
Closing the Gap Between the Upper Bound and the Lower Bound of Adam’s Iteration Complexity

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Recently, Arjevani et al. [1] establish a lower bound of iteration complexity for
2 the first-order optimization under an L -smooth condition and a bounded noise
3 variance assumption. However, a thorough review of existing literature on Adam’s
4 convergence reveals a noticeable gap: none of them meet the above lower bound. In
5 this paper, we close the gap by deriving a new convergence guarantee of Adam, with
6 only an L -smooth condition and a bounded noise variance assumption. Our results
7 remain valid across a broad spectrum of hyperparameters. Especially with properly
8 chosen hyperparameters, we derive an upper bound of iteration complexity of
9 Adam and show that it meets the lower bound for first-order optimizers. To the best
10 of our knowledge, this is the first to establish such a tight upper bound for Adam’s
11 convergence. Our proof utilizes novel techniques to handle the entanglement
12 between momentum and adaptive learning rate and to convert the first-order term in
13 the Descent Lemma to the gradient norm, which may be of independent interest.

14 1 Introduction

15 First-order optimizers, also known as gradient-based methods, make use of gradient (first-order
16 derivative) information to find the minimum of a function. They have become a cornerstone of
17 many machine learning algorithms due to the efficiency as only gradient information is required, and
18 the flexibility as gradients can be easily computed for any function represented as directed acyclic
19 computational graph via auto-differentiation [2, 19].

20 Therefore, it is fundamental to theoretically understand the properties of these first-order methods.
21 Recently, Arjevani et al. [1] establish a lower bound on the iteration complexity of stochastic first-
22 order methods. Formally, for a well-studied setting where the objective is L -smooth and a stochastic
23 oracle can query the gradient unbiasedly with bounded variance (see Assumption 1 and 2), any stochastic
24 first-order algorithm requires at least ε^{-4} queries (in the worst case) to find an ε -stationary point, i.e.,
25 a point with gradient norm at most ε . Arjevani et al. [1] further show the above lower bound is tight
26 as it matches the existing upper bound of iteration complexity of SGD [1].

27 On the other hand, among first-order optimizers, Adam [16] becomes dominant in training state-
28 of-the-art machine learning models [3, 15, 4, 11]. Compared to vanilla stochastic gradient descent
29 (SGD), Adam consists of two more key components: (i) momentum to accumulate historical gradient
30 information and (ii) adaptive learning rate to rectify coordinate-wise step sizes. The pseudo-code
31 of Adam is given as Algorithm 1. While the sophisticated design of Adam enables its empirical
32 superiority, it brings great challenges for the theoretical analysis. After examining a series of
33 theoretical works on the upper bound of iteration complexity of Adam [24, 9, 10, 27, 14, 21, 25], we
34 find that none of them match the lower bound for first-order optimizers: they not only consume more
35 queries than the lower bound to reach ε -stationary iterations but also requires additional assumptions.

36 This theoretical mismatch becomes even more unnatural given the great empirical advantage of Adam
 37 over SGD, which incites us to think:

38 *Is the gap between the upper and lower bounds for Adam a result of the inherent complexity induced*
 39 *by Adam’s design, or could it be attributed to the proof techniques not being sharp enough?*

40 This paper answers the above question, validating the latter hypothesis, by establishing a new upper
 41 bound on iteration complexity of Adam for a wide range of hyperparameters that cover typical
 42 choices. Specifically, our contribution can be summarized as follows:

- 43 • We examine existing works that analyze the iteration complexity of Adam, and find that
 44 none of them meets the lower bound of first-order optimization algorithms;
- 45 • We derive a new convergence guarantee of Adam with only assuming L -smooth condition
 46 and bounded variance assumption (Theorem 1), which holds for a wide range of hyperpa-
 47 rameters covering typical choices;
- 48 • With chosen hyperparameters, we further tighten Theorem 1 and show that the upper bound
 49 on the iteration complexity of Adam meets the lower bound, closing the gap (Theorem 2).
 50 Our upper bound is tighter than existing results by a logarithmic factor, in spite of weaker
 51 assumption.

52 To the best of our knowledge, this work provide the first upper bound on the iteration complexity
 53 of Adam without additional assumptions other than L -smooth condition and bounded variance
 54 assumption. It is also the first upper bound matching the lower bound of first-order optimizers.

55 **Organization of this paper.** The rest of the paper is organized as follows: in Section 2, we first
 56 present the notations and setup of analysis in this paper ; in Section 3, we revisit the existing works
 57 on the iteration complexity of Adam; in Section 4, we present a convergence analysis of Adam
 58 with general hyperparameters (Theorem 1); in Section 5, we tighten Theorem 1 with a chosen
 59 hyperparameter, and derive an upper bound of Adam’s iteration complexity which meets the lower
 60 bound; in Section 6, we discuss the limitation of our results; in Section 7, we discuss the related
 61 works.

62 2 Preliminary

63 The Adam algorithm is restated in Algorithm 1 for convenient reference. Note that compared to the
 64 original version of Adam in Kingma and Ba [16], the bias-correction terms are omitted to simplify
 65 the analysis, and our analysis can be immediately extended to the original version of Adam because
 66 the effect of bias-correction term decays exponentially. Also, in the original version of Adam, the
 67 adaptive learning rate is $\frac{\eta}{\sqrt{\nu_t + \epsilon \mathbf{1}_d}}$ instead of $\frac{\eta}{\sqrt{\nu_t}}$. However, our setting is more challenging and our
 68 result can be easily extend to the original version of Adam, since the ϵ term makes the adaptive
 69 learning rate upper bounded and eases the analysis.

Algorithm 1 Adam

Input: Stochastic oracle \mathcal{O} , learning rate $\eta > 0$, initial point $\mathbf{w}_1 \in \mathbb{R}^d$, initial conditioner $\nu_0 \in \mathbb{R}^+$,
 initial momentum \mathbf{m}_0 , momentum parameter β_1 , conditioner parameter β_2 , number of epoch T

- 1: Sample $r \sim \text{Unif}\{1, \dots, T\}$
- 2: **For** $t = 1 \rightarrow T$:
- 3: Generate a random z_t , and query stochastic oracle $\mathbf{g}_t = \mathcal{O}_f(\mathbf{w}_t, z_t)$
- 4: Calculate $\nu_t = \beta_2 \nu_{t-1} + (1 - \beta_2) \mathbf{g}_t^{\odot 2}$
- 5: Calculate $\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t$
- 6: Update $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \frac{1}{\sqrt{\nu_t}} \odot \mathbf{m}_t$
- 7: **EndFor**

Output: \mathbf{w}_r

70 **Notations.** For $a, b \in \mathbb{Z}^{\geq 0}$ and $a \leq b$, denote $[a, b] = \{a, a + 1, \dots, b - 1, b\}$. For any two vectors
 71 $\mathbf{w}, \mathbf{v} \in \mathbb{R}^d$, denote $\mathbf{w} \odot \mathbf{v}$ as the Hadamard product (i.e., coordinate-wise multiplication) between
 72 \mathbf{w} and \mathbf{v} . When analyzing Adam, we denote the true gradient at iteration t as $\mathbf{G}_t = \nabla f(\mathbf{w}_t)$, and

73 the sigma algebra before iteration t as $\mathcal{F}_t = \sigma(\mathbf{g}_1, \dots, \mathbf{g}_{t-1})$. We denote conditional expectation as
 74 $\mathbb{E}^{|\mathcal{F}_t}[\ast] = \mathbb{E}[\ast|\mathcal{F}_t]$. We also use asymptotic notations \mathfrak{o} , \mathcal{O} , Ω , and Θ , where $h_2(x) = \mathfrak{o}_{x \rightarrow x_0}(h_1(x))$
 75 means that $\lim_{x \rightarrow x_0} \frac{h_2(x)}{h_1(x)} = 0$ (when the context is clear, we abbreviate $x \rightarrow x_0$ and only use
 76 $\mathfrak{o}(h_1(x))$); $h_2(x) = \mathcal{O}(h_1(x))$ means that there exists constant γ independent of x such that $h_2(x) \leq$
 77 $\gamma h_1(x)$; $h_2(x) = \Omega(h_1(x))$ means that $h_1(x) = \mathcal{O}(h_2(x))$; and $h_2(x) = \Theta(h_1(x))$ means that
 78 $h_2(x) = \mathcal{O}(h_1(x))$ and $h_2(x) = \Omega(h_1(x))$.

79 **Objective function.** In this paper, we consider solving the following optimization problem:
 80 $\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$. We make the following assumption on the objective function f .

81 **Assumption 1** (On objective function). *We assume f is differentiable, and the gradient of f is*
 82 *L -Lipschitz.*

83 We denote the set of all objective functions satisfying Assumption 1 as $\mathcal{F}(L)$.

84 **Stochastic oracle.** As f is differentiable, we can utilize the gradient of f (i.e., ∇f) to solve the
 85 above optimization problem. However, the ∇f is usually expensive to compute. Instead, we query
 86 a stochastic estimation of ∇f through a stochastic oracle \mathcal{O} . Specifically, the stochastic oracle \mathcal{O}
 87 consists of a distribution \mathcal{P} over a measurable space \mathcal{Z} and a mapping $\mathcal{O}_f : \mathbb{R}^d \times \mathcal{Z} \rightarrow \mathbb{R}^d$. We make
 88 the following assumption on \mathcal{O} .

89 **Assumption 2** (On stochastic oracle). *We assume that \mathcal{O} is unbiased, i.e., $\forall \mathbf{w} \in \mathbb{R}^d,$
 90 $\mathbb{E}_{z \sim \mathcal{P}} \mathcal{O}_f(\mathbf{w}, z) = \nabla f(\mathbf{w})$. We further assume \mathcal{O} has bounded variance, i.e., $\forall \mathbf{w} \in \mathbb{R}^d,$
 91 $\mathbb{E}_{z \sim \mathcal{P}} [\|\mathcal{O}_f(\mathbf{w}, z) - \nabla f(\mathbf{w})\|^2] \leq \sigma^2$.*

92 We denote the set of all stochastic oracles satisfying Assumption 2 with variance bound σ^2 as $\mathfrak{D}(\sigma^2)$.

93 Adam belongs to first-order optimization algorithms, which is defined as follows:

94 **Definition 1** (First-order optimization algorithm). *An algorithm \mathbf{A} is called a first-order optimization*
 95 *algorithm, if it takes an input \mathbf{w}_1 and hyperparameter θ , and produces a sequence of parameters as*
 96 *follows: first sample a random seed r from some distribution \mathcal{P}_r^1 , set $\mathbf{w}_1^{\mathbf{A}(\theta)} = \mathbf{w}_1$ and then update*
 97 *the parameters as*

$$\mathbf{w}_{t+1}^{\mathbf{A}(\theta)} = \mathbf{A}_\theta^t(r, \mathbf{w}_1^{\mathbf{A}(\theta)}, \mathcal{O}_f(\mathbf{w}_1^{\mathbf{A}(\theta)}, z_1), \dots, \mathcal{O}_f(\mathbf{w}_t^{\mathbf{A}(\theta)}, z_t)),$$

98 where z_1, z_2, \dots, z_t are sampled i.i.d. from \mathcal{P} .

99 Denote the set of all first-order optimization algorithms as $\mathcal{A}_{\text{first}}$. We next introduce *iteration*
 100 *complexity* to measure the convergence rate of optimization algorithms.

101 **Definition 2** (Iteration complexity). *The iteration complexity of first-order optimization algorithm \mathbf{A}*
 102 *is defined as*

$$\mathcal{C}_\varepsilon(\mathbf{A}, \Delta, L, \sigma^2) = \sup_{\mathcal{O} \in \mathfrak{D}(\sigma^2)} \sup_{f \in \mathcal{F}(L)} \sup_{\mathbf{w}_1 : f(\mathbf{w}_1) = \Delta} \inf_{\theta} \{T : \mathbb{E} \|\nabla f(\mathbf{w}_T^{\mathbf{A}(\theta)})\| \leq \varepsilon\}.$$

103 Furthermore, the iteration complexity of the family of first-order optimization algorithms $\mathcal{A}_{\text{first}}$ is

$$\mathcal{C}_\varepsilon(\Delta, L, \sigma^2) = \sup_{\mathcal{O} \in \mathfrak{D}(\sigma^2)} \sup_{f \in \mathcal{F}(L)} \sup_{\mathbf{w}_1 : f(\mathbf{w}_1) = \Delta} \inf_{\mathbf{A} \in \mathcal{A}_{\text{first}}} \inf_{\theta} \{T : \mathbb{E} \|\nabla f(\mathbf{w}_T^{\mathbf{A}(\theta)})\| \leq \varepsilon\}.$$

104 It should be noticed that the iteration complexity of the family of first-order optimization algorithms
 105 is a lower bound of the iteration complexity of a specific first-order optimization algorithm, i.e.,
 106 $\forall \mathbf{A} \in \mathcal{A}_{\text{first}}, \mathcal{C}_\varepsilon(\mathbf{A}, \Delta, L, \sigma^2) \geq \mathcal{C}_\varepsilon(\Delta, L, \sigma^2)$.

107 3 None of existing upper bounds match the lower bound

108 In this section, we examine existing works that study the iteration complexity of Adam, and defer a
 109 discussion of other existing works to Appendix A. We find that none of them match the lower bound
 110 for first-order algorithms provided in [1] (restated as follows).

¹Such a random seed allows sampling from all iterations to generate the final output of the optimization algorithm. As an example, Algorithm 1 set \mathcal{P}_r .

111 **Proposition 1** (Theorem 3, [1]). $\forall L, \Delta, \sigma^2 > 0$, we have $\mathcal{C}_\varepsilon(\Delta, L, \sigma^2) = \Omega(\frac{1}{\varepsilon^4})$.

112 Note that in the above bound, we omit the dependence of the lower bound over Δ, L , and σ^2 , which
 113 is a standard practice in existing works (see Cutkosky and Mehta [8], Xie et al. [23], Faw et al. [13]
 114 as examples) because the dependence over the accuracy ε can be used to derive how much additional
 115 iterations is required for a smaller target accuracy and is thus of more interest. In this paper, when we
 116 say "match the lower bound", we always mean that the upper bound has the same order of ε as the
 117 lower bound.

118 Generally speaking, existing works on the iteration complexity of Adam can be divided into two cate-
 119 gories: they either (i) assume that gradient is universally bounded or (ii) make stronger assumptions
 120 on smoothness. Below we respectively explain how these two categories of works do not match the
 121 lower bound in [1].

