (L_{0} + L_{1} \| \nabla f(\boldsymbol{w}_{t}) \|) (\| \boldsymbol{u}_{t+1} - \boldsymbol{w}_{t} \| + \| \boldsymbol{u}_{t} - \boldsymbol{w}_{t} \|) \| \boldsymbol{u}_{t+1} - \boldsymbol{u}_{t} \| \\
& = \frac{L_{0} \eta^{2}}{2} \left(\frac{32}{(1 - \beta_{1})^{2}} \left\| \frac{\boldsymbol{g}_{t}}{\sqrt{\boldsymbol{\nu}_{t}}} \right\|^{2} + \frac{4}{(1 - \beta_{1})^{2}} \left\| \frac{\boldsymbol{m}_{t-1}}{\sqrt{\boldsymbol{\nu}_{t-1}}} \right\|^{2} \right) + \frac{3}{16} \eta \frac{\| \boldsymbol{G}_{t} \|^{2}}{\sqrt{\beta_{2} \boldsymbol{\nu}_{t-1}}} \\
& = \frac{\eta^{3}}{1 - \beta_{2}} \operatorname{Poly}(L_{0}, L_{1}, \sigma_{0}, \sigma_{1}, \frac{1}{1 - \beta_{1}}, f(\boldsymbol{w}_{1}) - f^{*}) \mathbb{E}^{|\mathcal{F}_{t}} \left(\sum_{s=1}^{t} \sqrt[8]{\beta_{1}^{t-s}} \left(\frac{1}{\sqrt{\beta_{2} \boldsymbol{\nu}_{s-1}}} - \frac{1}{\sqrt{\boldsymbol{\nu}_{s}}} \right) \right) \\
& = \frac{\eta^{3}}{1 - \beta_{2}} \operatorname{Poly}(L_{0}, L_{1}, \sigma_{0}, \sigma_{1}, \frac{1}{1 - \beta_{1}}, f(\boldsymbol{w}_{1}) - f^{*}) \mathbb{E}^{|\mathcal{F}_{t}} \left(\frac{\| \boldsymbol{G}_{t} \|^{2}}{\sqrt{\beta_{2} \boldsymbol{\nu}_{t-1}}} - \frac{\| \boldsymbol{G}_{t+1} \|^{2}}{\sqrt{\boldsymbol{\nu}_{t}}} \right) \\
& = \frac{\eta^{3}}{1 - \beta_{2}} \operatorname{Poly}(L_{0}, L_{1}, \sigma_{0}, \sigma_{1}, \frac{1}{1 - \beta_{1}}, f(\boldsymbol{w}_{1}) - f^{*}) \mathbb{E}^{|\mathcal{F}_{t}} \left(\frac{\| \boldsymbol{G}_{t} \|^{2}}{\sqrt{\beta_{2} \boldsymbol{\nu}_{t-1}}} - \frac{\| \boldsymbol{G}_{t+1} \|^{2}}{\sqrt{\boldsymbol{\nu}_{t}}} \right) \\
& = \frac{\eta^{3}}{1 - \beta_{2}} \operatorname{Poly}(L_{0}, L_{1}, \sigma_{0}, \sigma_{1}, \frac{1}{1 - \beta_{1}}, f(\boldsymbol{w}_{1}) - f^{*}) \mathbb{E}^{|\mathcal{F}_{t}} \left(\frac{\| \boldsymbol{G}_{t} \|^{2}}{\sqrt{\beta_{2} \boldsymbol{\nu}_{t-1}}} - \frac{\| \boldsymbol{G}_{t+1} \|^{2}}{\sqrt{\boldsymbol{\nu}_{t}}} \right) \\
& = \frac{\eta^{3}}{1 - \beta_{2}} \operatorname{Poly}(L_{0}, L_{1}, \sigma_{0}, \sigma_{1}, \frac{1}{1 - \beta_{1}}, f(\boldsymbol{w}_{1}) - f^{*}) \mathbb{E}^{|\mathcal{F}_{t}} \left(\frac{\| \boldsymbol{G}_{t} \|^{2}}{\sqrt{\beta_{2} \boldsymbol{\nu}_{t-1}}} - \frac{\| \boldsymbol{G}_{t+1} \|^{2}}{\sqrt{\boldsymbol{\nu}_{t}}} \right) \\
& = \frac{\eta^{3}}{1 - \beta_{2}} \operatorname{Poly}(L_{0}, L_{1}, \sigma_{0}, \sigma_{1}, \frac{1}{1 - \beta_{1}}, f(\boldsymbol{w}_{1}) - f^{*}) \mathbb{E}^{|\mathcal{F}_{t}} \left(\frac{\| \boldsymbol{G}_{t} \|^{2}}{\sqrt{\beta_{2} \boldsymbol{\nu}_{t-1}}} - \frac{\| \boldsymbol{G}_{t+1} \|^{2}}{\sqrt{\boldsymbol{\nu}_{t}}} \right) \\
& = \frac{\theta^{3}}{1 - \beta_{2}} \operatorname{Poly}(L_{0}, L_{1}, \sigma_{1}) \frac{1}{1 - \beta_{1}} \left(\sum_{s=1}^{t} \frac{\vartheta^{3}}{\sqrt{\beta_{1}^{t-s}} \| \boldsymbol{g}_{s} \|^{2} \| \boldsymbol{G}_{s} \|^{2}} \right) \\
& = \frac{\theta^{3}}{1 - \beta_{2}} \operatorname{Poly}(L_{0}, L_{1}, \sigma_{1}) \frac{\theta^{3}}{1 - \beta_{1}} \left(\sum_{s=1}^{t} \frac{\vartheta^{3}}{\sqrt{\beta_{1}^{t-s}} \| \boldsymbol{g}_{s} \|^{2} \| \boldsymbol{G}_{s} \|^{2}} \right) \\
& = \frac{\theta^{3}}{1 - \beta_{2}} \operatorname{Poly}(L_{0}, L_{1}, \sigma_{1}) \frac{\theta^{3}}{1 - \beta_{1}} \left(\sum_{s=1}^{t} \frac{\vartheta^{3}}{\sqrt{\beta_{1}^{t-s}} \| \boldsymbol{g}_{s} \|^{2} \| \boldsymbol{G}_{s} \|^{2}} \right) \\
& = \frac{\theta^{3}}{1 - \beta_{1}} \left(\sum_{s$$

Here in the second inequality we use $\beta_2 \geq \beta_1$, and in the last inequality we use $\frac{\eta}{\sqrt{1-\beta_2}}$ $\frac{\sqrt{1 - \frac{\beta_1^2}{\beta_2} (1 - \frac{\beta_1}{\sqrt{\beta_2}})^2}}{1024\sigma_1^2 (L_1 + L_0)(1 - \beta_1)}$