122 The first line of works, including Zaheer et al. [24], De et al. [9], Défossez et al. [10], Zou et al.
 123 [27], Guo et al. [14], assume that the gradient norm of f is universally bounded, i.e., $\|\nabla f(\mathbf{w})\| \leq G$,
 124 $\forall \mathbf{w} \in \mathbb{R}^d$. In other words, what they consider is another iteration complexity defined as follows:

$$\mathcal{C}_\varepsilon(\mathbf{A}, \Delta, L, \sigma^2, G) \triangleq \sup_{\mathcal{O} \in \mathfrak{D}(\sigma^2)} \sup_{f \in \mathcal{F}(L), \|\nabla f\| \leq G} \sup_{\mathbf{w}_1: f(\mathbf{w}_1) = \Delta} \inf_{\theta} \{T : \mathbb{E} \|\nabla f(\mathbf{w}_T^{\mathbf{A}(\theta)})\| \leq \varepsilon\}.$$

125 This line of works do not match the lower bound due to the following two reasons: First of all, the
 126 upper bound they derive is $O(\frac{\log 1/\varepsilon}{\varepsilon^4})$, which has an additional $\log \varepsilon$ factor more than the lower bound;
 127 secondly, the bound they derive is for $\mathcal{C}_\varepsilon(\mathbf{A}, \Delta, L, \sigma^2, G)$. Note that $\mathcal{F}(L) \cap \{f : \|\nabla f\| \leq G\}$ is a
 128 proper subset of $\mathcal{F}(L)$ for any G , where a simple example in $\mathcal{F}(L)$ but without bounded gradient is
 129 the quadratic function $f(x) = \|x\|^2$. Therefore, we have that

$$\mathcal{C}_\varepsilon(\mathbf{A}, \Delta, L, \sigma^2) \geq \mathcal{C}_\varepsilon(\mathbf{A}, \Delta, L, \sigma^2, G), \quad \forall G \geq 0, \quad (1)$$

130 and thus the upper bound on $\mathcal{C}_\varepsilon(\mathbf{A}, \Delta, L, \sigma^2, G)$ does not apply to $\mathcal{C}_\varepsilon(\mathbf{A}, \Delta, L, \sigma^2)$. Moreover, their
 131 upper bound of $\mathcal{C}_\varepsilon(\mathbf{A}, \Delta, L, \sigma^2, G)$ tends to ∞ as $G \rightarrow \infty$, which indicates that if following their
 132 analysis the upper bound of $\mathcal{C}_\varepsilon(\mathbf{A}, \Delta, L, \sigma^2)$ would be infinity based on Eq. (1).

133 The second line of works includes Shi et al. [21], Zhang et al. [25], Wang et al. [22], which additionally
 134 assume a mean-squared smoothness property besides Assumption 1 and 2, i.e., $\mathbb{E}_{z \sim \mathcal{P}} \|\mathbf{O}_f(\mathbf{w}, z) -$
 135 $\mathbf{O}_f(\mathbf{v}, z)\|^2 \leq L \|\mathbf{w} - \mathbf{v}\|^2$. Denote $\tilde{\mathfrak{D}}(\sigma^2, L) \triangleq \{\mathcal{O} : \mathbb{E}_{z \sim \mathcal{P}} \|\mathbf{O}_f(\mathbf{w}, z) - \mathbf{O}_f(\mathbf{v}, z)\|^2 \leq L \|\mathbf{w} -$
 136 $\mathbf{v}\|^2, \forall \mathbf{w}, \mathbf{v} \in \mathbb{R}^d\} \cap \mathfrak{D}(\sigma^2)$. The iteration complexity that they consider is defined as follows:

$$\tilde{\mathcal{C}}_\varepsilon(\mathbf{A}, \Delta, L, \sigma^2) = \sup_{\mathcal{O} \in \tilde{\mathfrak{D}}(\sigma^2, L)} \sup_{f \in \mathcal{F}(L)} \sup_{\mathbf{w}_1: f(\mathbf{w}_1) = \Delta} \inf_{\theta} \{T : \mathbb{E} \|\nabla f(\mathbf{w}_T^{\mathbf{A}(\theta)})\| \leq \varepsilon\}.$$

137 The rate derived in [21, 25, 22] is $O(\frac{\log 1/\varepsilon}{\varepsilon^6})$, which is derived by minimizing the upper bounds in
 138 [21, 25, 22] with respect to the hyperparameter of adaptive learning rate β_2 . According to [1], the
 139 lower bound of iteration complexity of $\tilde{\mathcal{C}}_\varepsilon(\mathbf{A}, \Delta, L, \sigma^2)$ is $\Omega(\frac{1}{\varepsilon^3})$ and smaller than the original lower
 140 bound $\Omega(\frac{1}{\varepsilon^4})$, resulting in an even larger gap between the upper bound and lower bound.

141 On the other hand, a concurrent work [17] which does not require bounded gradient assumption
 142 and mean-squared smoothness property but poses a stronger assumption on the stochastic ora-
 143 cle: the set of stochastic oracles they consider is $\tilde{\tilde{\mathfrak{D}}} = \{\mathcal{O} : \forall \mathbf{w} \in \mathbb{R}^d, \mathbb{E}_{z \sim \mathcal{P}} \mathbf{O}_f(\mathbf{w}, z) =$
 144 $\nabla f(\mathbf{w}), \mathbb{P}(\|\mathbf{O}_f(\mathbf{w}, z) - \nabla f(\mathbf{w})\|^2 \leq \sigma^2) = 1\}$. $\tilde{\tilde{\mathfrak{D}}}$ is a proper subset of \mathfrak{D} because a simple
 145 example is that $\mathbf{O}_f(\mathbf{w}, z) = \nabla f(\mathbf{w}) + z$ where z is a standard gaussian variable. Therefore, their
 146 result does not provide a valid upper bound of $\mathcal{C}_\varepsilon(\mathbf{A}, \Delta, L, \sigma^2)$.

147 4 Convergence analysis of Adam with only Assumptions 1 and 2

148 As discussed in Section 3, existing works on analyzing Adam require additional assumptions besides
 149 Assumption 1 and 2. In this section, we provide the first convergence analysis of Adam with only As-
 150 sumption 1 and 2, which naturally gives an upper bound on the iteration complexity $\mathcal{C}_\varepsilon(\mathbf{A}, \Delta, L, \sigma^2)$.
 151 Specifically, we present the following theorem.

152 **Theorem 1.** Let \mathbf{A} be by Adam (Algorithm 1) and $\theta = (\eta, \beta_1, \beta_2)$ are the hyperparameters of \mathbf{A} .
 153 Let Assumption 1 and 2 hold. Then, if $0 \leq \beta_1 < \beta_2 < 1$, we have

$$\begin{aligned} \mathbb{E} \sum_{t=1}^T \|\nabla f(\mathbf{w}_t)\| &\leq \sqrt{1-\beta_2} C_2 + \frac{2\sqrt{1-\beta_2}}{(1-\beta_1)\eta} C_1 d \ln \left(12C_2 + 2T \sum_{l=1}^d \sqrt{\nu_{0,l} + (3-\beta_2)\sigma^2} + 4dC_1 \ln dC_1 \right) \\ &\quad + \sqrt{C_2 + \frac{2}{(1-\beta_1)\eta} C_1 d \ln \left(12C_2 + 2T \sum_{l=1}^d \sqrt{\nu_{0,l} + (3-\beta_2)\sigma^2} + 4dC_1 \ln dC_1 \right)} \\ &\quad \times \sqrt{12C_2 + 2T \sum_{l=1}^d \sqrt{\nu_{0,l} + (3-\beta_2)\sigma^2} + 4dC_1 \ln dC_1}. \end{aligned} \quad (2)$$

154 where $\nu_{0,l}$ is the l -th coordinate of ν_0 ,

$$C_1 = \left(\frac{L}{2} \eta^2 + 2 \frac{\sqrt{1-\beta_2}}{(1-\beta_1)^2} \eta \sigma + \frac{\eta^2 \beta_1}{\sqrt{\beta_2} (1 - \frac{\beta_1}{\sqrt{\beta_2}})} + L^2 \frac{\beta_1 \eta^3 (1-\beta_1)}{\beta_2 (1-\beta_2)^{\frac{1}{2}} (1 - \frac{\beta_1^2}{\beta_2}) (1 - \frac{\beta_1}{\beta_2})^2} \frac{d}{\sigma} \frac{(1-\beta_1)^2}{(1 - \frac{\beta_1}{\sqrt{\beta_2}})^2} \right) \frac{1}{1-\beta_2}.$$

155 and

$$C_2 = \frac{2}{(1-\beta_1)\eta} \left(f(\mathbf{w}_1) + \sum_{l=1}^d 2C_1 \left(\mathbb{E} \ln \left(\frac{1}{\nu_{0,l}} \right) - T \ln \beta_2 \right) \right).$$

156 A proof sketch is given in Section 4.2 and the full proof is deferred to Appendix.

157 The right-hand side in Eq. (2) looks messy at the first glance. We next explain Theorem 1 in detail
 158 and make the upper bound's dependence over hyperparameters crystally clear.

159 4.1 Discussion on Theorem 1

160 **Required assumptions and conditions.** As mentioned previously, Theorem 1 only requires Assump-
 161 tion 1 and 2, which aligns with the setting of the lower bound (Proposition 1). To our best knowledge,
 162 this is the first analysis of Adam without additional assumptions. Also, Theorem 1 holds for general
 163 choices of hyperparameters since the only condition posed on hyperparameters is $\beta_1 < \beta_2$. Such
 164 condition covers a wide range of hyperparameters, e.g., the default setting $\beta_1 = 0.9$ and $\beta_2 = 0.999$
 165 in PyTorch [19].

166 **Dependence over β_2 , η , and T .** Here we consider the influence of β_2 , η , and T while fixing
 167 β_1 constant (we will discuss the effect of β_1 in Section 6). With logarithmic factors ignored and
 168 coefficients hidden, C_1 , C_2 and the right-hand-side of Eq. (2) can be rewritten with asymptotic
 169 notations as

$$\begin{aligned} C_1 &= \tilde{\mathcal{O}} \left(\frac{\eta}{\sqrt{1-\beta_2}} + \frac{\eta^3}{(1-\beta_2)^{\frac{3}{2}}} \right), \\ C_2 &= \tilde{\mathcal{O}} \left(\frac{1}{\sqrt{1-\beta_2}} + \frac{\eta^2}{(1-\beta_2)^{\frac{3}{2}}} + \frac{1}{\eta} + T \sqrt{1-\beta_2} + \frac{\eta^2 T}{(1-\beta_2)^{\frac{1}{2}}} \right), \\ \mathbb{E} \sum_{t=1}^T \|\nabla f(\mathbf{w}_t)\| &= \tilde{\mathcal{O}} \left(\sqrt{1-\beta_2} C_2 + \frac{\sqrt{1-\beta_2}}{\eta} C_1 + \sqrt{C_2 + \frac{C_1}{\eta} \sqrt{C_2 + T + C_1}} \right), \end{aligned}$$

170 where $\tilde{\mathcal{O}}$ denotes \mathcal{O} with logarithmic terms ignored. Consequently, the dependence of Eq. (2) over
 171 β_2, η and T becomes

$$\begin{aligned} \mathbb{E} \sum_{t=1}^T \|\nabla f(\mathbf{w}_t)\| &= \tilde{\mathcal{O}} \left(\frac{1}{\sqrt{1-\beta_2}} + \frac{\eta^2}{(1-\beta_2)^{\frac{3}{2}}} + \frac{1}{\eta} + \frac{\eta^2 T}{(1-\beta_2)^{\frac{1}{2}}} \right) \\ &\quad + \tilde{\mathcal{O}} \left(\frac{\sqrt{T}}{\sqrt[4]{1-\beta_2}} + \frac{\eta \sqrt{T}}{(1-\beta_2)^{\frac{3}{4}}} + \frac{\sqrt{T}}{\sqrt{\eta}} + T \sqrt[4]{1-\beta_2} + \frac{\eta T}{(1-\beta_2)^{\frac{1}{4}}} \right). \end{aligned}$$

172 Therefore, in order to ensure convergence, $\min_{t \in [T]} \mathbb{E} \|\mathbf{G}_t\|_1 \rightarrow 0$ as $T \rightarrow \infty$, a sufficient condition
 173 is that the right-hand-side of the above equation is $\mathcal{O}(T)$. Specifically, by choosing $\eta = \Theta(T^{-a})$ and
 174 $1 - \beta_2 = \Theta(T^{-b})$, we obtain that

$$\frac{1}{T} \mathbb{E} \sum_{t=1}^T \|\nabla f(\mathbf{w}_t)\| = \tilde{\mathcal{O}} \left(T^{\frac{b}{2}-1} + T^{-2a+\frac{3}{2}b-1} + T^{a-1} + T^{-2a+\frac{1}{2}b} + T^{\frac{b}{4}-\frac{1}{2}} + T^{-a+\frac{3}{4}b-\frac{1}{2}} + T^{\frac{1}{2}a-\frac{1}{2}} + T^{-a+\frac{1}{4}b} \right).$$

175 By simple calculation, we obtain that the right-hand side of the above inequality is $\mathcal{O}(1)$ as $T \rightarrow \infty$
176 if and only if $0 < \frac{b}{4} < a < 1$ and $3b - 4a < 2$. Moreover, the minimum of the right-hand side of the
177 above inequality is $\tilde{\mathcal{O}}(\frac{1}{T^{\frac{3}{4}}})$, which is achieved at $a = \frac{1}{2}$ and $b = 1$. Such a minimum implies an upper
178 bound of the iteration complexity which at most differs from the lower bound by logarithmic factors
179 as solving $\tilde{\mathcal{O}}(\frac{1}{T^{\frac{3}{4}}}) = \varepsilon$ gives $T = \tilde{\mathcal{O}}(\frac{1}{\varepsilon^{\frac{4}{3}}})$. In Theorem 2, we will further remove the logarithmic
180 factor by giving a refined proof when $a = \frac{1}{2}$ and $b = 1$ and close the gap between the upper and
181 lower bounds.

182 4.2 Proof Sketch of Theorem 1

183 In this section, we demonstrate the proof idea of Theorem 1. Concretely, we sketch the proof by
184 identifying two key challenges in the proof and provide our solutions respectively.

185 **Challenge I: Disentangle the stochasticity in momentum and adaptive learning rate.** According
186 to the standard descent lemma, we have that

$$\begin{aligned} \mathbb{E}f(\mathbf{w}_{t+1}) &= f(\mathbf{w}_t) + \mathbb{E} \left[\langle \mathbf{G}_t, \mathbf{w}_{t+1} - \mathbf{w}_t \rangle + \frac{L}{2} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|^2 \right] \\ &\leq \underbrace{\mathbb{E}f(\mathbf{w}_t) + \mathbb{E} \left[\left\langle \mathbf{G}_t, -\eta \frac{1}{\sqrt{\nu_t}} \odot \mathbf{m}_t \right\rangle \right]}_{\text{First Order}} + \underbrace{\frac{L}{2} \eta^2 \mathbb{E} \left\| \frac{1}{\sqrt{\nu_t}} \odot \mathbf{m}_t \right\|^2}_{\text{Second Order}} \end{aligned} \quad (3)$$

187 The first challenge arises from bounding the "First Order" term above. To facilitate the understanding
188 of the difficulty, we compare the "First Order" term of Adam to the corresponding "First Order" term
189 of SGD, i.e., $-\eta \mathbb{E} \langle \mathbf{G}_t, \mathbf{g}_t \rangle$. By directly applying $\mathbb{E}^{\mathcal{F}_t} g_t = \mathbf{G}_t$, we obtain that the "First-Order" term
190 of SGD equals to $-\eta \mathbb{E} \|\mathbf{G}_t\|^2$. However, as for Adam, there are two folds of trouble: firstly, we
191 do not know what $\mathbb{E}^{\mathcal{F}_t} \frac{1}{\sqrt{\nu_t}} \odot \mathbf{m}_t$ is, as the stochasticity in \mathbf{m}_t and ν_t entangles. Secondly, even
192 without ν_t , it is unclear how $\mathbb{E}^{\mathcal{F}_t} \mathbf{m}_t$ aligns with \mathbf{G}_t given the existence of $\mathbf{g}_{t-1}, \dots, \mathbf{g}_1$ in \mathbf{m}_t .