Applying the estimations of the first-order term and the second-order term back into the descent lemma, we derive that

$$\mathbb{E}^{|\mathcal{F}_{t}}f(\boldsymbol{u}_{t+1}) \leq f(\boldsymbol{u}_{t}) - \frac{3}{16}\eta \mathbb{E}^{|\mathcal{F}_{t}} \frac{\|\boldsymbol{G}_{t}\|^{2}}{\sqrt{\beta_{2}\boldsymbol{\nu}_{t-1}}} + \frac{L_{0}\eta^{2}}{2} \left(\frac{32}{(1-\beta_{1})^{2}} \left\|\frac{\boldsymbol{g}_{t}}{\sqrt{\boldsymbol{\nu}_{t}}}\right\|^{2} + \frac{4}{(1-\beta_{1})^{2}} \left\|\frac{\boldsymbol{m}_{t-1}}{\sqrt{\boldsymbol{\nu}_{t-1}}}\right\|^{2}\right) \\ + \frac{\eta^{3}}{1-\beta_{2}} \operatorname{Poly}(L_{0}, L_{1}, \sigma_{0}, \sigma_{1}, \frac{1}{1-\beta_{1}}, f(\boldsymbol{w}_{1}) - f^{*}) \mathbb{E}^{|\mathcal{F}_{t}}\left(\sum_{s=1}^{t} \sqrt[s]{\beta_{1}^{t-s}} \left(\frac{1}{\sqrt{\beta_{2}\boldsymbol{\nu}_{s-1}}} - \frac{1}{\sqrt{\boldsymbol{\nu}_{s}}}\right)\right)$$

$$+\left(\frac{\eta^{3}}{1-\beta_{2}}+\eta\right)\operatorname{Poly}(L_{0},L_{1},\sigma_{0},\sigma_{1},\frac{1}{1-\beta_{1}},f(\boldsymbol{w}_{1})-f^{*})\mathbb{E}^{|\mathcal{F}_{t}}\left(\frac{\|\boldsymbol{G}_{t}\|^{2}}{\sqrt{\beta_{2}\boldsymbol{\nu}_{t-1}}}-\frac{\|\boldsymbol{G}_{t+1}\|^{2}}{\sqrt{\boldsymbol{\nu}_{t}}}\right)$$

+ $\eta \operatorname{Poly}(L_0, L_1, \sigma_0, \sigma_1, \frac{1}{1 - \beta_1}, f(\boldsymbol{w}_1) - f^*) \mathbb{E}^{|\mathcal{F}_t} \left(\frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\beta_2 \boldsymbol{\nu}_{t-1}}} - \frac{\|\boldsymbol{G}_{t+1}\|^2}{\sqrt{\boldsymbol{\nu}_t}} \right)$

1995
1996
$$64\eta\sigma_1^2 L_1 = 32L_1\eta^2 \quad \left(\sum_{s=1}^t \sqrt[8]{\beta_1^{t-s}} \|g_s\|^2 \|G_s\|^2\right)$$

Constructing stopping time as $\tau \triangleq \min\{t : \|\nabla f(\boldsymbol{w}_{t+1})\| \le \frac{\sqrt[4]{L_0\sigma_0^2(f(\boldsymbol{w}_1) - f^*)}}{\sqrt[4]{T}}\} \land T$. Then, denote

$$\begin{aligned} & 2001 \\ & 2002 \\ & x_t \triangleq f(\boldsymbol{u}_t) - f(\boldsymbol{u}_{t+1}) - \frac{3}{16} \eta \mathbb{E} \frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\beta_2 \boldsymbol{\nu}_{t-1}}} + \frac{L_0 \eta^2}{2} \left(\frac{32}{(1-\beta_1)^2} \left\| \frac{\boldsymbol{g}_t}{\sqrt{\boldsymbol{\nu}_t}} \right\|^2 + \frac{4}{(1-\beta_1)^2} \left\| \frac{\boldsymbol{m}_{t-1}}{\sqrt{\boldsymbol{\nu}_{t-1}}} \right\|^2 \right) \\ & 2003 \\ & + \frac{\eta^3}{1-\beta_2} \operatorname{Poly}(L_0, L_1, \sigma_0, \sigma_1, \frac{1}{1-\beta_1}, f(\boldsymbol{w}_1) - f^*) \left(\sum_{s=1}^t \sqrt[8]{\beta_1^{t-s}} \left(\frac{1}{\sqrt{\beta_2 \boldsymbol{\nu}_{s-1}}} - \frac{1}{\sqrt{\boldsymbol{\nu}_s}} \right) \right) \right) \\ & 2006 \\ & 2006 \\ & + \left(\frac{\eta^3}{1-\beta_2} + \eta \right) \operatorname{Poly}(L_0, L_1, \sigma_0, \sigma_1, \frac{1}{1-\beta_1}, f(\boldsymbol{w}_1) - f^*) \left(\frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\beta_2 \boldsymbol{\nu}_{t-1}}} - \frac{\|\boldsymbol{G}_{t+1}\|^2}{\sqrt{\boldsymbol{\nu}_t}} \right) \right) \\ & + \frac{64\eta \sigma_1^2 L_1}{\sqrt{(1-\beta_1)^3}} \frac{32L_1 \eta^2}{(1-\beta_1)^{\frac{1}{2}}} \left(\sum_{s=1}^t \frac{\sqrt[8]{\beta_1^{t-s}}}{\sqrt{\beta_2 \boldsymbol{\nu}_{s-1}} \sqrt{\boldsymbol{\nu}_s^2}} \right), \end{aligned}$$

2013

and due to Eq. (26), we have $\mathbb{E}^{|\mathcal{F}_t} x_t \ge 0$, and thus $S_t \triangleq \sum_{s=1}^t x_s$ ($S_0 = 0$) is a submartingale with respect to $\{\mathcal{F}_t\}_t$. Also, as τ is a bounding stopping theorem, by optional stopping time, we obtain that $\mathbb{E}S_{\tau} \ge 0$, which gives

$$\begin{aligned} \frac{3}{16} \eta \mathbb{E} \sum_{t=1}^{\tau} \frac{\|\boldsymbol{G}_{t}\|^{2}}{\sqrt{\beta_{2}\nu_{t-1}}} \leq f(\boldsymbol{u}_{1}) - \mathbb{E}f(\boldsymbol{u}_{\tau+1}) + \frac{L_{0}\eta^{2}}{2} \mathbb{E} \sum_{t=1}^{\tau} \left(\frac{32}{(1-\beta_{1})^{2}} \left\| \frac{\boldsymbol{g}_{t}}{\sqrt{\nu_{t}}} \right\|^{2} + \frac{4}{(1-\beta_{1})^{2}} \left\| \frac{\boldsymbol{m}_{t-1}}{\sqrt{\nu_{t-1}}} \right\|^{2} \right) \\ + \frac{\eta^{3}}{1-\beta_{2}} \operatorname{Poly}(L_{0}, L_{1}, \sigma_{0}, \sigma_{1}, \frac{1}{1-\beta_{1}}, f(\boldsymbol{w}_{1}) - f^{*}) \mathbb{E} \sum_{t=1}^{\tau} \left(\frac{1}{\sqrt{\beta_{1}^{L-s}}} \left(\frac{1}{\sqrt{\beta_{2}\nu_{t-1}}} - \frac{1}{\sqrt{\nu_{s}}} \right) \right) \\ + \mathbb{E} \sum_{t=1}^{\tau} \left(\frac{\eta^{3}}{1-\beta_{2}} + \eta \right) \operatorname{Poly}(L_{0}, L_{1}, \sigma_{0}, \sigma_{1}, \frac{1}{1-\beta_{1}}, f(\boldsymbol{w}_{1}) - f^{*}) \left(\frac{\|\boldsymbol{G}_{t}\|^{2}}{\sqrt{\beta_{2}\nu_{t-1}}} - \frac{\|\boldsymbol{G}_{t+1}\|^{2}}{\sqrt{\nu_{t}}} \right) \\ + \mathbb{E} \sum_{t=1}^{\tau} \left(\frac{64\eta\sigma_{1}^{2}L_{1}}{\sqrt{(1-\beta_{1})^{3}}} \frac{32L_{1}\eta^{2}}{(1-\beta_{1})^{\frac{1}{2}}} \left(\sum_{s=1}^{s} \frac{\sqrt{\beta_{1}^{L-s}}\|\boldsymbol{g}_{s}\|^{2}}{\sqrt{\beta_{2}\nu_{s-1}}\sqrt{\nu_{s}^{2}}} \right) \\ \leq f(\boldsymbol{u}_{1}) - \mathbb{E}f(\boldsymbol{u}_{\tau+1}) + \frac{L_{0}\eta^{2}}{2} \mathbb{E} \sum_{t=1}^{\tau} \left(\frac{32}{(1-\beta_{1})^{2}} \right) \frac{\boldsymbol{g}_{t}}{\sqrt{\mu_{t}}} \right\|^{2} + \frac{4}{(1-\beta_{1})^{2}} \left\| \frac{\boldsymbol{m}_{t-1}}{\sqrt{\nu_{t-1}}} \right\|^{2} \right) \\ + \frac{\eta^{3}}{1-\beta_{2}} \operatorname{Poly}(L_{0}, L_{1}, \sigma_{0}, \sigma_{1}, \frac{1}{1-\beta_{1}}, f(\boldsymbol{w}_{1}) - f^{*}) \mathbb{E} \sum_{t=1}^{\tau} \left(\frac{1}{\sqrt{\beta_{2}\nu_{t-1}}} - \frac{1}{\sqrt{\nu_{t}}} \right) \\ + \mathbb{E} \sum_{t=1}^{\tau} \left(\frac{\eta^{3}}{1-\beta_{2}} + \eta \right) \operatorname{Poly}(L_{0}, L_{1}, \sigma_{0}, \sigma_{1}, \frac{1}{1-\beta_{1}}, f(\boldsymbol{w}_{1}) - f^{*}) \mathbb{E} \sum_{t=1}^{\tau} \left(\frac{1}{\sqrt{\beta_{2}\nu_{t-1}}} - \frac{\|\boldsymbol{G}_{t+1}\|^{2}}{\sqrt{\nu_{t}}} \right) \\ + \mathbb{E} \sum_{t=1}^{\tau} \frac{64\eta\sigma_{1}^{2}L_{1}}{\sqrt{(1-\beta_{1})^{3}}} \frac{256L_{1}\eta^{2}}{(1-\beta_{1})^{\frac{3}{2}}} \left(\frac{|\boldsymbol{g}_{t}|^{2}||^{2}|\boldsymbol{g}_{t}||^{2}}{\sqrt{\nu_{t}}} \right) \\ + \mathbb{E} \sum_{t=1}^{\tau} \frac{64\eta\sigma_{1}^{2}L_{1}}{\sqrt{(1-\beta_{1})^{3}}} \frac{256L_{1}\eta^{2}}{(1-\beta_{1})^{\frac{3}{2}}} \mathbb{E} \sum_{t=1}^{\tau} \left(\frac{32}{(1-\beta_{1})^{2}} \left\| \frac{\boldsymbol{g}_{t}}{\boldsymbol{y}_{t}} \right\|^{2} + \frac{4}{(1-\beta_{1})^{2}} \left\| \frac{\boldsymbol{m}_{t-1}}{\sqrt{\nu_{t}}} \right\|^{2} \right) \\ + \mathbb{E} \sum_{t=1}^{\tau} \frac{64\eta\sigma_{1}^{2}L_{1}}{\sqrt{(1-\beta_{1})^{3}}} \frac{256L_{1}\eta^{2}}{(1-\beta_{1})^{\frac{3}{2}}} \left(\frac{|\boldsymbol{g}_{t}|^{2}|^{2}}{(1-\beta_{1})^{2}} \right) \\ \leq f(\boldsymbol{u}_{1}) - \mathbb{E} f(\boldsymbol{u}_{1}) - \mathbb{E} f(\boldsymbol{u}_{1}) - \frac{1}{\sqrt{\nu_{t}}} \right) \\ + \frac{1}{\eta} - \frac{\beta_{2}}{\beta_{2}} \operatorname{Poly}(L_{0}, L_{1}, \sigma_{0}, \sigma_{1}, \frac{1}{(1-\beta_{1})^{2}} \right) \left(\frac{$$

where the last inequality is because due to $\frac{\eta}{\sqrt{1-\beta_2}} \leq \frac{1-\beta_1}{128L_1\sigma_1}$, following the similar reasoning of Eq. (25), we have

By rearranging the inequality and due to Lemma 3, we obtain

 $+\frac{1}{16}\eta \mathbb{E}\sum_{t=1}^{\tau}\frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\beta_2\boldsymbol{\nu}_{t-1}}}.$

 $\tau \parallel \mathbf{C} \parallel^2$

 $\frac{64\eta\sigma_1^2 L_1}{\sqrt{(1-\beta_1)^3}} \frac{256L_1\eta^2}{(1-\beta_1)^{\frac{3}{2}}} \left(\frac{\|\boldsymbol{g}_t\|^2 \|\boldsymbol{G}_t\|^2}{\sqrt{\beta_2\boldsymbol{\nu}_{t-1}}\sqrt{\boldsymbol{\nu}_t^2}}\right)$

 $\leq \frac{128\eta\sigma_1^2 L_1}{\sqrt{(1-\beta_1)^3}} \frac{256L_1\eta^2}{(1-\beta_2)(1-\beta_1)^{\frac{3}{2}}} \left(\frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\beta_2\boldsymbol{\nu}_{t-1}}} - \frac{\|\boldsymbol{G}_{t+1}\|^2}{\sqrt{\boldsymbol{\nu}_t}}\right)$

$$\frac{1}{8}\eta \mathbb{E} \sum_{t=1}^{7} \frac{\|\mathbf{G}_{t}\|}{\sqrt{\beta_{2}\nu_{t-1}}} = \frac{1}{8}\eta \mathbb{E} \sum_{t=1}^{7} \frac{1}{\sqrt{\beta_{2}\nu_{t-1}}} = \frac{1}{8}\eta \mathbb{E} \sum_{t=1}^{7} \frac{1}{1-\beta_{1}} \int_{t=1}^{7} \frac{1}{1-\beta$$