193 **Solution to Challenge I.** For $i \in [1, t]$, we define a set of surrogate conditioner $\tilde{\nu}_t^i \triangleq \beta_2^i \nu_{t-i} +$
194 $\sum_{j=0}^{i-1} \beta_2^j (1 - \beta_2) \mathbf{G}_{t-i+1}^{\odot 2} + (1 - \beta_2) \sigma^2$, and $\tilde{\nu}_t^0 \triangleq \nu_t$. Note that $\tilde{\nu}_t^i$ is measurable with respect to
195 \mathcal{F}_{t-i+1} . The key idea of our solution is the following *peeling-off strategy*: starting from $\mathbb{E}[\langle \mathbf{G}_t, \frac{1}{\sqrt{\nu_t}} \odot$
196 $\mathbf{m}_t \rangle]$, we replace $\nu_t = \tilde{\nu}_t^0$ by $\tilde{\nu}_t^1$ (of course, such a replacement will bring a error term, which we
197 temporarily ignore and will consider it in the formal proof) and obtain $\mathbb{E}[\langle \mathbf{G}_t, \frac{1}{\sqrt{\tilde{\nu}_t^1}} \odot \mathbf{m}_t \rangle]$. As
198 $\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t$, we further have $\mathbb{E}[\langle \mathbf{G}_t, \frac{1}{\sqrt{\tilde{\nu}_t^1}} \odot \mathbf{m}_t \rangle] = \mathbb{E}[\langle \mathbf{G}_t, \frac{1}{\sqrt{\tilde{\nu}_t^1}} \odot (1 - \beta_1) \mathbf{g}_t \rangle] +$
199 $\mathbb{E}[\langle \mathbf{G}_t - \mathbf{G}_{t-1}, \frac{1}{\sqrt{\tilde{\nu}_t^1}} \odot \beta_1 \mathbf{m}_{t-1} \rangle] + \mathbb{E}[\langle \mathbf{G}_{t-1}, \frac{1}{\sqrt{\tilde{\nu}_t^1}} \odot \beta_1 \mathbf{m}_{t-1} \rangle]$. As $\tilde{\nu}_t^1$ is measurable w.r.t. \mathcal{F}_t , we
200 can then disentangle the stochasticity in \mathbf{g}_t and ν_t , and the term $\mathbb{E}[\langle \mathbf{G}_t, \frac{1}{\sqrt{\tilde{\nu}_t^1}} \odot (1 - \beta_1) \mathbf{g}_t \rangle]$ equals to
201 $\mathbb{E}[\langle \mathbf{G}_t, \frac{1}{\sqrt{\tilde{\nu}_t^1}} \odot (1 - \beta_1) \mathbf{G}_t \rangle]$, which is desired. The term $\mathbb{E}[\langle \mathbf{G}_t - \mathbf{G}_{t-2}, \frac{1}{\sqrt{\tilde{\nu}_t^1}} \odot \beta_1 \mathbf{m}_{t-1} \rangle]$ is small
202 due to L -smooth condition. The term $\mathbb{E}[\langle \mathbf{G}_{t-1}, \frac{1}{\sqrt{\tilde{\nu}_t^1}} \odot \beta_1 \mathbf{m}_{t-1} \rangle]$ resembles $\mathbb{E}[\langle \mathbf{G}_t, \frac{1}{\sqrt{\nu_t}} \odot \mathbf{m}_t \rangle]$,
203 and we can apply the methodology recursively to get $\mathbb{E}[\langle \mathbf{G}_{t-2}, \frac{1}{\sqrt{\tilde{\nu}_t^2}} \odot \beta_1^2 \mathbf{m}_{t-2} \rangle]$, $\mathbb{E}[\langle \mathbf{G}_{t-3}, \frac{1}{\sqrt{\tilde{\nu}_t^3}} \odot$
204 $\beta_1^3 \mathbf{m}_{t-3} \rangle]$, and so on. All in all, the above methodology can be summarized as the following lemma.
205 **Lemma 1.** *Let all conditions in Theorem 1 hold. Denote $F_t^i \triangleq \mathbb{E}[\langle \mathbf{G}_{t-i}, \frac{1}{\sqrt{\tilde{\nu}_t^i}} \odot \mathbf{m}_{t-i} \rangle]$. Set $\mathbf{G}_0 \triangleq \mathbf{G}_1$*

206 *Then, $\forall t \geq 1$ and $i \in [0, t - 1]$,*

$$\begin{aligned} F_t^i &\geq \beta_1 F_t^{i+1} + \frac{(1 - \beta_1)}{2} \mathbb{E} \left[\left\| \frac{1}{\sqrt{\tilde{\nu}_t^{i+1}}} \odot \mathbf{G}_{t-i} \right\|^2 \right] - \beta_1 L \mathbb{E} \left[\left\| \mathbf{w}_{t-i} - \mathbf{w}_{t-i-1} \right\| \left\| \frac{1}{\sqrt{\tilde{\nu}_t^{i+1}}} \odot \mathbf{m}_{t-i-1} \right\| \right] \\ &\quad - \left(2 \frac{\sqrt{1 - \beta_2}}{1 - \beta_1} \sigma + L^2 \frac{\eta^2 (1 - \beta_1)}{(1 - \beta_2)^{\frac{1}{2}} (1 - \frac{\beta_2^2}{\beta_2}) \beta_2^2 \sigma} \frac{i}{d} \right) \mathbb{E} \left\| \frac{1}{\sqrt{\nu_{t-i}}} \odot \mathbf{m}_{t-i} \right\|^2. \end{aligned}$$

207 The proof is deferred to Appendix C.1. We highlight here that despite the simple methodology
 208 above, the proof itself is highly non-trivial and technical. The core difficulty lies in handling the error
 209 introduced by approximating $\tilde{\nu}_t^i$ with $\tilde{\nu}_t^{i+1}$, where we need to bound the gap both between g_{t-i} and
 210 \mathbf{G}_{t-i} and between \mathbf{G}_{t-i} and \mathbf{G}_{t-i+1} .

211 **Remark 1.** *Our surrogate conditioners $\tilde{\nu}_t^i$ are novel. Previously, there are other surrogate conditioners*
 212 *in Défossez et al. [10], Zou et al. [27] which help to disentangle the stochasticity in \mathbf{g}_t and ν_t .*
 213 *However, none of them can be applied in our setting because the bounded gradient assumption is*
 214 *required to use them, which is missed in our setting. Therefore, our surrogate conditioners may also*
 215 *shed light on the other analysis of Adam where no bounded gradient is assumed.*

216 Based on Lemma 1, we can estimate the "First-Order" term recursively. Combining the estimation of
 217 the "First-Order" term back to the descent lemma (Eq. (3)) and summing the descent lemma over t
 218 from 1 to T , we obtain

$$\sum_{t=1}^T \frac{(1-\beta_1)\eta}{2} \mathbb{E} \left[\left\| \frac{1}{\sqrt[4]{\tilde{\nu}_t^1}} \odot \mathbf{G}_t \right\|^2 \right] \leq f(\mathbf{w}_1) - \mathbb{E}f(\mathbf{w}_{T+1}) + \sum_{l=1}^d C_1 \left(\mathbb{E} \ln \left(\frac{\nu_{T,l}}{\nu_{0,l}} \right) - T \ln \beta_2 \right). \quad (4)$$

219 We then encounter the second challenge.

220 **Challenge II: Convert Eq. (4) to a bound of gradient norm.** Although we have bounded the sum
 221 of $\mathbb{E}[\|\frac{1}{\sqrt[4]{\tilde{\nu}_t^1}} \odot \mathbf{G}_t\|^2]$, we need to convert it into a bound of $\mathbb{E}[\|\mathbf{G}_t\|^2]$. In existing works [27, 10, 14]

222 which assumes bounded gradient, such a conversion is straightforward because (their version of) $\tilde{\nu}_t^1$
 223 is upper bounded. However, we do not assume bounded gradient and $\tilde{\nu}_t^1$ can be arbitrarily large,
 224 making $\mathbb{E}[\|\frac{1}{\sqrt[4]{\tilde{\nu}_t^1}} \odot \mathbf{G}_t\|^2]$ arbitrarily small than $\mathbb{E}[\|\mathbf{G}_t\|^2]$.

225 **Solution to Challenge II.** As this part involves coordinate-wise analysis, we define $\mathbf{g}_{t,l}$, $\mathbf{G}_{t,l}$, $\nu_{t,l}$,
 226 and $\tilde{\nu}_{t,l}^1$ respectively as the l -th coordinate of \mathbf{g}_t , \mathbf{G}_t , ν_t , and $\tilde{\nu}_t^1$. To begin with, note that due to
 227 Cauchy's inequality and Hölder's inequality,

$$\left(\mathbb{E} \sum_{t=1}^T \|\mathbf{G}_t\| \right)^2 \leq \left(\sum_{t=1}^T \mathbb{E} \left[\left\| \frac{1}{\sqrt[4]{\tilde{\nu}_t^1}} \odot \mathbf{G}_t \right\|^2 \right] \right) \left(\sum_{t=1}^T \mathbb{E} \left[\left\| \sqrt[4]{\tilde{\nu}_t^1} \right\|^2 \right] \right). \quad (5)$$

228 Therefore, we only need to derive an upper bound of $\sum_{t=1}^T \mathbb{E}[\|\sqrt[4]{\tilde{\nu}_t^1}\|^2]$, which is achieved by the
 229 following divide-and-conquer methodology. Firstly, when $|\mathbf{G}_{t,l}| \geq \sigma$, we can show $2\mathbb{E}^{|\mathcal{F}_t} |\mathbf{g}_{t,l}|^2 \geq$
 230 $2|\mathbf{G}_{t,l}|^2 \geq \mathbb{E}^{|\mathcal{F}_t} |\mathbf{g}_{t,l}|^2$. Then, by the concavity of $f(x) = \frac{x}{\sqrt{a+x}}$ ($a > 0$) and through a massive
 231 calculation, we obtain that

$$\mathbb{E} \left[\frac{|\mathbf{G}_{t,l}|^2}{\sqrt[4]{\tilde{\nu}_{t,l}^1}} \mathbf{1}_{|\mathbf{G}_{t,l}| \geq \sigma} \right] \geq \frac{1}{3(1-\beta_2)} \mathbb{E} \left(\sqrt{\nu_{t,l} + (1-\beta_2)\sigma^2} - \sqrt{\beta_2(\nu_{t-1,l} + (1-\beta_2)\sigma^2)} \right) \mathbf{1}_{|\mathbf{G}_{t,l}| \geq \sigma},$$

232 and thus

$$\sum_{t=1}^T \mathbb{E} \left[\frac{|\mathbf{G}_{t,l}|^2}{\sqrt[4]{\tilde{\nu}_{t,l}^1}} \right] \geq \sum_{t=1}^T \frac{1}{3(1-\beta_2)} \mathbb{E} \left(\sqrt{\nu_{t,l} + (1-\beta_2)\sigma^2} - \sqrt{\beta_2(\nu_{t-1,l} + (1-\beta_2)\sigma^2)} \right) \mathbf{1}_{|\mathbf{G}_{t,l}| \geq \sigma}.$$

233 Secondly, when $|\mathbf{G}_{t,l}| < \sigma$, define $\{\bar{\nu}_{t,l}\}_{t=0}^\infty$ as $\bar{\nu}_{0,l} = \nu_{0,l}$, $\bar{\nu}_{t,l} = \bar{\nu}_{t-1,l} + |\mathbf{g}_{t,l}|^2 \mathbf{1}_{|\mathbf{G}_{t,l}| < \sigma}$. One can
 234 easily observe that $\bar{\nu}_{t,l} \leq \nu_{t,l}$, and thus

$$\begin{aligned} & \sum_{t=1}^T \mathbb{E} \left(\sqrt{\nu_{t,l} + (1-\beta_2)\sigma^2} - \sqrt{\beta_2(\nu_{t-1,l} + (1-\beta_2)\sigma^2)} \right) \mathbf{1}_{|\mathbf{G}_{t,l}| \geq \sigma} \\ & \leq \sum_{t=1}^T \mathbb{E} \left(\sqrt{\bar{\nu}_{t,l} + (1-\beta_2)\sigma^2} - \sqrt{\beta_2(\bar{\nu}_{t-1,l} + (1-\beta_2)\sigma^2)} \right) \\ & = \mathbb{E} \sqrt{\bar{\nu}_{T,l} + (1-\beta_2)\sigma^2} + (1-\sqrt{\beta_2}) \sum_{t=1}^{T-1} \mathbb{E} \sqrt{\bar{\nu}_{t,l} + (1-\beta_2)\sigma^2} - \mathbb{E} \sqrt{\beta_2(\bar{\nu}_{0,l} + (1-\beta_2)\sigma^2)}. \end{aligned}$$

235 Putting the above two estimations together, we derive that

$$\sum_{t=1}^T \sum_{l=1}^d \mathbb{E} \sqrt{\boldsymbol{\nu}_{t,l} + (1 - \beta_2)\sigma^2} \leq 3(1 + \sqrt{\beta_2}) \sum_{t=1}^T \sum_{l=1}^d \mathbb{E} \left[\frac{|\mathbf{G}_{t,l}|^2}{\sqrt{\tilde{\boldsymbol{\nu}}_{t,l}^1}} \right] + T \sum_{l=1}^d \sqrt{\boldsymbol{\nu}_{0,l} + (3 - \beta_2)\sigma^2}.$$

236 The above methodology can be summarized as the following lemma.

237 **Lemma 2.** *Let all conditions in Theorem 1 hold. Then,*

$$\sum_{t=1}^T \sum_{l=1}^d \mathbb{E} \sqrt{\boldsymbol{\nu}_{t,l} + (1 - \beta_2)\sigma^2} \leq 2T \sum_{l=1}^d \sqrt{\boldsymbol{\nu}_{0,l} + (3 - \beta_2)\sigma^2} + 4dC_1 \ln dC_1 + 12C_2.$$

238 Based on Lemma 2, we can derive the estimation of $\sum_{t=1}^T \mathbb{E}[\|\sqrt[4]{\tilde{\boldsymbol{\nu}}_t^1}\|^2]$ since $\tilde{\boldsymbol{\nu}}_t^1$ is close to $\boldsymbol{\nu}_t$. The
 239 proof is then completed by combining the estimation of $\sum_{t=1}^T \mathbb{E}[\|\sqrt[4]{\tilde{\boldsymbol{\nu}}_t^1}\|^2]$ and Eq. (5).

240 5 Gap-closing upper bound on the iteration complexity of Adam

241 In this section, based on a refined proof of Stage II of Theorem 1 (see Appendix C) under the specific
 242 case $\eta = \Theta(1/\sqrt{T})$ and $\beta_2 = 1 - \Theta(1/T)$, we show that the logarithmic factor in Theorem 1 can be
 243 removed and the lower bound can be achieved. Specifically, we have the following theorem.

244 **Theorem 2.** *Let Assumption 1 and Assumption 2 hold. Then, select the hyperparameters of Adam as
 245 $\eta = \frac{a}{\sqrt{T}}$, $\beta_2 = 1 - \frac{b}{T}$ and $\beta_1 = c\beta_2$, where $a, b > 0$ and $0 \leq c < 1$ are independent of T . Then, let
 246 \mathbf{w}_τ be the output of Adam in Algorithm 1, and we have*

$$\begin{aligned} \mathbb{E}\|\nabla f(\mathbf{w}_\tau)\| &\leq \frac{1}{\sqrt[4]{T}} \sqrt{\frac{2}{\sqrt{b}} \left(D_1 + 2D_2 \ln \left(\frac{2\sqrt{b}}{\sqrt{T}} D_1 + \frac{4b}{T} D_2^2 + \sum_{l=1}^d \sqrt{\boldsymbol{\nu}_{0,l} + 3b\sigma^2} \right) \right)} \\ &\times \sqrt{\frac{2\sqrt{b}}{\sqrt{T}} D_1 + \frac{4b}{T} D_2^2 + \sum_{l=1}^d \sqrt{\boldsymbol{\nu}_{0,l} + 3b\sigma^2} + \frac{1}{T} \left(D_1 + 2D_2 \ln \left(\frac{2\sqrt{b}}{\sqrt{T}} D_1 + \frac{4b}{T} D_2^2 + \sum_{l=1}^d \sqrt{\boldsymbol{\nu}_{0,l} + 3b\sigma^2} \right) \right)}, \end{aligned}$$

247 where

$$\begin{aligned} D_1 &\triangleq \frac{4\sqrt{b}}{a(1-c)} f(\mathbf{w}_1) + \sum_{l=1}^d \frac{2}{ab\sqrt{b}} \left(La^2 + 4 \frac{a\sqrt{b}\sigma}{(1-c)^2} + 2 \frac{a^2c}{1-c} + 2 \frac{L^2ca^3d}{\sqrt{b}(1-c)^5\sigma} \right) (-\ln(\boldsymbol{\nu}_{0,l}) + b), \\ D_2 &\triangleq d \frac{2}{ab\sqrt{b}} \left(La^2 + 4 \frac{a\sqrt{b}\sigma}{(1-c)^2} + 2 \frac{a^2c}{1-c} + 4 \frac{L^2ca^3d}{\sqrt{b}(1-c)^5\sigma} \right). \end{aligned}$$

248 As a result, let \mathbf{A} be Adam in Algorithm 1, we have $\mathcal{C}_\varepsilon(\mathbf{A}, \Delta, L, \sigma^2) = \mathcal{O}(\frac{1}{\varepsilon^4})$.

249 The proof of Theorem 2 is based on a refined solution of Challenge II in the proof of Theorem 1
 250 under the specific hyperparameter settings, and we defer the concrete proof to Appendix D. Below
 251 we discuss on Theorem 2, comparing it with practice, with Theorem 1 and existing convergence rate
 252 of Adam, and with the convergence rate of AdaGrad.

253 **Alignment with the practical hyperparameter choice.** The hyperparameter setting in Theorem
 254 2 indicates that to achieve the lower bound of iteration complexity, we need to select small η and
 255 close-to-1 β_2 , with less requirement over β_1 . This agrees with the hyperparameter setting in deep
 256 learning libraries, for example, $\eta = 10^{-3}$, $\beta_2 = 0.999$, and $\beta_1 = 0.9$ in PyTorch.

257 **Comparison with Theorem 1 and existing works.** To our best knowledge, Theorem 2 is the first to
 258 derive the iteration complexity $\mathcal{O}(\frac{1}{\varepsilon^4})$. Previously, the state-of-art iteration complexity is $\mathcal{O}(\frac{\log 1/\varepsilon}{\varepsilon^4})$
 259 [10] where they additionally assume bounded gradient. Theorem 2 is also tight than Theorem 1 (while
 260 Theorem 1 holds for more general hyperparameter settings). As discussed in Section 4.1, if applying
 261 the hyperparameter setting in Theorem 2 (i.e., $\eta = \frac{a}{\sqrt{T}}$, $\beta_2 = 1 - \frac{b}{T}$ and $\beta_1 = c\beta_2$) to Theorem 1,
 262 we will obtain that $\mathbb{E}\|\nabla f(\mathbf{w}_\tau)\| \leq \mathcal{O}(\text{poly}(\log T)/\sqrt[4]{T})$ and $\mathcal{C}_\varepsilon(\mathbf{A}, \Delta, L, \sigma^2) = \mathcal{O}(\frac{\log 1/\varepsilon}{\varepsilon^4})$, which

263 is worse than the upper bound in Theorem 2 and the lower bound in Proposition 1 by a logarithmic
 264 factor.

265 **Comparison with AdaGrad.** AdaGrad [12] is another popular adaptive optimizer. Under Assump-
 266 tions 1 and 2, the state-of-art iteration complexity of AdaGrad is $\mathcal{O}(\frac{\log 1/\varepsilon}{\varepsilon^4})$ [13], which is worse
 267 than Adam by a logarithmic factor. Here we show that such a gap may be not due to the limitation
 268 of analysis, and can be explained by analogizing AdaGrad to Adam without momentum as SGD
 269 with diminishing learning rate to SGD with constant learning rate. To start with, the update rule of
 270 AdaGrad is given as

$$\boldsymbol{\nu}_t = \boldsymbol{\nu}_{t-1} + \mathbf{g}_t^{\odot 2}, \mathbf{w}_{t+1} = \mathbf{w}_t - \eta \frac{1}{\sqrt{\boldsymbol{\nu}_t}} \odot \mathbf{g}_t. \quad (6)$$

271 We first show that in Algorithm 1, if we allow the hyperparameters to be dynamical, i.e.,

$$\boldsymbol{\nu}_t = \beta_{2,t} \boldsymbol{\nu}_{t-1} + (1 - \beta_{2,t}) \mathbf{g}_t^{\odot 2}, \mathbf{m}_t = \beta_{1,t} \mathbf{m}_{t-1} + (1 - \beta_{1,t}) \mathbf{g}_t, \mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \frac{1}{\sqrt{\boldsymbol{\nu}_t}} \odot \mathbf{m}_t, \quad (7)$$

272 then Adam is equivalent to AdaGrad by setting $\eta_t = \frac{\eta}{\sqrt{t}}$, $\beta_{1,t} = 0$, and $\beta_{2,t} = 1 - \frac{1}{t}$. Specifically, by
 273 setting $\boldsymbol{\mu}_t = t \boldsymbol{\nu}_t$ in Eq. (7), we have Eq. (7) is equivalent to with Eq. (6) (by replacing $\boldsymbol{\nu}_t$ by $\boldsymbol{\mu}_t$ in
 274 Eq. (6)). Comparing the above hyperparameter setting with that in Theorem 2, we see that the above
 275 hyperparameter setting can be obtained by changing T to t and setting $c = 0$ in Theorem 2. This
 276 is similar to the relationship between SGD with diminishing learning rate $\Theta(1/\sqrt{t})$ and SGD with
 277 diminishing learning rate $\Theta(1/\sqrt{T})$. Moreover, the iteration complexity of SGD with diminishing
 278 learning rate $\Theta(1/\sqrt{t})$ also has an additional logarithmic factor than SGD with constant learning rate,
 279 which may explain the gap between AdaGrad and Adam.

280 6 Limitations

281 Despite that our work provide the first result closing the upper bound and lower bound of the iteration
 282 complexity of Adam, there are several limitations listed as follows:

283 **Dependence over the dimension d .** The bounds in Theorem 1 and Theorem 2 is monotonously
 284 increasing with respect to d . This is undesired since the upper bound of iteration complexity of SGD
 285 is invariant with respect to d . Nevertheless, removing such an dependence over d is technically hard
 286 since we need to deal with every coordinate separately due to coordinate-wise learning rate, while the
 287 descent lemma does not hold for a single coordinate but combines all coordinates together. To our
 288 best knowledge, all existing works on the convergence of Adam also suffers from the same problem.
 289 We leave removing the dependence over d as an important future work.

290 **No better result with momentum.** It can be observed that in Theorem 1 and Theorem 2, the tightest
 291 bound is achieved when $\beta_1 = 0$ (i.e., no momentum is applied). This contradicts with the common
 292 wisdom that momentum helps to accelerate. Although the benefit of momentum is not very clear for
 293 simple optimizer SGD with momentum, we view this as a limitation of our work and defer proving
 294 the benefit of momentum in Adam as a future work.

295 7 Related works

296 Section 3 has provided a detailed discussion over existing convergence analysis of Adam. In this
 297 section, we briefly review other related works. Adam is proposed with a convergence analysis in
 298 online optimization [16]. The proof, however, is latter shown to be flawed in Reddi et al. [20] as it
 299 requires the adaptive learning rate of Adam to be non-increasing. This motivates a line of works
 300 modifying Adam to ensure convergence. The modifications include enforcing the adaptive learning
 301 rate to be non-increasing [20, 5], imposing upper bound and lower bound of the adaptive learning
 302 rate [18], and using different approach to estimate second-order momentum [26, 7]. Recently, Chen
 303 et al. [6] discover a new optimizer Lion through Symbolic Discovery, which uses sign operation to
 304 replace the adaptive learning rate in Adam, achieving comparable performance of Adam with less
 305 memory costs.

References

- 306
- 307 [1] Y. Arjevani, Y. Carmon, J. C. Duchi, D. J. Foster, N. Srebro, and B. Woodworth. Lower bounds
308 for non-convex stochastic optimization. *Mathematical Programming*, pages 1–50, 2022.
- 309 [2] J. Bradbury, R. Frostig, P. Hawkins, M. J. Johnson, C. Leary, D. Maclaurin, G. Necula, A. Paszke,
310 J. VanderPlas, S. Wanderman-Milne, and Q. Zhang. JAX: composable transformations of
311 Python+NumPy programs, 2018. URL <http://github.com/google/jax>.
- 312 [3] A. Brock, J. Donahue, and K. Simonyan. Large scale gan training for high fidelity natural image
313 synthesis. In *International Conference on Learning Representations*, 2018.
- 314 [4] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam,
315 G. Sastry, A. Askell, et al. Language models are few-shot learners. *Advances in neural
316 information processing systems*, 33:1877–1901, 2020.
- 317 [5] X. Chen, S. Liu, R. Sun, and M. Hong. On the convergence of a class of Adam-type algorithms
318 for non-convex optimization. *arXiv preprint arXiv:1808.02941*, 2018.
- 319 [6] X. Chen, C. Liang, D. Huang, E. Real, K. Wang, Y. Liu, H. Pham, X. Dong, T. Luong, C.-J.
320 Hsieh, et al. Symbolic discovery of optimization algorithms. *arXiv preprint arXiv:2302.06675*,
321 2023.
- 322 [7] M. Crawshaw, M. Liu, F. Orabona, W. Zhang, and Z. Zhuang. Robustness to unbounded
323 smoothness of generalized signSGD. *arXiv preprint arXiv:2208.11195*, 2022.
- 324 [8] A. Cutkosky and H. Mehta. Momentum improves normalized SGD. In *International conference
325 on machine learning*, pages 2260–2268. PMLR, 2020.
- 326 [9] S. De, A. Mukherjee, and E. Ullah. Convergence guarantees for RMSProp and ADAM in
327 non-convex optimization and an empirical comparison to Nesterov acceleration. *arXiv preprint
328 arXiv:1807.06766*, 2018.
- 329 [10] A. Défossez, L. Bottou, F. Bach, and N. Usunier. A simple convergence proof of Adam and
330 Adagrad. *Transactions on Machine Learning Research*, 2022.
- 331 [11] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani,
332 M. Minderer, G. Heigold, S. Gelly, et al. An image is worth 16x16 words: Transformers for
333 image recognition at scale. In *International Conference on Learning Representations*, 2020.
- 334 [12] J. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and
335 stochastic optimization. *Journal of machine learning research*, 12(7), 2011.
- 336 [13] M. Faw, I. Tziotis, C. Caramanis, A. Mokhtari, S. Shakkottai, and R. Ward. The power of
337 adaptivity in SGD: Self-tuning step sizes with unbounded gradients and affine variance. In
338 *Conference on Learning Theory*, pages 313–355. PMLR, 2022.
- 339 [14] Z. Guo, Y. Xu, W. Yin, R. Jin, and T. Yang. A novel convergence analysis for algorithms of the
340 Adam family. *arXiv preprint arXiv:2112.03459*, 2021.
- 341 [15] J. D. M.-W. C. Kenton and L. K. Toutanova. Bert: Pre-training of deep bidirectional transformers
342 for language understanding. In *Proceedings of NAACL-HLT*, pages 4171–4186, 2019.
- 343 [16] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint
344 arXiv:1412.6980*, 2014.
- 345 [17] H. Li, A. Jadbabaie, and A. Rakhlin. Convergence of Adam under relaxed assumptions. *arXiv
346 preprint arXiv:2304.13972*, 2023.
- 347 [18] L. Luo, Y. Xiong, Y. Liu, and X. Sun. Adaptive gradient methods with dynamic bound of
348 learning rate. *arXiv preprint arXiv:1902.09843*, 2019.
- 349 [19] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin,
350 N. Gimelshein, L. Antiga, et al. Pytorch: An imperative style, high-performance deep learning
351 library. volume 32, 2019.
- 352 [20] S. J. Reddi, S. Kale, and S. Kumar. On the convergence of Adam and beyond. *arXiv preprint
353 arXiv:1904.09237*, 2019.
- 354 [21] N. Shi, D. Li, M. Hong, and R. Sun. RMSprop converges with proper hyper-parameter. In
355 *International Conference on Learning Representations*, 2021.
- 356 [22] B. Wang, Y. Zhang, H. Zhang, Q. Meng, Z.-M. Ma, T.-Y. Liu, and W. Chen. Provable adaptivity
357 in Adam. *arXiv preprint arXiv:2208.09900*, 2022.

- 358 [23] X. Xie, P. Zhou, H. Li, Z. Lin, and S. Yan. Adan: Adaptive nesterov momentum algorithm for
359 faster optimizing deep models. *arXiv preprint arXiv:2208.06677*, 2022.
- 360 [24] M. Zaheer, S. Reddi, D. Sachan, S. Kale, and S. Kumar. Adaptive methods for nonconvex
361 optimization. *Advances in neural information processing systems*, 31, 2018.
- 362 [25] Y. Zhang, C. Chen, N. Shi, R. Sun, and Z.-Q. Luo. Adam can converge without any modification
363 on update rules. *arXiv preprint arXiv:2208.09632*, 2022.
- 364 [26] Z. Zhou, Q. Zhang, G. Lu, H. Wang, W. Zhang, and Y. Yu. Adashift: Decorrelation and
365 convergence of adaptive learning rate methods. *arXiv preprint arXiv:1810.00143*, 2018.
- 366 [27] F. Zou, L. Shen, Z. Jie, W. Zhang, and W. Liu. A sufficient condition for convergences of Adam
367 and RMSProp. In *Proceedings of the IEEE/CVF conference on computer vision and pattern
368 recognition*, pages 11127–11135, 2019.