+ $\frac{\eta^3}{1-\beta_2}$ Poly $(L_0, L_1, \sigma_0, \sigma_1, \frac{1}{1-\beta_1}, f(\boldsymbol{w}_1) - f^*) \sum_{i=1}^t \sqrt[8]{\beta_1^{t-s}} \left(\frac{1}{\sqrt{\beta_2 \boldsymbol{\nu}_{t-1}}} - \frac{1}{\sqrt{\boldsymbol{\nu}_t}} \right)$

Furthermore, as $T \ge \operatorname{Poly}(L_0, L_1, \sigma_0, \sigma_1, \frac{1}{1-\beta_1}, f(\boldsymbol{w}_1) - f^*), \eta^3 = \frac{\eta}{T} \operatorname{Poly}(L_0, L_1, \sigma_0, \sigma_1, \frac{1}{1-\beta_1}, f(\boldsymbol{w}_1) - f^*), \eta(1-\beta_2) = \frac{\eta}{T} \operatorname{Poly}(L_0, L_1, \sigma_0, \sigma_1, \frac{1}{1-\beta_1}, f(\boldsymbol{w}_1) - f^*),$ and $\|\boldsymbol{G}_t\| \ge \frac{\operatorname{Poly}(L_0, L_1, \sigma_0, \sigma_1, \frac{1}{1-\beta_1}, f(\boldsymbol{w}_1) - f^*)}{\sqrt[4]{T}}$, we have

$$\eta^{3} \operatorname{Poly}(L_{0}, L_{1}, \sigma_{0}, \sigma_{1}, \frac{1}{1 - \beta_{1}}, f(\boldsymbol{w}_{1}) - f^{*}) \mathbb{E} \sum_{t=1}^{\tau-1} \frac{1}{\sqrt{\beta_{2}\boldsymbol{\nu}_{t}}} \leq \frac{1}{32} \eta \mathbb{E} \sum_{t=1}^{\tau} \frac{\|\boldsymbol{G}_{t}\|^{2}}{\sqrt{\beta_{2}\boldsymbol{\nu}_{t-1}}},$$
$$\mathbb{E} \sum_{t=1}^{\tau} \left(\eta^{3} + \eta(1 - \beta_{2})\right) \operatorname{Poly}(L_{0}, L_{1}, \sigma_{0}, \sigma_{1}, \frac{1}{1 - \beta_{1}}, f(\boldsymbol{w}_{1}) - f^{*}) \left(\frac{\|\boldsymbol{G}_{t}\|^{2}}{\sqrt{\beta_{2}\boldsymbol{\nu}_{t-1}}}\right) \leq \frac{1}{32} \eta \mathbb{E} \sum_{t=1}^{\tau} \frac{\|\boldsymbol{G}_{t}\|^{2}}{\sqrt{\beta_{2}\boldsymbol{\nu}_{t-1}}}$$

which thus leads to

$$\frac{1}{16} \eta \mathbb{E} \sum_{t=1}^{\tau} \frac{\|\boldsymbol{G}_{t}\|^{2}}{\sqrt{\beta_{2}\boldsymbol{\nu}_{t-1}}} \leq f(\boldsymbol{u}_{1}) - \mathbb{E}f(\boldsymbol{u}_{\tau+1}) + \frac{64L_{0}\eta^{2}}{(1-\beta_{2})(1-\beta_{1})^{2}} \mathbb{E} \left(\ln \frac{\boldsymbol{\nu}_{\tau}}{\boldsymbol{\nu}_{0}} - \tau \ln \beta_{2} \right) \\
+ \frac{\eta^{3}}{1-\beta_{2}} \operatorname{Poly}(L_{0}, L_{1}, \sigma_{0}, \sigma_{1}, \frac{1}{1-\beta_{1}}, f(\boldsymbol{w}_{1}) - f^{*}) \frac{1}{\sqrt{\beta_{2}\boldsymbol{\nu}_{0}}} + \left(\frac{\eta^{3}}{1-\beta_{2}} + \eta \right) \operatorname{Poly}(L_{0}, L_{1}, \sigma_{0}, \sigma_{1}, \frac{1}{1-\beta_{1}}, f(\boldsymbol{w}_{1}) - f^{*}) \left(\frac{\|\boldsymbol{G}_{1}\|^{2}}{\sqrt{\beta_{2}\boldsymbol{\nu}_{0}}} \right). \tag{27}$$

Similar to the proof of Theorem 3, we then transfer the above bound to the bound of $\sum_{t=1}^{\tau} \|G_t\|$ by two rounds of divide-and-conquer. In the first round, we will bound $\mathbb{E} \ln \nu_{\tau}$. To start with, we have that

2104
2105
$$\frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\beta_2 \boldsymbol{\nu}_{t-1}}} \mathbb{1}_{\|\boldsymbol{G}_t\| \ge \frac{\sigma_0}{\sigma_1}} \ge \frac{\frac{1}{2\sigma_1^2} \mathbb{E}^{|\mathcal{F}_t|} \|\boldsymbol{g}_t\|^2}{\sqrt{\beta_2 \boldsymbol{\nu}_{t-1}}} \mathbb{1}_{\|\boldsymbol{G}_t\| \ge \frac{\sigma_0}{\sigma_1}}$$

Furthermore, we have $\mathbb{E} \frac{\frac{\nu_0}{1-\beta_2}}{\sqrt{\sum_{t=1}^{\tau} \|\boldsymbol{g}_t\|^2 + \frac{\nu_0}{1-\beta_2}}} + \mathbb{E} \sum_{t=1}^{\tau} \frac{\|\boldsymbol{g}_t\|^2}{\sqrt{\sum_{t=1}^{\tau} \|\boldsymbol{g}_t\|^2 + \frac{\nu_0}{1-\beta_2}}} \mathbb{1}_{\|\boldsymbol{G}_t\| < \frac{\sigma_0}{\sigma_1}}$ $\leq \mathbb{E} \frac{\frac{\nu_0}{1-\beta_2}}{\sqrt{\sum_{t=1}^{\tau} \|\boldsymbol{g}_t\|^2 \mathbb{1}_{\|\boldsymbol{G}_t\| < \frac{\sigma_0}{\sigma_1}} + \frac{\nu_0}{1-\beta_2}}} + \mathbb{E} \sum_{t=1}^{\tau} \frac{\|\boldsymbol{g}_t\|^2}{\sqrt{\sum_{t=1}^{\tau} \|\boldsymbol{g}_t\|^2 \mathbb{1}_{\|\boldsymbol{G}_t\| < \frac{\sigma_0}{\sigma_1}} + \frac{\nu_0}{1-\beta_2}}} \mathbb{1}_{\|\boldsymbol{G}_t\| < \frac{\sigma_0}{\sigma_1}}$ $= \mathbb{E}_{\sqrt{\sum_{t=1}^{\tau} \|\boldsymbol{g}_t\|^2 \mathbb{1}_{\|\boldsymbol{G}_t\| < \frac{\sigma_0}{\sigma_1}} + \frac{\boldsymbol{\nu}_0}{1 - \beta_2}} \leq \sqrt{\mathbb{E}_{t=1}^{\tau} \|\boldsymbol{g}_t\|^2 \mathbb{1}_{\|\boldsymbol{G}_t\| < \frac{\sigma_0}{\sigma_1}} + \frac{\boldsymbol{\nu}_0}{1 - \beta_2}}$ $\stackrel{(\star)}{\leq} \sqrt{2\sigma_0^2 \mathbb{E}\tau + \frac{\boldsymbol{\nu}_0}{1-\beta_2}} \leq \sqrt{2\sigma_0^2 T + \frac{\boldsymbol{\nu}_0}{1-\beta_2}},$

where inequality (*) is due to $E^{|\mathcal{F}_t|} \|\boldsymbol{g}_t\|^2 \mathbb{1}_{\|\boldsymbol{G}_t\| < \frac{\sigma_0}{\sigma_1}} \le 2\sigma_0^2$ and optimal stopping theorem.