369 **A Related Works**

370 **B Auxilliary Lemmas**

371 The following two lemmas are useful when bounding the second-order term.

372 **Lemma 3.** Assume we have $0 < \beta_2 < 1$ and a sequence of real numbers $(a_n)_{n=1}^\infty$. Let $b_0 > 0$ and
 373 $b_n = \beta_2 b_{n-1} + (1 - \beta_2) a_n^2$. Then, we have

$$\sum_{n=1}^T \frac{a_n^2}{b_n} \leq \frac{1}{1 - \beta_2} \left(\ln \left(\frac{b_T}{b_0} \right) - T \ln \beta_2 \right).$$

374 **Lemma 4.** Assume we have $0 < \beta_1^2 < \beta_2 < 1$ and a sequence of real numbers $(a_n)_{n=1}^\infty$. Let $b_0 > 0$,
 375 $b_n = \beta_2 b_{n-1} + (1 - \beta_2) a_n^2$, $c_0 = 0$, and $c_n = \beta_1 c_{n-1} + (1 - \beta_1) a_n$. Then, we have

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

376 *Proof.* To begin with,

$$\frac{|c_n|}{\sqrt{b_n}} \leq (1 - \beta_1) \sum_{i=1}^n \frac{\beta_1^{n-i} |a_i|}{\sqrt{b_n}} \leq (1 - \beta_1) \sum_{i=1}^n \frac{\beta_1^{n-i} |a_i|}{\sqrt{b_n}} \leq (1 - \beta_1) \sum_{i=1}^n \left(\frac{\beta_1}{\sqrt{\beta_2}} \right)^{n-i} \frac{|a_i|}{\sqrt{b_i}}.$$

377 Applying Cauchy's inequality, we obtain

$$\begin{aligned} \frac{|c_n|^2}{b_n} &\leq (1 - \beta_1)^2 \left(\sum_{i=1}^n \left(\frac{\beta_1}{\sqrt{\beta_2}} \right)^{n-i} \frac{|a_i|}{\sqrt{b_i}} \right)^2 \\ &\leq (1 - \beta_1)^2 \left(\sum_{i=1}^n \left(\frac{\beta_1}{\sqrt{\beta_2}} \right)^{n-i} \right) \left(\sum_{i=1}^n \left(\frac{\beta_1}{\sqrt{\beta_2}} \right)^{n-i} \frac{|a_i|^2}{b_i} \right) \leq \frac{(1 - \beta_1)^2}{1 - \frac{\beta_1}{\sqrt{\beta_2}}} \left(\sum_{i=1}^n \left(\frac{\beta_1}{\sqrt{\beta_2}} \right)^{n-i} \frac{|a_i|^2}{b_i} \right). \end{aligned}$$

378 Summing the above inequality over n from 1 to T then leads to

$$\begin{aligned} \sum_{n=1}^T \frac{|c_n|^2}{b_n} &\leq \frac{(1 - \beta_1)^2}{1 - \frac{\beta_1}{\sqrt{\beta_2}}} \sum_{n=1}^T \left(\sum_{i=1}^n \left(\frac{\beta_1}{\sqrt{\beta_2}} \right)^{n-i} \frac{|a_i|^2}{b_i} \right) = \frac{(1 - \beta_1)^2}{1 - \frac{\beta_1}{\sqrt{\beta_2}}} \sum_{n=1}^T \frac{|a_n|^2}{b_n} \left(\sum_{i=0}^{T-n} \left(\frac{\beta_1}{\sqrt{\beta_2}} \right)^i \right) \\ &\leq \frac{(1 - \beta_1)^2}{(1 - \frac{\beta_1}{\sqrt{\beta_2}})^2} \sum_{n=1}^T \frac{|a_n|^2}{b_n} \leq \frac{(1 - \beta_1)^2}{(1 - \frac{\beta_1}{\sqrt{\beta_2}})^2 (1 - \beta_2)} \left(\ln \left(\frac{b_T}{b_0} \right) - T \ln \beta_2 \right). \end{aligned}$$

379 The proof is completed. □

380 The following lemma bound the update norm of Adam.

381 **Lemma 5.** We have $\forall t \geq 1$, $|\mathbf{w}_{t+1,l} - \mathbf{w}_{t,l}| \leq \eta \frac{1 - \beta_1}{\sqrt{1 - \beta_2} \sqrt{1 - \frac{\beta_1^2}{\beta_2}}}$.

382 *Proof.* We have that

$$\begin{aligned} |\mathbf{w}_{t+1,l} - \mathbf{w}_{t,l}| &= \eta \left| \frac{\mathbf{m}_{t,l}}{\sqrt{\nu_{t,l}}} \right| \leq \eta \frac{\sum_{i=0}^{t-1} (1 - \beta_1) \beta_1^i |\mathbf{g}_{t-i,l}|}{\sqrt{\sum_{i=0}^{t-1} (1 - \beta_2) \beta_2^i |\mathbf{g}_{t-i,l}|^2 + \beta_2^t \nu_{0,l}}} \\ &\leq \eta \frac{1 - \beta_1}{\sqrt{1 - \beta_2}} \frac{\sqrt{\sum_{i=0}^{t-1} \beta_2^i |\mathbf{g}_{t-i,l}|^2} \sqrt{\sum_{i=0}^{t-1} \frac{\beta_1^{2i}}{\beta_2^i}}}{\sqrt{\sum_{i=0}^{t-1} \beta_2^i |\mathbf{g}_{t-i,l}|^2}} \leq \eta \frac{1 - \beta_1}{\sqrt{1 - \beta_2} \sqrt{1 - \frac{\beta_1^2}{\beta_2}}}. \end{aligned}$$

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

384 **C Proof of Theorem 1**

385 **C.1 Proof of Lemma 1 and Lemma 2**

386 *Proof of Lemma 1.* $\forall i \in [0, t-1]$, we have the following decomposition:

$$F_t^i = \underbrace{\mathbb{E} \left[\left\langle \mathbf{G}_{t-i}, \frac{1}{\sqrt{\tilde{\nu}_t^{i+1}}} \odot \mathbf{m}_{t-i} \right\rangle \right]}_{(i)_t^i} + \underbrace{\mathbb{E} \left[\left\langle \mathbf{G}_{t-i}, \left(\frac{1}{\sqrt{\tilde{\nu}_t^i}} - \frac{1}{\sqrt{\tilde{\nu}_t^{i+1}}} \right) \odot \mathbf{m}_{t-i} \right\rangle \right]}_{(ii)_t^i}.$$

387 As for $(i)_t^i$, according to the definition of \mathbf{m}_{t-i} , it can be lower bounded as

$$\begin{aligned} & \mathbb{E} \left[\left\langle \mathbf{G}_{t-i}, \frac{1}{\sqrt{\tilde{\nu}_t^{i+1}}} \odot \mathbf{m}_{t-i} \right\rangle \right] = \mathbb{E} \left[\left\langle \mathbf{G}_{t-i}, (1-\beta_1) \frac{1}{\sqrt{\tilde{\nu}_t^{i+1}}} \odot \mathbf{g}_{t-i} \right\rangle \right] + \mathbb{E} \left[\left\langle \mathbf{G}_{t-i}, \beta_1 \frac{1}{\sqrt{\tilde{\nu}_t^{i+1}}} \odot \mathbf{m}_{t-i-1} \right\rangle \right] \\ &= \mathbb{E} \left[(1-\beta_1) \left\| \frac{1}{\sqrt{\tilde{\nu}_t^{i+1}}} \odot \mathbf{G}_{t-i} \right\|^2 \right] + \mathbb{E} \left[\left\langle \mathbf{G}_{t-i-1}, \beta_1 \frac{1}{\sqrt{\tilde{\nu}_t^{i+1}}} \odot \mathbf{m}_{t-i-1} \right\rangle \right] + \mathbb{E} \left[\left\langle \mathbf{G}_{t-i} - \mathbf{G}_{t-i-1}, \beta_1 \frac{1}{\sqrt{\tilde{\nu}_t^{i+1}}} \odot \mathbf{m}_{t-i-1} \right\rangle \right] \\ &\geq \mathbb{E} \left[(1-\beta_1) \left\| \frac{1}{\sqrt{\tilde{\nu}_t^{i+1}}} \odot \mathbf{G}_{t-i} \right\|^2 \right] + \mathbb{E} \left[\left\langle \mathbf{G}_{t-i-1}, \beta_1 \frac{1}{\sqrt{\tilde{\nu}_t^{i+1}}} \odot \mathbf{m}_{t-i-1} \right\rangle \right] - \beta_1 L \mathbb{E} \left[\|\mathbf{w}_{t-i} - \mathbf{w}_{t-i-1}\| \left\| \frac{1}{\sqrt{\tilde{\nu}_t^{i+1}}} \odot \mathbf{m}_{t-i-1} \right\| \right], \end{aligned}$$

388 where the last inequality is due to Assumption 1. As for $(ii)_t^i$, if $i = 0$, we have

$$\begin{aligned} & \left| \mathbb{E}^{\mathcal{F}_t} \left[\left\langle \mathbf{G}_t, \left(\frac{1}{\sqrt{\tilde{\nu}_t^0}} - \frac{1}{\sqrt{\tilde{\nu}_t^1}} \right) \odot \mathbf{m}_t \right\rangle \right] \right| \leq \sum_{l=1}^d |\mathbf{G}_{t,l}| \mathbb{E}^{\mathcal{F}_t} \left[|\mathbf{m}_{t,l}| \left| \frac{1}{\sqrt{\nu_{t,l}}} - \frac{1}{\sqrt{\tilde{\nu}_{t,l}^1}} \right| \right] \\ &= \sum_{l=1}^d |\mathbf{G}_{t,l}| \mathbb{E}^{\mathcal{F}_t} \left[|\mathbf{m}_{t,l}| \frac{(1-\beta_2) \left(|\mathbf{G}_{t,l}|^2 - |\mathbf{g}_{t,l}|^2 \right) + (1-\beta_2)\sigma^2}{\sqrt{\nu_{t,l} \tilde{\nu}_{t,l}^1} (\sqrt{\nu_{t,l}} + \sqrt{\tilde{\nu}_{t,l}^1})} \right] \\ &\stackrel{(*)}{\leq} \sum_{l=1}^d |\mathbf{G}_{t,l}| \mathbb{E}^{\mathcal{F}_t} \left[|\mathbf{m}_{t,l}| \frac{(1-\beta_2) |\mathbf{G}_{t,l} - \mathbf{g}_{t,l}| (|\mathbf{G}_{t,l}| + |\mathbf{g}_{t,l}|) + (1-\beta_2)\sigma^2}{\sqrt{\nu_{t,l} \tilde{\nu}_{t,l}^1} (\sqrt{\nu_{t,l}} + \sqrt{\tilde{\nu}_{t,l}^1})} \right] \\ &\leq \sum_{l=1}^d |\mathbf{G}_{t,l}| \mathbb{E}^{\mathcal{F}_t} \left[|\mathbf{m}_{t,l}| \frac{\sqrt{1-\beta_2} |\mathbf{G}_{t,l} - \mathbf{g}_{t,l}| + \sqrt{1-\beta_2} \sigma}{\sqrt{\nu_{t,l} \tilde{\nu}_{t,l}^1}} \right] \\ &\stackrel{(*)}{\leq} \sum_{l=1}^d \frac{\sqrt{1-\beta_2} (1-\beta_1) |\mathbf{G}_{t,l}|^2}{4\sigma \tilde{\nu}_{t,l}^1} \left(\mathbb{E}^{\mathcal{F}_t} |\mathbf{G}_{t,l} - \mathbf{g}_{t,l}|^2 + \sigma^2 \right) + \sum_{l=1}^d 2 \frac{\sqrt{1-\beta_2}}{1-\beta_1} \mathbb{E}^{\mathcal{F}_t} \sigma \frac{|\mathbf{m}_{t,l}|^2}{\nu_{t,l}} \\ &\leq \sum_{l=1}^d \frac{(1-\beta_1) |\mathbf{G}_{t,l}|^2}{2\sqrt{\tilde{\nu}_{t,l}^1}} + \sum_{l=1}^d 2 \frac{\sqrt{1-\beta_2}}{1-\beta_1} \mathbb{E}^{\mathcal{F}_t} \sigma \frac{|\mathbf{m}_{t,l}|^2}{\nu_{t,l}}, \end{aligned}$$

389 where inequality $(*)$ is due to the triangle inequality, and inequality $(*)$ is due to the mean-value
390 inequality, and the last inequality is due to Assumption 2.

391 If $i > 0$, then

$$\begin{aligned}
& \left| \mathbb{E}^{|\mathcal{F}_{t-i}} \left[\left\langle \mathbf{G}_{t-i}, \left(\frac{1}{\sqrt{\tilde{\nu}_t^i}} - \frac{1}{\sqrt{\tilde{\nu}_t^{i+1}}} \right) \odot \mathbf{m}_{t-i} \right\rangle \right] \right| \\
& \leq \sum_{l=1}^d |\mathbf{G}_{t-i,l}| \mathbb{E}^{|\mathcal{F}_{t-i}} \left[|\mathbf{m}_{t-i,l}| \left| \frac{1}{\sqrt{\tilde{\nu}_{t,l}^i}} - \frac{1}{\sqrt{\tilde{\nu}_{t,l}^{i+1}}} \right| \right] \\
& \leq \sum_{l=1}^d |\mathbf{G}_{t-i,l}| \mathbb{E}^{|\mathcal{F}_{t-i}} \left[|\mathbf{m}_{t-i,l}| \frac{(1-\beta_2)\beta_2^i \|\mathbf{G}_{t-i,l}\|^2 - |\mathbf{g}_{t-i,l}|^2 + \sum_{j=0}^{i-1} \beta_2^j (1-\beta_2) \|\mathbf{G}_{t-i,l}\|^2 - \|\mathbf{G}_{t-i+1,l}\|^2}{\sqrt{\tilde{\nu}_{t,l}^i \tilde{\nu}_{t,l}^{i+1}} (\sqrt{\tilde{\nu}_{t,l}^i} + \sqrt{\tilde{\nu}_{t,l}^{i+1}})} \right] \\
& \leq \sum_{l=1}^d |\mathbf{G}_{t-i,l}| \mathbb{E}^{|\mathcal{F}_{t-i}} \left[|\mathbf{m}_{t-i,l}| \frac{(1-\beta_2)\beta_2^i \mathbf{G}_{t-i,l} - \mathbf{g}_{t-i,l} (|\mathbf{G}_{t-i,l}| + |\mathbf{g}_{t-i,l}|)}{\sqrt{\tilde{\nu}_{t,l}^i \tilde{\nu}_{t,l}^{i+1}} (\sqrt{\tilde{\nu}_{t,l}^i} + \sqrt{\tilde{\nu}_{t,l}^{i+1}})} \right] \\
& \quad + \sum_{l=1}^d |\mathbf{G}_{t-i,l}| \mathbb{E}^{|\mathcal{F}_{t-i}} \left[|\mathbf{m}_{t-i,l}| \frac{\sum_{j=0}^{i-1} \beta_2^j (1-\beta_2) \|\mathbf{G}_{t-i,l}\| - \|\mathbf{G}_{t-i+1,l}\| (|\mathbf{G}_{t-i,l}| + \|\mathbf{G}_{t-i+1,l}\|)}{\sqrt{\tilde{\nu}_{t,l}^i \tilde{\nu}_{t,l}^{i+1}} (\sqrt{\tilde{\nu}_{t,l}^i} + \sqrt{\tilde{\nu}_{t,l}^{i+1}})} \right]
\end{aligned}$$

392 Applying Cauchy's inequality, we obtain the RHS of the above inequality is smaller than

$$\begin{aligned}
& \sum_{l=1}^d |\mathbf{G}_{t-i,l}| \mathbb{E}^{|\mathcal{F}_{t-i}} \left[|\mathbf{m}_{t-i,l}| \frac{\sqrt{1-\beta_2} \sqrt{\beta_2^i} |\mathbf{G}_{t-i,l} - \mathbf{g}_{t-i,l}|}{\sqrt{\tilde{\nu}_{t,l}^i \tilde{\nu}_{t,l}^{i+1}}} \right] \\
& \quad + \sum_{l=1}^d |\mathbf{G}_{t-i,l}| \mathbb{E}^{|\mathcal{F}_{t-i}} \left[|\mathbf{m}_{t-i,l}| \frac{\sqrt{\sum_{j=0}^{i-1} \beta_2^j (1-\beta_2) \|\mathbf{G}_{t-i,l} - \mathbf{G}_{t-i+1,l}\|^2}}{\sqrt{\tilde{\nu}_{t,l}^i \tilde{\nu}_{t,l}^{i+1}}} \right] \\
& \stackrel{(\star)}{\leq} \sum_{l=1}^d \frac{\sqrt{1-\beta_2} (1-\beta_1) \|\mathbf{G}_{t-i,l}\|^2}{4\sigma \tilde{\nu}_{t,l}^{i+1}} \left(\mathbb{E}^{|\mathcal{F}_t} |\mathbf{G}_{t,l} - \mathbf{g}_{t,l}|^2 \right) + \sum_{l=1}^d \frac{\sqrt{1-\beta_2}}{(1-\beta_1)} \beta_2^i \mathbb{E}^{|\mathcal{F}_{t-i}} \sigma \frac{|\mathbf{m}_{t-i,l}|^2}{\tilde{\nu}_{t,l}^i} \\
& \quad + \sum_{l=1}^d \frac{\sqrt{1-\beta_2} \sigma (1-\beta_1) \|\mathbf{G}_{t-i,l}\|^2}{4\tilde{\nu}_{t,l}^{i+1}} + \sum_{l=1}^d \frac{1}{(1-\beta_1)} \mathbb{E}^{|\mathcal{F}_{t-i}} \frac{1}{\sigma} \frac{|\mathbf{m}_{t-i,l}|^2}{\tilde{\nu}_{t,l}^i} \left(\sum_{j=0}^{i-1} \beta_2^j \sqrt{1-\beta_2} \|\mathbf{G}_{t-i,l} - \mathbf{G}_{t-i+1,l}\|^2 \right) \\
& \leq \sum_{l=1}^d \frac{(1-\beta_1) \|\mathbf{G}_{t-i,l}\|^2}{2\sqrt{\tilde{\nu}_{t,l}^{i+1}}} + \sum_{l=1}^d \frac{\sqrt{1-\beta_2}}{1-\beta_1} \mathbb{E}^{|\mathcal{F}_{t-i}} \sigma \beta_2^i \frac{|\mathbf{m}_{t-i,l}|^2}{\nu_{t,l}} + L^2 \frac{\eta^2 \sqrt{1-\beta_2}}{(1-\beta_1) \beta_2^i \sigma} \mathbb{E}^{|\mathcal{F}_{t-i}} \left(\sum_{l=1}^d \frac{|\mathbf{m}_{t-i,l}|^2}{\nu_{t-i,l}} \right)^2 \\
& \stackrel{(\circ)}{\leq} \sum_{l=1}^d \frac{(1-\beta_1) \|\mathbf{G}_{t-i,l}\|^2}{2\sqrt{\tilde{\nu}_{t,l}^{i+1}}} + \sum_{l=1}^d \frac{\sqrt{1-\beta_2}}{(1-\beta_1)} \mathbb{E}^{|\mathcal{F}_{t-i}} \sigma \frac{|\mathbf{m}_{t-i,l}|^2}{\nu_{t-i,l}} + L^2 \frac{\eta^2 (1-\beta_1)}{(1-\beta_2)^{\frac{1}{2}} (1-\frac{\beta_2^i}{\beta_2}) \beta_2^i \sigma} \mathbb{E}^{|\mathcal{F}_{t-i}} \left(\sum_{l=1}^d \frac{|\mathbf{m}_{t-i,l}|^2}{\nu_{t-i,l}} \right).
\end{aligned}$$

393 Here inequality (\star) is due to the mean-value inequality, and inequality (\circ) is due to Lemma 5. Putting
394 the estimation of $(i)_t^i$ and $(ii)_t^i$ together completes the proof. \square

395 *Proof of Lemma 2.* To begin with, we have that

$$\sum_{t=1}^T \mathbb{E} \left[\frac{|\mathbf{G}_{t,l}|^2}{\sqrt{\tilde{\nu}_{t,l}^1}} \mathbf{1}_{|\mathbf{G}_{t,l}| \geq \sigma} \right] \leq \sum_{t=1}^T \mathbb{E} \left[\frac{|\mathbf{G}_{t,l}|^2}{\sqrt{\tilde{\nu}_{t,l}^1}} \right]. \quad (8)$$

396 On the other hand, we have that

$$\begin{aligned}
& \frac{|\mathbf{G}_{t,l}|^2}{\sqrt{\tilde{\nu}_{t,l}^1}} \mathbf{1}_{|\mathbf{G}_{t,l}| \geq \sigma} \geq \frac{\frac{2}{3} \|\mathbf{G}_{t,l}\|^2 + \frac{1}{3} \sigma^2}{\sqrt{\tilde{\nu}_{t,l}^1}} \mathbf{1}_{|\mathbf{G}_{t,l}| \geq \sigma} \geq \frac{\frac{1}{3} \mathbb{E}^{|\mathcal{F}_t} |\mathbf{g}_{t,l}|^2 + \frac{1-\beta_2}{3} \sigma^2}{\sqrt{\tilde{\nu}_{t,l}^1}} \mathbf{1}_{|\mathbf{G}_{t,l}| \geq \sigma} \\
& \geq \frac{\frac{1}{3} \mathbb{E}^{|\mathcal{F}_t} |\mathbf{g}_{t,l}|^2 + \frac{1-\beta_2}{3} \sigma^2}{\sqrt{\beta_2 \nu_{t-1,l} + (1-\beta_2) \mathbb{E}^{|\mathcal{F}_t} |\mathbf{g}_{t,l}|^2 + (1-\beta_2) \sigma^2}} \mathbf{1}_{|\mathbf{G}_{t,l}| \geq \sigma} \\
& \geq \mathbb{E}^{|\mathcal{F}_t} \frac{\frac{1}{3} |\mathbf{g}_{t,l}|^2 + \frac{1-\beta_2}{3} \sigma^2}{\sqrt{\beta_2 \nu_{t-1,l} + (1-\beta_2) |\mathbf{g}_{t,l}|^2 + (1-\beta_2) \sigma^2}} \mathbf{1}_{|\mathbf{G}_{t,l}| \geq \sigma}.
\end{aligned}$$

397 Here the last inequality is due to the concavity of $\frac{x}{\sqrt{x+a}}$ with respect to x . As a conclusion,

$$\begin{aligned} \sum_{t=1}^T \mathbb{E} \left[\frac{|\mathbf{G}_{t,l}|^2}{\sqrt{\tilde{\nu}_{t,l}^1}} \mathbf{1}_{|\mathbf{G}_{t,l}| \geq \sigma} \right] &\geq \sum_{t=1}^T \mathbb{E} \left[\frac{\left(\frac{1}{3} |\mathbf{g}_{t,l}|^2 + \frac{1-\beta_2}{3} \sigma^2 \right)}{\sqrt{\beta_2 \boldsymbol{\nu}_{t-1,l} + (1-\beta_2) |\mathbf{g}_{t,l}|^2 + (1-\beta_2) \sigma^2}} \mathbf{1}_{|\mathbf{G}_{t,l}| \geq \sigma} \right] \\ &\geq \frac{1}{3(1-\beta_2)} \sum_{t=1}^T \mathbb{E} \left(\sqrt{\boldsymbol{\nu}_{t,l} + (1-\beta_2) \sigma^2} - \sqrt{\beta_2 (\boldsymbol{\nu}_{t-1,l} + (1-\beta_2) \sigma^2)} \right) \mathbf{1}_{|\mathbf{G}_{t,l}| \geq \sigma}. \end{aligned}$$

398 On the other hand, as stated in Section 4.2, we define $\{\bar{\nu}_{t,l}\}_{t=0}^\infty$ as $\bar{\nu}_{0,l} = \boldsymbol{\nu}_{0,l}$, $\bar{\nu}_{t,l} = \bar{\nu}_{t-1,l} +$
399 $|g_{t,l}|^2 \mathbf{1}_{|\mathbf{G}_{t,l}| < \sigma}$. One can easily observe that $\bar{\nu}_{t,l} \leq \boldsymbol{\nu}_{t,l}$, and thus

$$\begin{aligned} &\sum_{t=1}^T \mathbb{E} \left(\sqrt{\boldsymbol{\nu}_{t,l} + (1-\beta_2) \sigma^2} - \sqrt{\beta_2 (\boldsymbol{\nu}_{t-1,l} + (1-\beta_2) \sigma^2)} \right) \mathbf{1}_{|\mathbf{G}_{t,l}| < \sigma} \\ &= \sum_{t=1}^T \mathbb{E} \left(\sqrt{\beta_2 \boldsymbol{\nu}_{t-1,l} + (1-\beta_2) |g_{t,l}|^2 + (1-\beta_2) \sigma^2} - \sqrt{\beta_2 (\boldsymbol{\nu}_{t-1,l} + (1-\beta_2) \sigma^2)} \right) \mathbf{1}_{|\mathbf{G}_{t,l}| < \sigma} \\ &\leq \sum_{t=1}^T \mathbb{E} \left(\sqrt{\beta_2 \bar{\nu}_{t-1,l} + (1-\beta_2) |g_{t,l}|^2 + (1-\beta_2) \sigma^2} - \sqrt{\beta_2 (\bar{\nu}_{t-1,l} + (1-\beta_2) \sigma^2)} \right) \mathbf{1}_{|\mathbf{G}_{t,l}| < \sigma} \\ &\leq \sum_{t=1}^T \mathbb{E} \left(\sqrt{\beta_2 \bar{\nu}_{t-1,l} + (1-\beta_2) |g_{t,l}|^2 \mathbf{1}_{|\mathbf{G}_{t,l}| < \sigma} + (1-\beta_2) \sigma^2} - \sqrt{\beta_2 (\bar{\nu}_{t-1,l} + (1-\beta_2) \sigma^2)} \right) \\ &= \sum_{t=1}^T \mathbb{E} \left(\sqrt{\bar{\nu}_{t,l} + (1-\beta_2) \sigma^2} - \sqrt{\beta_2 (\bar{\nu}_{t-1,l} + (1-\beta_2) \sigma^2)} \right) \\ &= \mathbb{E} \sqrt{\bar{\nu}_{T,l} + (1-\beta_2) \sigma^2} + (1-\sqrt{\beta_2}) \sum_{t=1}^{T-1} \mathbb{E} \sqrt{\bar{\nu}_{t,l} + (1-\beta_2) \sigma^2} - \mathbb{E} \sqrt{\beta_2 (\bar{\nu}_{0,l} + (1-\beta_2) \sigma^2)}. \end{aligned}$$

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

$$\begin{aligned} &\mathbb{E} \sqrt{\boldsymbol{\nu}_{T,l} + (1-\beta_2) \sigma^2} + (1-\sqrt{\beta_2}) \sum_{t=1}^{T-1} \mathbb{E} \sqrt{\boldsymbol{\nu}_{t,l} + (1-\beta_2) \sigma^2} - \mathbb{E} \sqrt{\beta_2 (\boldsymbol{\nu}_{0,l} + (1-\beta_2) \sigma^2)} \\ &= \sum_{t=1}^T \mathbb{E} \left(\sqrt{\boldsymbol{\nu}_{t,l} + (1-\beta_2) \sigma^2} - \sqrt{\beta_2 (\boldsymbol{\nu}_{t-1,l} + (1-\beta_2) \sigma^2)} \right) \\ &\leq \sum_{t=1}^T \mathbb{E} \left(\sqrt{\boldsymbol{\nu}_{t,l} + (1-\beta_2) \sigma^2} - \sqrt{\beta_2 (\boldsymbol{\nu}_{t-1,l} + (1-\beta_2) \sigma^2)} \right) \mathbf{1}_{|\mathbf{G}_{t,l}| \geq \sigma} \\ &\quad + \sum_{t=1}^T \mathbb{E} \left(\sqrt{\boldsymbol{\nu}_{t,l} + (1-\beta_2) \sigma^2} - \sqrt{\beta_2 (\boldsymbol{\nu}_{t-1,l} + (1-\beta_2) \sigma^2)} \right) \mathbf{1}_{|\mathbf{G}_{t,l}| < \sigma} \\ &= 3(1-\beta_2) \sum_{t=1}^T \mathbb{E} \left[\frac{|\mathbf{G}_{t,l}|^2}{\sqrt{\tilde{\nu}_{t,l}^1}} \right] + \mathbb{E} \sqrt{\bar{\nu}_{T,l} + (1-\beta_2) \sigma^2} + (1-\sqrt{\beta_2}) \sum_{t=1}^{T-1} \mathbb{E} \sqrt{\bar{\nu}_{t,l} + (1-\beta_2) \sigma^2} - \mathbb{E} \sqrt{\beta_2 (\bar{\nu}_{0,l} + (1-\beta_2) \sigma^2)}. \end{aligned}$$

401 As $\mathbb{E} \sqrt{\boldsymbol{\nu}_{T,l} + (1-\beta_2) \sigma^2} \geq \mathbb{E} \sqrt{\bar{\nu}_{T,l} + (1-\beta_2) \sigma^2}$ and $\mathbb{E} \sqrt{\boldsymbol{\nu}_{0,l} + (1-\beta_2) \sigma^2} =$
402 $\mathbb{E} \sqrt{\bar{\nu}_{0,l} + (1-\beta_2) \sigma^2}$, we obtain that

$$\begin{aligned} (1-\sqrt{\beta_2}) \sum_{t=1}^T \mathbb{E} \sqrt{\boldsymbol{\nu}_{t,l} + (1-\beta_2) \sigma^2} &\leq 3(1-\beta_2) \sum_{t=1}^T \mathbb{E} \left[\frac{|\mathbf{G}_{t,l}|^2}{\sqrt{\tilde{\nu}_{t,l}^1}} \right] + (1-\sqrt{\beta_2}) \sum_{t=1}^T \mathbb{E} \sqrt{\bar{\nu}_{t,l} + (1-\beta_2) \sigma^2} \\ &\leq 3(1-\beta_2) \sum_{t=1}^T \mathbb{E} \left[\frac{|\mathbf{G}_{t,l}|^2}{\sqrt{\tilde{\nu}_{t,l}^1}} \right] + (1-\sqrt{\beta_2}) \sum_{t=1}^T \sqrt{\mathbb{E} \bar{\nu}_{t,l} + (1-\beta_2) \sigma^2} \\ &\leq 3(1-\beta_2) \sum_{t=1}^T \mathbb{E} \left[\frac{|\mathbf{G}_{t,l}|^2}{\sqrt{\tilde{\nu}_{t,l}^1}} \right] + (1-\sqrt{\beta_2}) \sum_{t=1}^T \sqrt{\boldsymbol{\nu}_{0,l} + (1-\beta_2) \sigma^2}. \end{aligned} \tag{9}$$

403 Leveraging Eq. (4), we then obtain that

$$\begin{aligned}
& \sum_{t=1}^T \sum_{l=1}^d \mathbb{E} \sqrt{\boldsymbol{\nu}_{t,l} + (1 - \beta_2)\sigma^2} \\
& \leq 3(1 + \sqrt{\beta_2}) \sum_{t=1}^T \mathbb{E} \left[\frac{|\mathbf{G}_{t,l}|^2}{\sqrt{\tilde{\boldsymbol{\nu}}_{t,l}^1}} \right] + \sum_{t=1}^T \sum_{l=1}^d \sqrt{\boldsymbol{\nu}_{0,l} + (3 - \beta_2)\sigma^2} \\
& \leq \frac{12}{(1 - \beta_1)\eta} \left(f(\mathbf{w}_1) + 2 \sum_{l=1}^d C_1 \left(\mathbb{E} \ln \left(\frac{\sqrt{\boldsymbol{\nu}_{t,l} + (1 - \beta_2)\sigma^2}}{\boldsymbol{\nu}_{0,l}} \right) - T \ln \beta_2 \right) \right) + T \sum_{l=1}^d \sqrt{\boldsymbol{\nu}_{0,l} + (3 - \beta_2)\sigma^2} \\
& \leq \frac{12}{(1 - \beta_1)\eta} \left(f(\mathbf{w}_1) + 2 \sum_{l=1}^d C_1 \left(\mathbb{E} \ln \left(\frac{\sum_{t=1}^T \sqrt{\boldsymbol{\nu}_{t,l} + (1 - \beta_2)\sigma^2}}{\boldsymbol{\nu}_{0,l}} \right) - T \ln \beta_2 \right) \right) + T \sum_{l=1}^d \sqrt{\boldsymbol{\nu}_{0,l} + (3 - \beta_2)\sigma^2} \\
& \leq \frac{12}{(1 - \beta_1)\eta} \left(f(\mathbf{w}_1) + 2 \sum_{l=1}^d C_1 \left(\ln \left(\frac{\mathbb{E} \sum_{t=1}^T \sum_{m=1}^d \sqrt{\boldsymbol{\nu}_{t,m} + (1 - \beta_2)\sigma^2}}{\boldsymbol{\nu}_{0,l}} \right) - T \ln \beta_2 \right) \right) + T \sum_{l=1}^d \sqrt{\boldsymbol{\nu}_{0,l} + (3 - \beta_2)\sigma^2},
\end{aligned}$$

404 where in the last inequality we use the concavity of $h(x) = \ln x$. Solving the above inequality with
405 respect to $\sum_{t=1}^T \sum_{l=1}^d \mathbb{E} \sqrt{\boldsymbol{\nu}_{t,l} + (1 - \beta_2)\sigma^2}$ then gives

$$\begin{aligned}
\sum_{t=1}^T \sum_{l=1}^d \mathbb{E} \sqrt{\boldsymbol{\nu}_{t,l} + (1 - \beta_2)\sigma^2} & \leq 2T \sum_{l=1}^d \sqrt{\boldsymbol{\nu}_{0,l} + (3 - \beta_2)\sigma^2} + 4dC_1 \ln dC_1 \\
& \quad + \frac{24}{(1 - \beta_1)\eta} \left(f(\mathbf{w}_1) + 2 \sum_{l=1}^d C_1 \left(\ln \left(\frac{1}{\boldsymbol{\nu}_{0,l}} \right) - T \ln \beta_2 \right) \right).
\end{aligned}$$

406 The proof is then completed.

407 □

408 C.2 Proof of Theorem 1

409 *Proof of Theorem 1.* As stated in Section 4.2, the proof involves solving two key challenges. We
410 respectively divide the proof into two stages according to the challenges.