Conclusively, we obtain

$$\begin{split} \mathbb{E} \sqrt{\sum_{t=1}^{\tau} \|\boldsymbol{g}_{t}\|^{2} + \frac{\boldsymbol{\nu}_{0}}{1-\beta_{2}}} \\ &= \left(\mathbb{E} \frac{\frac{\boldsymbol{\nu}_{0}}{1-\beta_{2}}}{\sqrt{\sum_{t=1}^{\tau} \|\boldsymbol{g}_{t}\|^{2} + \frac{\boldsymbol{\nu}_{0}}{1-\beta_{2}}}} + \mathbb{E} \sum_{t=1}^{\tau} \frac{\|\boldsymbol{g}_{t}\|^{2}}{\sqrt{\sum_{t=1}^{\tau} \|\boldsymbol{g}_{t}\|^{2} + \frac{\boldsymbol{\nu}_{0}}{1-\beta_{2}}}} \mathbb{1}_{\|\boldsymbol{G}_{t}\| < \frac{\sigma_{0}}{\sigma_{1}}} \right) \\ &+ \mathbb{E} \sum_{t=1}^{\tau} \frac{\|\boldsymbol{g}_{t}\|^{2}}{\sqrt{\sum_{t=1}^{\tau} \|\boldsymbol{g}_{t}\|^{2} + \frac{\boldsymbol{\nu}_{0}}{1-\beta_{2}}}} \mathbb{1}_{\|\boldsymbol{G}_{t}\| \geq \frac{\sigma_{0}}{\sigma_{1}}} \right) \\ &\leq \sqrt{\frac{\boldsymbol{\nu}_{0}}{1-\beta_{2}} + 2\sigma_{0}^{2}T} + 2\sqrt{1-\beta_{2}} \mathbb{E} \sum_{t=1}^{\tau} \frac{\|\boldsymbol{g}_{t}\|^{2}}{\sqrt{\beta_{2}\boldsymbol{\nu}_{t-1}}} \mathbb{1}_{\|\boldsymbol{G}_{t}\| \geq \frac{\sigma_{0}}{\sigma_{1}}} \\ &\stackrel{(o)}{\leq} \sqrt{\frac{\boldsymbol{\nu}_{0}}{1-\beta_{2}} + 2\sigma_{0}^{2}T} + 2\sigma_{1}^{2}\sqrt{1-\beta_{2}} \mathbb{E} \sum_{t=1}^{\tau} \frac{\|\boldsymbol{G}_{t}\|^{2}}{\sqrt{\beta_{2}\boldsymbol{\nu}_{t-1}}}, \end{split}$$

2139 where inequality (\circ) is due to optimal stopping theorem.

Then by substituting $\mathbb{E} \sum_{t=1}^{\tau} \frac{\|G_t\|^2}{\sqrt{\beta_2 \nu_{t-1}}}$ we obtain that

$$\begin{aligned} & = 142 \\ & = 144 \\ & = 144 \\ & = \sqrt{\sum_{t=1}^{\tau} \|g_t\|^2 + \frac{\nu_0}{1-\beta_2}} \\ & = \sqrt{\frac{\nu_0}{1-\beta_2} + 2\sigma_0^2 T} + \frac{32\sigma_1^2\sqrt{1-\beta_2}}{\eta} \frac{1}{16}\eta \mathbb{E}\sum_{t=1}^{\tau} \frac{\|G_t\|^2}{\sqrt{\beta_2\nu_{t-1}}} \\ & = 146 \\ & = \sqrt{\frac{\nu_0}{1-\beta_2} + 2\sigma_0^2 T} + \frac{32\sigma_1^2\sqrt{1-\beta_2}\sigma_1^2}{\eta} \left(f(\mathbf{u}_1) - f^* + \frac{64L_0\eta^2}{(1-\beta_2)(1-\beta_1)^2} \mathbb{E}\left(\ln\frac{\nu_\tau}{\nu_0} - T\ln\beta_2\right) \right) \\ & = 148 \\ & = \sqrt{\frac{\nu_0}{1-\beta_2} + 2\sigma_0^2 T} + \frac{32\sqrt{1-\beta_2}\sigma_1^2}{\eta} \left(f(\mathbf{u}_1) - f^* + \frac{64L_0\eta^2}{(1-\beta_2)(1-\beta_1)^2} \mathbb{E}\left(\ln\frac{\nu_\tau}{\nu_0} - T\ln\beta_2\right) \right) \\ & = 148 \\ & = \sqrt{\frac{\nu_0}{1-\beta_2} + 2\sigma_0^2 T} + \frac{32\sqrt{1-\beta_2}\sigma_1^2}{\eta} \left(f(\mathbf{w}_1) - f^* \right) \frac{1}{\sqrt{\beta_2\nu_0}} \left(\frac{\eta^3}{1-\beta_2} + \eta \right) \operatorname{Poly}(L_0, L_1, \sigma_0, \sigma_1, \frac{1}{1-\beta_1}, f(\mathbf{w}_1) - f^*) \left(\frac{\|G_1\|^2}{\sqrt{\beta_2\nu_0}} \right) \right) \\ & = 151 \\ & = \sqrt{\frac{\nu_0}{1-\beta_2} + 2\sigma_0^2 T} + \frac{32\sqrt{1-\beta_2}\sigma_1^2}{\eta} \left(f(\mathbf{w}_1) - f^* + \frac{128L_0\eta^2}{(1-\beta_2)(1-\beta_1)^2} \left(\ln\frac{\mathbb{E}\sqrt{1-\beta_2}\sqrt{\sum_{t=1}^{\tau} \|g_t\|^2 + \frac{\nu_0}{1-\beta_2}}}{\sqrt{\nu_0}} - T\ln\beta_2 \right) \right) \\ & = 151 \\ & = \sqrt{\frac{1}{1-\beta_2}} \operatorname{Poly}(L_0, L_1, \sigma_0, \sigma_1, \frac{1}{1-\beta_1}, f(\mathbf{w}_1) - f^*) \frac{1}{\sqrt{\beta_2\nu_0}}} + \left(\frac{\eta^3}{1-\beta_2} + \eta \right) \operatorname{Poly}(L_0, L_1, \sigma_0, \sigma_1, \frac{1}{1-\beta_1}, f(\mathbf{w}_1) - f^*) \left(\frac{\|G_1\|^2}{\sqrt{\beta_2\nu_0}} \right) \right) \\ & = 156 \\ & \text{Multiplying both sides of the above inequality by $\sqrt{1-\beta_2} \text{ then gives} \end{aligned}$$$

2158
2159
$$\sqrt{1-\beta_2}\mathbb{E}\sqrt{\sum_{t=1}^{\tau}\|\boldsymbol{g}_t\|^2 + \frac{\boldsymbol{\nu}_0}{1-\beta_2}} \le 3\sqrt{\boldsymbol{\nu}_0 + 2\sigma_0^2 T(1-\beta_2)} + \frac{1}{4}\ln\mathbb{E}\sqrt{1-\beta_2}\sqrt{\sum_{t=1}^{\tau}\|\boldsymbol{g}_t\|^2 + \frac{\boldsymbol{\nu}_0}{1-\beta_2}}$$

2160
2160
2161 where we use
$$\eta = \frac{\operatorname{Poly}(L_0, L_1, \sigma_0, \sigma_1, \frac{1}{1-\beta_1}, f(\boldsymbol{w}_1) - f^*)}{\sqrt{T}}, 1 - \beta_2 = \frac{\operatorname{Poly}(L_0, L_1, \sigma_0, \sigma_1, \frac{1}{1-\beta_1}, f(\boldsymbol{w}_1) - f^*)}{T},$$

2162 and $T \ge \operatorname{Poly}(L_0, L_1, \sigma_0, \sigma_1, \frac{1}{1-\beta_1}, f(\boldsymbol{w}_1) - f^*).$ Therefore, we obtain
2163 $\sqrt{1-\beta_2}\mathbb{E}\sqrt{\sum_{t=1}^{\tau} \|\boldsymbol{g}_t\|^2 + \frac{\boldsymbol{\nu}_0}{1-\beta_2}} \le 6\sqrt{\boldsymbol{\nu}_0 + 2\sigma_0^2 T(1-\beta_2)}.$

Therefore, Eq. (27) can be rewritten as

We then execute the second round of divide-and-conquer. To begin with, we have that

$$\mathbb{E}\sum_{t=1}^{\tau}\left[\frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\beta_2\boldsymbol{\nu}_{t-1}}}\mathbb{1}_{\|\boldsymbol{G}_t\|\geq\frac{\sigma_0}{\sigma_1}}\right] \leq \mathbb{E}\sum_{t=1}^{\tau}\left[\frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\beta_2\boldsymbol{\nu}_{t-1}}}\right].$$

On the other hand, we have that

$$\frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\beta_2 \boldsymbol{\nu}_{t-1}}} \mathbb{1}_{\|\boldsymbol{G}_t\| \ge \frac{\sigma_0}{\sigma_1}} \ge \frac{1}{2\sigma_1^2} \mathbb{E}^{|\mathcal{F}_t} \frac{\|\boldsymbol{g}_t\|^2}{\sqrt{\boldsymbol{\nu}_t} + \sqrt{\beta_2 \boldsymbol{\nu}_{t-1}}} \mathbb{1}_{\|\boldsymbol{G}_t\| \ge \frac{\sigma_0}{\sigma_1}} \ge \frac{1}{2(1-\beta_2)\sigma_1^2} \mathbb{E}^{|\mathcal{F}_t} \left[\left(\sqrt{\boldsymbol{\nu}_t} - \sqrt{\beta_2 \boldsymbol{\nu}_{t-1}}\right) \mathbb{1}_{\|\boldsymbol{G}_t\| \ge \frac{\sigma_0}{\sigma_1}} \right]$$

Meanwhile, recall $\{\bar{\boldsymbol{\nu}}_t\}_{t=0}^{\infty}$ as $\bar{\boldsymbol{\nu}}_0 = \boldsymbol{\nu}_0$, $\bar{\boldsymbol{\nu}}_t = \beta_2 \bar{\boldsymbol{\nu}}_{t-1} + (1-\beta_2)|g_t|^2 \mathbb{1}_{\|\boldsymbol{G}_t\| < \frac{\sigma_0^2}{\sigma_1^2}}$ and $\bar{\boldsymbol{\nu}}_t \leq \boldsymbol{\nu}_t$. Therefore

$$\begin{split} & \sum_{t=1}^{2193} \qquad \mathbb{E}\sum_{t=1}^{\tau} \left[\left(\sqrt{\nu_t} - \sqrt{\beta_2 \nu_{t-1}} \right) \mathbb{1}_{\|G_t\| < \frac{\sigma_t^2}{\sigma_t^2}} \right] \\ & = \mathbb{E}\sum_{t=1}^{\tau} \left(\sqrt{\beta_2 \nu_{t-1} + (1 - \beta_2) \|g_t\|^2} - \sqrt{\beta_2 \nu_{t-1}} \right) \mathbb{1}_{\|G_t\| < \frac{\sigma_t^2}{\sigma_t^2}} \\ & = \mathbb{E}\sum_{t=1}^{\tau} \left(\sqrt{\beta_2^2 \bar{\nu}_{t-1} + \beta_2 (1 - \beta_2) \|g_t\|^2} + (1 - \beta_2) \sigma_0^2 - \sqrt{\beta_2 (\beta_2 \bar{\nu}_{t-1} + (1 - \beta_2) \sigma_0^2)} \right) \mathbb{1}_{\|G_t\| < \frac{\sigma_t^2}{\sigma_t^2}} \\ & = \mathbb{E}\sum_{t=1}^{\tau} \left(\sqrt{\beta_2 \bar{\nu}_{t-1} + (1 - \beta_2) \|g_t\|^2} \mathbb{1}_{\|G_t\| < \frac{\sigma_t^2}{\sigma_t^1}} - \sqrt{\beta_2 \bar{\nu}_{t-1}} \right) \\ & = \mathbb{E}\sum_{t=1}^{\tau} \left(\sqrt{\bar{\nu}_t} - \sqrt{\beta_2 \bar{\nu}_{t-1}} \right) \\ & = \mathbb{E}\sqrt{\bar{\nu}_\tau} + (1 - \sqrt{\beta_2}) \sum_{t=1}^{\tau-1} \mathbb{E}\sqrt{\bar{\nu}_t} - \mathbb{E}\sqrt{\beta_2 \nu_0} \\ & = \mathbb{E}\sqrt{\bar{\nu}_\tau} + (1 - \sqrt{\beta_2}) \sum_{t=1}^{\tau} \mathbb{E}\sqrt{\bar{\nu}_t} - \mathbb{E}\sqrt{\beta_2 \nu_0} \\ & \leq \mathbb{E}\sqrt{\bar{\nu}_\tau} + (1 - \sqrt{\beta_2}) T \sigma_0 - \mathbb{E}\sqrt{\beta_2 \nu_0}. \end{split}$$