411 **Stage I.** Based on Lemma 1, we can estimate $\mathbb{E} \langle \mathbf{G}_t, \frac{1}{\sqrt{\tilde{\boldsymbol{\nu}}_t}} \odot \mathbf{m}_t \rangle = F_t^0$ recursively. Specifically, we
412 have

$$\begin{aligned}
F_t^0 & \geq \sum_{i=0}^{t-1} \beta_1^i \left(\frac{(1 - \beta_1)}{2} \mathbb{E} \left[\left\| \frac{1}{\sqrt[4]{\tilde{\boldsymbol{\nu}}_t^{i+1}}} \odot \mathbf{G}_{t-i} \right\|^2 \right] - \beta_1 \mathbb{E} \left[\left\| \mathbf{w}_{t-i} - \mathbf{w}_{t-i-1} \right\| \left\| \frac{1}{\sqrt{\tilde{\boldsymbol{\nu}}_t^{i+1}}} \odot \mathbf{m}_{t-i-1} \right\|^2 \right] \right. \\
& \quad \left. - \left(2 \frac{\sqrt{1 - \beta_2}}{1 - \beta_1} \sigma + L^2 \frac{\eta^2(1 - \beta_1)}{(1 - \beta_2)^{\frac{1}{2}}(1 - \frac{\beta_1^2}{\beta_2})\beta_2^i \sigma} d \right) \mathbb{E} \left\| \frac{1}{\sqrt{\boldsymbol{\nu}_{t-i}}} \odot \mathbf{m}_{t-i} \right\|^2 \right) \\
& \geq \frac{(1 - \beta_1)}{2} \mathbb{E} \left[\left\| \frac{1}{\sqrt[4]{\tilde{\boldsymbol{\nu}}_t^1}} \odot \mathbf{G}_t \right\|^2 \right] - \sum_{i=0}^{t-1} \beta_1^i \left(\beta_1 \mathbb{E} \left[\left\| \mathbf{w}_{t-i} - \mathbf{w}_{t-i-1} \right\| \left\| \frac{1}{\sqrt{\tilde{\boldsymbol{\nu}}_t^{i+1}}} \odot \mathbf{m}_{t-i-1} \right\|^2 \right] \right. \\
& \quad \left. + \left(2 \frac{\sqrt{1 - \beta_2}}{1 - \beta_1} \sigma + L^2 \frac{\eta^2(1 - \beta_1)}{(1 - \beta_2)^{\frac{1}{2}}(1 - \frac{\beta_1^2}{\beta_2})\beta_2^i \sigma} d \right) \mathbb{E} \left\| \frac{1}{\sqrt{\boldsymbol{\nu}_{t-i}}} \odot \mathbf{m}_{t-i} \right\|^2 \right)
\end{aligned}$$

413 Applying the above inequality back to Eq. (3) then gives

$$\begin{aligned} & \mathbb{E}f(\mathbf{w}_{t+1}) \\ & \leq \mathbb{E}f(\mathbf{w}_t) - \frac{(1-\beta_1)\eta}{2} \mathbb{E} \left[\left\| \frac{1}{\sqrt[4]{\tilde{\nu}_t^1}} \odot \mathbf{G}_t \right\|^2 \right] + \frac{L}{2} \eta^2 \mathbb{E} \left\| \frac{1}{\sqrt{\nu_t}} \odot \mathbf{m}_t \right\|^2 + \eta \sum_{i=0}^{t-1} \beta_1^i (\beta_1 \mathbb{E} [\|\mathbf{w}_{t-i} - \mathbf{w}_{t-i-1}\| \\ & \quad \times \left\| \frac{1}{\sqrt{\tilde{\nu}_t^{i+1}}} \odot \mathbf{m}_{t-i-1} \right\|] + \left(2 \frac{\sqrt{1-\beta_2}}{1-\beta_1} \sigma + L^2 \frac{\eta^2(1-\beta_1)}{(1-\beta_2)^{\frac{1}{2}}(1-\frac{\beta_1^2}{\beta_2})\beta_2^i} \sigma \right) \mathbb{E} \left\| \frac{1}{\sqrt{\nu_{t-i}}} \odot \mathbf{m}_{t-i} \right\|^2 \right). \end{aligned}$$

414 Summing the above inequality with respect to t then gives

$$\begin{aligned} & \mathbb{E}f(\mathbf{w}_{T+1}) \\ & \leq f(\mathbf{w}_1) - \sum_{t=1}^T \frac{(1-\beta_1)\eta}{2} \mathbb{E} \left[\left\| \frac{1}{\sqrt{\tilde{\nu}_t^1}} \odot \mathbf{G}_t \right\|^2 \right] + \left(\frac{L}{2} \eta^2 + 2 \frac{\sqrt{1-\beta_2}}{(1-\beta_1)^2} \eta \sigma + \frac{\eta^2 \beta_1}{\sqrt{\beta_2}(1-\frac{\beta_1}{\sqrt{\beta_2}})} \right. \\ & \quad \left. + L^2 \frac{\beta_1 \eta^3 (1-\beta_1)}{\beta_2(1-\beta_2)^{\frac{1}{2}}(1-\frac{\beta_1^2}{\beta_2})(1-\frac{\beta_1}{\beta_2})^2} \frac{d}{\sigma} \right) \sum_{t=1}^T \mathbb{E} \left\| \frac{1}{\sqrt{\nu_t}} \odot \mathbf{m}_t \right\|^2. \end{aligned}$$

415 Here the inequality is due to

$$\begin{aligned} 2 \frac{\sqrt{1-\beta_2}}{1-\beta_1} \eta \sum_{t=1}^T \sum_{i=0}^{t-1} \beta_1^i \mathbb{E} \sigma \left\| \frac{1}{\sqrt{\nu_{t-i}}} \odot \mathbf{m}_{t-i} \right\|^2 & = 2 \frac{\sqrt{1-\beta_2}}{1-\beta_1} \eta \sigma \sum_{i=1}^T \sum_{t=i}^T \beta_1^{t-i} \mathbb{E} \left\| \frac{1}{\sqrt{\nu_t}} \odot \mathbf{m}_t \right\|^2 \\ & \leq 2 \frac{\sqrt{1-\beta_2}}{(1-\beta_1)^2} \eta \sigma \sum_{i=1}^T \mathbb{E} \left\| \frac{1}{\sqrt{\nu_i}} \odot \mathbf{m}_i \right\|^2, \end{aligned}$$

416

$$\begin{aligned} & \eta \sum_{t=1}^T \sum_{i=0}^{t-1} \beta_1^{i+1} \mathbb{E} \left[\|\mathbf{w}_{t-i} - \mathbf{w}_{t-i-1}\| \left\| \frac{1}{\sqrt{\tilde{\nu}_t^{i+1}}} \odot \mathbf{m}_{t-i-1} \right\| \right] \\ & \leq \eta \sum_{t=1}^T \sum_{i=0}^{t-1} \frac{\beta_1^{i+1}}{\sqrt{\beta_2^{i+1}}} \mathbb{E} \left[\|\mathbf{w}_{t-i} - \mathbf{w}_{t-i-1}\| \left\| \frac{1}{\sqrt{\nu_{t-i-1}}} \odot \mathbf{m}_{t-i-1} \right\| \right] \\ & = \eta^2 \sum_{t=1}^T \sum_{i=0}^{t-1} \frac{\beta_1^{i+1}}{\sqrt{\beta_2^{i+1}}} \mathbb{E} \left[\left\| \frac{1}{\sqrt{\nu_{t-i-1}}} \odot \mathbf{m}_{t-i-1} \right\|^2 \right] = \eta^2 \sum_{i=0}^{T-1} \sum_{t=i+1}^T \frac{\beta_1^{t-i}}{\sqrt{\beta_2^{t-i}}} \mathbb{E} \left[\left\| \frac{1}{\sqrt{\nu_i}} \odot \mathbf{m}_i \right\|^2 \right] \\ & \leq \frac{\eta^2 \beta_1}{\sqrt{\beta_2}(1-\frac{\beta_1}{\sqrt{\beta_2}})} \sum_{i=0}^{T-1} \mathbb{E} \left[\left\| \frac{1}{\sqrt{\nu_i}} \odot \mathbf{m}_i \right\|^2 \right] = \frac{\eta^2 \beta_1}{\sqrt{\beta_2}(1-\frac{\beta_1}{\sqrt{\beta_2}})} \sum_{i=1}^{T-1} \mathbb{E} \left[\left\| \frac{1}{\sqrt{\nu_i}} \odot \mathbf{m}_i \right\|^2 \right], \end{aligned}$$

417 and

$$\begin{aligned} & L^2 \frac{\eta^3(1-\beta_1)}{(1-\beta_2)^{\frac{1}{2}}(1-\frac{\beta_1^2}{\beta_2})} \frac{d}{\sigma} \sum_{t=1}^T \sum_{i=0}^{t-1} \frac{\beta_1^i}{\beta_2^i} i \mathbb{E} \left\| \frac{1}{\sqrt{\nu_{t-i}}} \odot \mathbf{m}_{t-i} \right\|^2 \\ & = L^2 \frac{\eta^3(1-\beta_1)}{(1-\beta_2)^{\frac{1}{2}}(1-\frac{\beta_1^2}{\beta_2})} \frac{d}{\sigma} \sum_{i=1}^T \sum_{t=i}^T \frac{\beta_1^{t-i}}{\beta_2^{t-i}} (t-i) \mathbb{E} \left\| \frac{1}{\sqrt{\nu_i}} \odot \mathbf{m}_i \right\|^2 \leq L^2 \frac{\beta_1 \eta^3(1-\beta_1)}{\beta_2(1-\beta_2)^{\frac{1}{2}}(1-\frac{\beta_1^2}{\beta_2})(1-\frac{\beta_1}{\beta_2})^2} \frac{d}{\sigma} \sum_{i=1}^T \mathbb{E} \left\| \frac{1}{\sqrt{\nu_i}} \odot \mathbf{m}_i \right\|^2. \end{aligned}$$

418 Applying Lemma 4, we obtain that

$$\begin{aligned} & \mathbb{E}f(\mathbf{w}_{T+1}) \\ & \leq f(\mathbf{w}_1) + \sum_{t=1}^d \left(\frac{L}{2} \eta^2 + 2 \frac{\sqrt{1-\beta_2}}{(1-\beta_1)^2} \eta \sigma + \frac{\eta^2 \beta_1}{\sqrt{\beta_2}(1-\frac{\beta_1}{\sqrt{\beta_2}})} + L^2 \frac{\beta_1 \eta^3(1-\beta_1)}{\beta_2(1-\beta_2)^{\frac{1}{2}}(1-\frac{\beta_1^2}{\beta_2})(1-\frac{\beta_1}{\beta_2})^2} \frac{d}{\sigma} \frac{(1-\beta_1)^2}{(1-\frac{\beta_1}{\sqrt{\beta_2}})^2} \right) \frac{1}{1-\beta_2} \\ & \quad \times \left(\mathbb{E} \ln \left(\frac{\nu_{T,l}}{\nu_{0,l}} \right) - T \ln \beta_2 \right) - \sum_{t=1}^T \frac{(1-\beta_1)\eta}{2} \mathbb{E} \left[\left\| \frac{1}{\sqrt{\tilde{\nu}_t^1}} \odot \mathbf{G}_t \right\|^2 \right]. \end{aligned}$$

419 The proof of Stage I is completed.

420 **Stage II.** According to Cauchy's inequality, we have

$$\left(\mathbb{E} \sum_{t=1}^T \|\mathbf{G}_t\|_1 \right)^2 \leq \left(\sum_{t=1}^T \mathbb{E} \left[\left\| \frac{1}{\sqrt[4]{\tilde{\nu}_t^1}} \odot \mathbf{G}_t \right\|^2 \right] \right) \left(\sum_{t=1}^T \mathbb{E} \left[\left\| \sqrt[4]{\tilde{\nu}_t^1} \right\|^2 \right] \right). \quad (10)$$

421 Meanwhile, by Lemma 2, we have

$$\begin{aligned} \sum_{t=1}^T \mathbb{E} \left[\left\| \sqrt[4]{\tilde{\nu}_t^1} \right\|^2 \right] &= \mathbb{E} \left[\sum_{t=1}^T \sum_{l=1}^d \sqrt{\beta_2 \nu_{t-1,l} + (1-\beta_2) |\mathbf{G}_{t,l}|^2 + (1-\beta_2) \sigma^2} \right] \\ &\leq \mathbb{E} \left[\sum_{t=1}^T \sum_{l=1}^d \left(\sqrt{\beta_2 \nu_{t-1,l} + (1-\beta_2) \sigma^2} + \sqrt{1-\beta_2} |\mathbf{G}_{t,l}| \right) \right] \\ &= \mathbb{E} \left[\sum_{t=1}^T \sum_{l=1}^d \sqrt{\beta_2 \nu_{t-1,l} + (1-\beta_2) \sigma^2} + \sum_{t=1}^T \sqrt{1-\beta_2} \|\mathbf{G}_t\|_1 \right] \\ &\leq \mathbb{E} \left[\sum_{t=1}^T \sqrt{1-\beta_2} \|\mathbf{G}_t\|_1 \right] + 2T \sqrt{\nu_{0,l} + (3-\beta_2) \sigma^2} + 4dC_1 \ln dC_1 \\ &\quad + \frac{24}{(1-\beta_1)\eta} \left(f(\mathbf{w}_1) + 2 \sum_{l=1}^d C_1 \left(\ln \left(\frac{1}{\nu_{0,l}} \right) - T \ln \beta_2 \right) \right). \end{aligned}$$

422 Combining the above inequality and Eq. (10) gives

$$\begin{aligned} \left(\mathbb{E} \sum_{t=1}^T \|\mathbf{G}_t\|_1 \right)^2 &\leq \frac{2}{(1-\beta_1)\eta} \left(f(\mathbf{w}_1) + \sum_{l=1}^d C_1 \left(\mathbb{E} \ln \left(\frac{\nu_{T,l}}{\nu_{0,l}} \right) - T \ln \beta_2 \right) \right) \\ &\quad \times \left(\mathbb{E} \left[\sum_{t=1}^T \sqrt{1-\beta_2} \|\mathbf{G}_t\|_1 \right] + 2T \sqrt{\nu_{0,l} + (3-\beta_2) \sigma^2} + 4dC_1 \ln dC_1 \right. \\ &\quad \left. + \frac{24}{(1-\beta_1)\eta} \left(f(\mathbf{w}_1) + 2 \sum_{l=1}^d C_1 \left(\ln \left(\frac{1}{\nu_{0,l}} \right) - T \ln \beta_2 \right) \right) \right). \end{aligned}$$

423 Solving the above quadratic inequality with respect to $\mathbb{E} \sum_{t=1}^T \|\mathbf{G}_t\|_1$ then completes the proof.