All in all, summing the above two inequalities together, we obtain that

$$\mathbb{E}\sqrt{\boldsymbol{\nu}_{\tau}} + (1-\sqrt{\beta_2})\mathbb{E}\sum_{t=1}^{\tau-1}\sqrt{\boldsymbol{\nu}_t} - \sqrt{\beta_2\boldsymbol{\nu}_0}$$

$$= \mathbb{E} \sum_{t=1}^{\tau} \left(\sqrt{\boldsymbol{\nu}_t} - \sqrt{\beta_2 \boldsymbol{\nu}_{t-1}} \right)$$
$$\leq \mathbb{E} \sum_{t=1}^{\tau} \left(\sqrt{\boldsymbol{\nu}_t} - \sqrt{\beta_2 \boldsymbol{\nu}_{t-1}} \right) \mathbb{1}_{\|G_t\| \ge \frac{\sigma_0}{\sigma_1}} + \mathbb{E} \sum_{t=1}^{\tau} \left(\sqrt{\boldsymbol{\nu}_t} - \sqrt{\beta_2 \boldsymbol{\nu}_{t-1}} \right) \mathbb{1}_{\|G_t\| < \frac{\sigma_0^2}{\sigma_1^2}}$$

Since $\boldsymbol{\nu}_0 = \bar{\boldsymbol{\nu}}_0$ and $\mathbb{E}\sqrt{\boldsymbol{\nu}_{\tau}} \geq \mathbb{E}\sqrt{\bar{\boldsymbol{\nu}}_{\tau}}$, we obtain

$$(1 - \sqrt{\beta_2}) \mathbb{E} \sum_{t=0}^{\tau-1} \sqrt{\tilde{\nu}_t} \le 2(1 - \beta_2) \sigma_1^2 \mathbb{E} \sum_{t=1}^{\tau} \left[\frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\beta_2 \boldsymbol{\nu}_{t-1}}} \right] + (1 - \sqrt{\beta_2}) T \sigma_0.$$

 $\leq 2(1-\beta_2)\sigma_1^2 \sum_{t=1}^{\tau} \mathbb{E}\left[\frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\beta_2\boldsymbol{\nu}_{t-1}}}\right] + \mathbb{E}\sqrt{\boldsymbol{\nu}_{\tau}} + (1-\sqrt{\beta_2})T\sigma_0 - \mathbb{E}\sqrt{\beta_2\boldsymbol{\nu}_0}.$

2232 Dividing both sides of the above equation by $1 - \sqrt{\beta_2}$ then gives

$$\mathbb{E}\sum_{t=0}^{\tau-1}\sqrt{\widetilde{\boldsymbol{\nu}}_t} \leq 4\sigma_1^2 \mathbb{E}\sum_{t=1}^{\tau}\left[\frac{\|\boldsymbol{G}_t\|^2}{\sqrt{\beta_2\boldsymbol{\nu}_{t-1}}}\right] + T\sigma_0$$

By applying Eq. (21) and the constraint of τ , we obtain that

$$\mathbb{E}\sum_{t=0}^{\tau-1}\sqrt{\nu_{t}} \leq T\sigma_{0} + 64\frac{\sigma_{1}^{2}}{\eta} \left(f(\boldsymbol{u}_{1}) - \mathbb{E}f(\boldsymbol{u}_{\tau+1}) + \frac{64L_{0}\eta^{2}}{1-\beta_{2}} \left(\ln\frac{6\sqrt{\nu_{0} + 2\sigma_{0}^{2}T(1-\beta_{2})}}{\sqrt{\nu_{0}}} - T\ln\beta_{2} \right) \\ + \frac{\eta^{3}}{1-\beta_{2}}\operatorname{Poly}(L_{0}, L_{1}, \sigma_{0}, \sigma_{1}, \frac{1}{1-\beta_{1}}, f(\boldsymbol{w}_{1}) - f^{*})\frac{1}{\sqrt{\beta_{2}\nu_{0}}} \\ + \left(\frac{\eta^{3}}{1-\beta_{2}} + \eta\right)\operatorname{Poly}(L_{0}, L_{1}, \sigma_{0}, \sigma_{1}, \frac{1}{1-\beta_{1}}, f(\boldsymbol{w}_{1}) - f^{*}) \left(\frac{\|\boldsymbol{G}_{1}\|^{2}}{\sqrt{\beta_{2}\nu_{0}}}\right) \right) \\ \leq 4T\sigma_{0}.$$

2247 Here the last inequality is due to $\eta = \frac{\operatorname{Poly}(L_0, L_1, \sigma_0, \sigma_1, \frac{1}{1-\beta_1}, f(\boldsymbol{w}_1) - f^*)}{\sqrt{T}}, \ 1 - \beta_2 =$ 2248 $\frac{\operatorname{Poly}(L_0, L_1, \sigma_0, \sigma_1, \frac{1}{1-\beta_1}, f(\boldsymbol{w}_1) - f^*)}{T}, \text{ and } T \ge \operatorname{Poly}(L_0, L_1, \sigma_0, \sigma_1, \frac{1}{1-\beta_1}, f(\boldsymbol{w}_1) - f^*). \text{ Therefore, by}$ 2250 Cauchy's inequality, we obtain that

$$\begin{aligned} & 2251 \\ & \left(\mathbb{E}\sum_{t=1}^{\tau} \|\nabla f(\boldsymbol{w}_{t})\|\right)^{2} \leq \left(\mathbb{E}\sum_{t=0}^{\tau-1} \sqrt{\nu_{t}}\right) \left(\mathbb{E}\sum_{t=1}^{\tau} \left[\frac{\|\boldsymbol{G}_{t}\|^{2}}{\sqrt{\beta_{2}\nu_{t-1}}}\right]\right) \\ & \leq 4T\sigma_{0} \times \frac{16}{\eta} \left(f(\boldsymbol{w}_{1}) - f^{*} + \frac{64L_{0}\eta^{2}}{1 - \beta_{2}} \left(\ln \frac{6\sqrt{\nu_{0} + 2\sigma_{0}^{2}T(1 - \beta_{2})}}{\sqrt{\nu_{0}}} - T\ln\beta_{2}\right) \\ & + \frac{\eta^{3}}{1 - \beta_{2}} \operatorname{Poly}(L_{0}, L_{1}, \sigma_{0}, \sigma_{1}, \frac{1}{1 - \beta_{1}}, f(\boldsymbol{w}_{1}) - f^{*}) \frac{1}{\sqrt{\beta_{2}\nu_{0}}} \\ & + \left(\frac{\eta^{3}}{1 - \beta_{2}} + \eta\right) \operatorname{Poly}(L_{0}, L_{1}, \sigma_{0}, \sigma_{1}, \frac{1}{1 - \beta_{1}}, f(\boldsymbol{w}_{1}) - f^{*}) \left(\frac{\|\boldsymbol{G}_{1}\|^{2}}{\sqrt{\beta_{2}\nu_{0}}}\right) \right) \\ & 2261 \\ & 2262 \\ & \leq 4T\sigma_{0} \times \frac{16}{\eta} \left(3(f(\boldsymbol{w}_{1}) - f^{*}) + \frac{128L_{0}\eta^{2}}{(1 - \beta_{2})(1 - \beta_{1})^{2}} \left(\ln \frac{6\sqrt{\nu_{0} + 2\sigma_{0}^{2}T(1 - \beta_{2})}}{\sqrt{\nu_{0}}} - T\ln\beta_{2}\right) \right) \\ & \leq 64T\sigma_{0} \left(3\sqrt{TL_{0}(f(\boldsymbol{w}_{1}) - f^{*})} + \frac{128L_{0}\eta}{(1 - \beta_{2})(1 - \beta_{1})^{2}T}T \left(\ln \frac{6\sqrt{\nu_{0} + 2\sigma_{0}^{2}T(1 - \beta_{2})}}{\sqrt{\nu_{0}}}\right) \right) \\ & \leq 64T\sigma_{0} \left(387(1 - \beta_{1})\sqrt{TL_{0}(f(\boldsymbol{w}_{1}) - f^{*})} \right), \end{aligned}$$

2268 2269 where inequality (*) is due to $\eta = \frac{\operatorname{Poly}(L_0, L_1, \sigma_0, \sigma_1, \frac{1}{1-\beta_1}, f(\boldsymbol{w}_1) - f^*)}{\sqrt{T}}, \ 1 - \beta_2 = \frac{\operatorname{Poly}(L_0, L_1, \sigma_0, \sigma_1, \frac{1}{1-\beta_1}, f(\boldsymbol{w}_1) - f^*)}{T}, \ \text{and} \ T \ge \operatorname{Poly}(L_0, L_1, \sigma_0, \sigma_1, \frac{1}{1-\beta_1}, f(\boldsymbol{w}_1) - f^*), \ \text{inequality} \ (\bullet)$ 2271 is due to that $\eta = \frac{\sqrt{f(\boldsymbol{w}_1) - f^*}}{\sqrt{L_0}}, \ \text{and} \ \text{last inequality} \ \text{is due to} \ \frac{1}{(1-\beta_2)T} \ln \frac{6\sqrt{\nu_0 + 2\sigma_0^2 T(1-\beta_2)}}{\sqrt{\nu_0}} \le 6.$

We then consider two cases: $\tau < T$ and $\tau = T$: for the first case, we have that according to the definition of τ $\frac{4}{\sigma^2 T_{-1}(f(at)) - f(at)}$

$$\mathbb{E}\min_{t\in[1,T]} \|\boldsymbol{G}_t\| \mathbb{1}_{\tau < T} \le \frac{\sqrt[4]{\sigma_0^2 L_0(f(\boldsymbol{w}_1) - f^*)}}{\sqrt[4]{T}}.$$

2278 For the latter case, we have

2276 2277

2283 2284 2285

2286

$$\mathbb{E}\min_{t\in[1,T]} \|\boldsymbol{G}_t\| \mathbb{1}_{\tau=T} \le \frac{1}{T} \left(\mathbb{E}\sum_{t=1}^T \|\nabla f(\boldsymbol{w}_t)\| \mathbb{1}_{\tau=T} \right) \le \frac{1}{T} \left(\mathbb{E}\sum_{t=1}^\tau \|\nabla f(\boldsymbol{w}_t)\| \right) \le \frac{256\sqrt[4]{\sigma_0^2 L_0(f(\boldsymbol{w}_1) - f^*)}}{\sqrt[4]{T}}$$

Summing the two inequalities above complete the proof.

D.2 PROOF OF PARAMETER-AGNOSTIC ADAM

2287 D.2.1 RELATED WORKS ON PARAMETER AGNOSTIC OPTIMIZATION

2288 **Parameter-agnostic optimization.** The term "parameter-agnostic" implies that the optimizer is 2289 capable of converging without the need for extensive hyperparameter tuning or detailed knowledge 2290 of the task characteristics. Designing parameter-agnostic or parameter-free optimizers is a significant 2291 challenge, as it can help avoid the extensive cost associated with hyperparameter search. Existing 2292 works on parameter-agnostic optimization can be categorized into several streams based on the 2293 settings they are predicated upon. In the deterministic offline setting, it is widely acknowledged that GD is not parameter-agnostic, even under an L-smooth condition Nesterov et al. (2018). However, 2294 this can be rectified by combining the GD with the Backtracking Line Search technique Armijo 2295 (1966). In the stochastic offline setting, under the L-smooth condition, multiple algorithms have 2296 been shown to be parameter-agnostic Yang et al. (2023); Ward et al. (2020); Faw et al. (2022); Wang 2297 et al. (2023b); Cutkosky & Mehta (2020). More recently, Hübler et al. (2023) demonstrated that 2298 Normalized-SGDM can be parameter-agnostic even under an (L_0, L_1) -smooth condition. In the 2299 realm of online convex optimization, Orabona & Pál (2016); Orabona & Tommasi (2017) have shown 2300 there exist parameter-free algorithms achieving optimal dependence regarding not only the final error 2301 but also other problem specifics. 2302

2303 D.2.2 OUR RESULTS 2304

As we select $\eta = 1/\sqrt{T}$, choosing $1-\beta_2 = \Omega(1/T)$ has the advantage that the update norm decreases with respect to T. This makes Adam parameter-agnostic under the (L_0, L_1) -smooth condition, as the update norm will eventually become smaller than $\frac{1}{L_1}$ as T increases.

Theorem 12. Let Assumptions 1 and 2 hold. Then, at the t-th iteration, setting $\eta = \frac{1}{\sqrt{t}}$, $\beta_2 = 1 - \frac{1}{\sqrt[4]{t^3}}$, we have that Algorithm 1 satisfies

$$\mathbb{E}\min_{t\in[1,T]} \|
abla f(oldsymbol{w}_t)\| \leq ilde{\mathcal{O}}\left(rac{1}{\sqrt[4]{T}}
ight).$$

It is shown in Hübler et al. (2023) that Normed-SGDM is parameter-agnostic. Here we show thatAdam with a specific scheduler can achieve the same goal.

Proof. As described in Section 5, the proof immediately follows by several modifications of the proofof Theorem 11:

2319

2316

- 2320 2321
- We start our analysis for $t \ge \frac{L_1^4 128^4 \sigma_1^4}{(1-\beta_1)^4}$. For $t \le \frac{L_1^4 128^4 \sigma_1^4}{(1-\beta_1)^4}$, the function value can be bounded by constant as $\frac{L_1^4 128^4 \sigma_1^4}{(1-\beta_1)^4}$ is independent of t;

2322	• For $t \geq \frac{L_1^4 128^4 \sigma_1^4}{(1-\sigma_t)^4}$, we have $\max\{\ \boldsymbol{u}_t - \boldsymbol{w}_t\ , \ \boldsymbol{u}_{t+1} - \boldsymbol{w}_t\ \} \leq \frac{1}{L_1}$ since $\ \boldsymbol{w}_{t+1} - \boldsymbol{w}_t\ $	$\ - \boldsymbol{w}_t \ \leq \frac{1}{\frac{4}{t}}$
2323	which also can be used to prove Eq. (25) .	V L
2324		
2323	The proof is completed.	
2320		
2327		
2328		
2329		
2330		
2331		
2332		
2333		
2334		
2335		
2330		
2337		
2338		
2339		
2340		
2341		
2342		
2343		
2344		
2345		
2340		
2348		
2349		
2350		
2351		
2352		
2353		
2354		
2355		
2356		
2357		
2358		
2359		
2360		
2361		
2362		
2363		
2364		
2365		
2366		
2367		
2368		
2369		
2370		
2371		
2372		
2373		
2374		
2375		