424 □

425 D Proof of Theorem 2

426 *Proof.* According to Stage I in the proof of Theorem 1, we obtain

$$\begin{aligned} &\mathbb{E} f(\mathbf{w}_{T+1}) \\ &\leq f(\mathbf{w}_1) + \sum_{t=1}^d \left(\frac{L}{2} \eta^2 + 2 \frac{\sqrt{1-\beta_2}}{(1-\beta_1)^2} \eta \sigma + \frac{\eta^2 \beta_1}{\sqrt{\beta_2} (1-\frac{\beta_1}{\sqrt{\beta_2}})} + L^2 \frac{\beta_1 \eta^3 (1-\beta_1)}{\beta_2 (1-\beta_2)^{\frac{1}{2}} (1-\frac{\beta_1^2}{\beta_2}) (1-\frac{\beta_1}{\sqrt{\beta_2}})^2} \frac{d}{\sigma} \frac{(1-\beta_1)^2}{(1-\frac{\beta_1}{\sqrt{\beta_2}})^2} \right) \frac{1}{1-\beta_2} \\ &\quad \times \mathbb{E} \left(\ln \left(\frac{\nu_{T,l}}{\nu_{0,l}} \right) - T \ln \beta_2 \right) - \sum_{t=1}^T \frac{(1-\beta_1)\eta}{2} \mathbb{E} \left[\left\| \frac{1}{\sqrt[4]{\tilde{\nu}_t^1}} \odot \mathbf{G}_t \right\|^2 \right]. \end{aligned}$$

427 Applying the definition of η , β_1 , and β_2 , we obtain that

$$\sum_{t=1}^T \mathbb{E} \left[\left\| \frac{1}{\sqrt[4]{\tilde{\nu}_t^1}} \odot \mathbf{G}_t \right\|^2 \right] \leq \frac{2\sqrt{T}}{\sqrt{b}} \left(D_1 + \frac{D_2}{d} \sum_{l=1}^d \mathbb{E} \ln \nu_{T,l} \right). \quad (11)$$

428 Meanshile, we have that

$$\begin{aligned}
& \frac{|\mathbf{G}_{t,l}|^2}{\sqrt{\tilde{\nu}_{t,l}^1}} \mathbf{1}_{|\mathbf{G}_{t,l}| \geq \sigma} \geq \frac{\frac{1}{2} \mathbb{E}^{|\mathcal{F}_t} |\mathbf{g}_{t,l}|^2}{\sqrt{\tilde{\nu}_{t,l}^1}} \mathbf{1}_{|\mathbf{G}_{t,l}| \geq \sigma} \\
& = \frac{\frac{1}{2} \mathbb{E}^{|\mathcal{F}_t} |\mathbf{g}_{t,l}|^2}{\sqrt{\beta_2 \nu_{t-1,l} + (1-\beta_2) \mathbb{E}^{|\mathcal{F}_t} |\mathbf{g}_{t,l}|^2 + (1-\beta_2) \sigma^2}} \mathbf{1}_{|\mathbf{G}_{t,l}| \geq \sigma} \\
& \geq \frac{1}{2} \mathbb{E}^{|\mathcal{F}_t} \frac{|\mathbf{g}_{t,l}|^2}{\sqrt{\beta_2 \nu_{t-1,l} + (1-\beta_2) |\mathbf{g}_{t,l}|^2 + (1-\beta_2) \sigma^2}} \mathbf{1}_{|\mathbf{G}_{t,l}| \geq \sigma} \\
& \geq \frac{1}{2\sqrt{1-\beta_2}} \mathbb{E}^{|\mathcal{F}_t} \frac{|\mathbf{g}_{t,l}|^2}{\sqrt{\frac{\nu_{0,l}}{1-\beta_2} + \sum_{s=1}^T |g_{s,l}|^2 + \sigma^2}} \mathbf{1}_{|\mathbf{G}_{t,l}| \geq \sigma},
\end{aligned}$$

429 where the last inequality is due to that

$$\begin{aligned}
& \beta_2 \nu_{t-1,l} + (1-\beta_2) |\mathbf{g}_{t,l}|^2 = (1-\beta_2) \sum_{s=1}^t \beta_2^{t-s} |g_{s,l}|^2 + \beta_2^t \nu_{0,l} \\
& \leq (1-\beta_2) \sum_{s=1}^T |g_{s,l}|^2 + \nu_{0,l}. \tag{12}
\end{aligned}$$

430 Furthermore, we have

$$\begin{aligned}
& \frac{\sigma^2 + \frac{\nu_{0,l}}{1-\beta_2}}{\sqrt{\frac{\nu_{0,l}}{1-\beta_2} + \sum_{s=1}^T |g_{s,l}|^2 + \sigma^2}} + \sum_{t=1}^T \mathbb{E} \frac{|\mathbf{g}_{t,l}|^2}{\sqrt{\frac{\nu_{0,l}}{1-\beta_2} + \sum_{s=1}^T |g_{s,l}|^2 + \sigma^2}} \mathbf{1}_{|\mathbf{G}_{t,l}| < \sigma} \\
& \leq \frac{\sigma^2 + \frac{\nu_{0,l}}{1-\beta_2}}{\sqrt{\frac{\nu_{0,l}}{1-\beta_2} + \sum_{s=1}^T |g_{s,l}|^2 + \sigma^2}} + \sum_{t=1}^T \mathbb{E} \frac{|\mathbf{g}_{t,l}|^2}{\sqrt{\frac{\nu_{0,l}}{1-\beta_2} + \sum_{s=1}^T |g_{s,l}|^2} \mathbf{1}_{|\mathbf{G}_{s,l}| < \sigma} + \sigma^2}} \mathbf{1}_{|\mathbf{G}_{t,l}| < \sigma} \\
& = \mathbb{E} \sqrt{\frac{\nu_{0,l}}{1-\beta_2} + \sum_{s=1}^T |g_{s,l}|^2 \mathbf{1}_{|\mathbf{G}_{s,l}| < \sigma} + \sigma^2} \leq \sqrt{\frac{\nu_{0,l}}{1-\beta_2} + \mathbb{E} \sum_{s=1}^T |g_{s,l}|^2 \mathbf{1}_{|\mathbf{G}_{s,l}| < \sigma} + \sigma^2} \\
& \leq \sqrt{\frac{\nu_{0,l}}{1-\beta_2} + 2\sigma^2 T + \sigma^2}.
\end{aligned}$$

431 Conclusively, we obtain

$$\begin{aligned}
& \mathbb{E} \sqrt{\frac{\nu_{0,l}}{1-\beta_2} + \sum_{s=1}^T |g_{s,l}|^2 + \sigma^2} \\
& = \frac{\sigma^2 + \frac{\nu_{0,l}}{1-\beta_2}}{\sqrt{\frac{\nu_{0,l}}{1-\beta_2} + \sum_{s=1}^T |g_{s,l}|^2 + \sigma^2}} + \sum_{t=1}^T \mathbb{E} \frac{|\mathbf{g}_{t,l}|^2}{\sqrt{\frac{\nu_{0,l}}{1-\beta_2} + \sum_{s=1}^T |g_{s,l}|^2 + \sigma^2}} \mathbf{1}_{|\mathbf{G}_{t,l}| < \sigma} \\
& \quad + \sum_{t=1}^T \mathbb{E} \frac{|\mathbf{g}_{t,l}|^2}{\sqrt{\frac{\nu_{0,l}}{1-\beta_2} + \sum_{s=1}^T |g_{s,l}|^2 + \sigma^2}} \mathbf{1}_{|\mathbf{G}_{t,l}| \geq \sigma} \\
& \leq \sqrt{\frac{\nu_{0,l}}{1-\beta_2} + 2\sigma^2 T + \sigma^2} + 2\sqrt{1-\beta_2} \sum_{t=1}^T \frac{|\mathbf{G}_{t,l}|^2}{\sqrt{\tilde{\nu}_{t,l}^1}} \mathbf{1}_{|\mathbf{G}_{t,l}| \geq \sigma}.
\end{aligned}$$

432 Summing the above inequality with respect to l then gives

$$\begin{aligned}
& \sum_{l=1}^d \mathbb{E} \sqrt{\frac{\nu_{0,l}}{1-\beta_2} + \sum_{s=1}^T |g_{s,l}|^2 + \sigma^2} \\
& \leq \sum_{l=1}^d \sqrt{\frac{\nu_{0,l}}{1-\beta_2} + 2\sigma^2 T + \sigma^2} + 2\sqrt{1-\beta_2} \sum_{l=1}^d \sum_{t=1}^T \frac{|\mathbf{G}_{t,l}|^2}{\sqrt{\tilde{\nu}_{t,l}^1}} \mathbf{1}_{|\mathbf{G}_{t,l}| \geq \sigma} \\
& \leq \sum_{l=1}^d \sqrt{\frac{\nu_{0,l}}{1-\beta_2} + 2\sigma^2 T + \sigma^2} \\
& \quad + \frac{4\sqrt{b}}{a(1-c)} f(\mathbf{w}_1) + \sum_{l=1}^d \frac{2}{ab\sqrt{b}} \left(La^2 + 4\frac{a\sqrt{b}\sigma}{(1-c)^2} + 2\frac{a^2c}{1-c} + 2\frac{L^2ca^3d}{\sqrt{b}(1-c)^5\sigma} \right) \left(\mathbb{E} \ln \left(\frac{\nu_{T,l}}{\nu_{0,l}} \right) + b \right) \\
& = \sum_{l=1}^d \sqrt{\frac{\nu_{0,l}}{1-\beta_2} + 2\sigma^2 T + \sigma^2} + \sum_{l=1}^d \frac{2}{ab\sqrt{b}} \left(La^2 + 4\frac{a\sqrt{b}\sigma}{(1-c)^2} + 2\frac{a^2c}{1-c} + 4\frac{L^2ca^3d}{\sqrt{b}(1-c)^5\sigma} \right) \mathbb{E} \ln(\sqrt{\nu_{T,l}}) \\
& \quad + \frac{4\sqrt{b}}{a(1-c)} f(\mathbf{w}_1) + \sum_{l=1}^d \frac{2}{ab\sqrt{b}} \left(La^2 + 4\frac{a\sqrt{b}\sigma}{(1-c)^2} + 2\frac{a^2c}{1-c} + 2\frac{L^2ca^3d}{\sqrt{b}(1-c)^5\sigma} \right) (-\ln(\nu_{0,l}) + b) \\
& \leq \sum_{l=1}^d \sqrt{\frac{\nu_{0,l}}{1-\beta_2} + 2\sigma^2 T + \sigma^2} \\
& \quad + d \frac{2}{ab\sqrt{b}} \left(La^2 + 4\frac{a\sqrt{b}\sigma}{(1-c)^2} + 2\frac{a^2c}{1-c} + 4\frac{L^2ca^3d}{\sqrt{b}(1-c)^5\sigma} \right) \mathbb{E} \ln \left(\sum_{l=1}^d \sqrt{1-\beta_2} \sqrt{\frac{\nu_{0,l}}{1-\beta_2} + \sum_{s=1}^T |g_{s,l}|^2 + \sigma^2} \right) \\
& \quad + \frac{4\sqrt{b}}{a(1-c)} f(\mathbf{w}_1) + \sum_{l=1}^d \frac{2}{ab\sqrt{b}} \left(La^2 + 4\frac{a\sqrt{b}\sigma}{(1-c)^2} + 2\frac{a^2c}{1-c} + 2\frac{L^2ca^3d}{\sqrt{b}(1-c)^5\sigma} \right) (-\ln(\nu_{0,l}) + b) \\
& \leq \sum_{l=1}^d \sqrt{\frac{\nu_{0,l}}{1-\beta_2} + 3\sigma^2 T + D_1} + D_2 \ln \left(\mathbb{E} \sum_{l=1}^d \sqrt{1-\beta_2} \sqrt{\frac{\nu_{0,l}}{1-\beta_2} + \sum_{s=1}^T |g_{s,l}|^2 + \sigma^2} \right),
\end{aligned}$$

433 where the second inequality is due to Eq. (11), the second-to-last inequality is due to Eq. (12),

434 and the last inequality is due to Jensen's inequality. Solving the above inequality with respect to

435 $\sqrt{1-\beta_2} \sum_{l=1}^d \mathbb{E} \sqrt{\frac{\nu_{0,l}}{1-\beta_2} + \sum_{s=1}^T |g_{s,l}|^2 + \sigma^2}$ then gives

$$\begin{aligned}
\sqrt{1-\beta_2} \sum_{l=1}^d \mathbb{E} \sqrt{\frac{\nu_{0,l}}{1-\beta_2} + \sum_{s=1}^T |g_{s,l}|^2 + \sigma^2} & \leq 2\sqrt{1-\beta_2} D_1 + 4\sqrt{1-\beta_2} D_2 \ln(1 + \sqrt{1-\beta_2} D_2) \\
& \quad + \sum_{l=1}^d \sqrt{\nu_{0,l} + 3b\sigma^2}.
\end{aligned}$$

436 Therefore, by Cauchy's inequality, we have

$$\mathbb{E} \left[\sum_{t=1}^T \|\mathbf{G}_t\|_1 \right]^2 \leq \left(\sum_{t=1}^T \mathbb{E} \left[\left\| \frac{1}{\sqrt{\tilde{\nu}_t^1}} \odot \mathbf{G}_t \right\|^2 \right] \right) \left(\sum_{t=1}^T \sum_{l=1}^d \mathbb{E} \sqrt{\tilde{\nu}_{t,l}^1} \right).$$

437 Since

$$\begin{aligned}
& \sum_{t=1}^T \sum_{l=1}^d \sqrt{\tilde{\nu}_{t,l}^1} \leq \sum_{t=1}^T \sum_{l=1}^d \left(\sqrt{\beta_2 \nu_{t-1,l} + (1-\beta_2)\sigma^2} + \sqrt{(1-\beta_2)} |\mathbf{G}_{t,l}| \right) \\
& \leq T \sum_{l=1}^d \sqrt{1-\beta_2} \sqrt{\frac{\nu_{0,l}}{1-\beta_2} + \sum_{s=1}^T |g_{s,l}|^2 + \sigma^2} + \sum_{t=1}^T \sum_{l=1}^d \sqrt{(1-\beta_2)} |\mathbf{G}_{t,l}| \\
& \leq T \left(2\sqrt{1-\beta_2} D_1 + 4\sqrt{1-\beta_2} D_2 \ln(1 + \sqrt{1-\beta_2} D_2) + \sum_{l=1}^d \sqrt{\nu_{0,l} + 3b\sigma^2} \right) + \sum_{t=1}^T \sqrt{(1-\beta_2)} \|\mathbf{G}_t\|_1,
\end{aligned}$$

438 we have

$$\begin{aligned}
& \mathbb{E} \left[\sum_{t=1}^T \|\mathbf{G}_t\|_1 \right]^2 \\
& \leq \left(T \left(2\sqrt{1-\beta_2} D_1 + 4\sqrt{1-\beta_2} D_2 \ln(1 + \sqrt{1-\beta_2} D_2) + \sum_{l=1}^d \sqrt{\nu_{0,l} + 3b\sigma^2} \right) + \sum_{t=1}^T \sqrt{(1-\beta_2)} \mathbb{E} \|\mathbf{G}_t\|_1 \right) \\
& \quad \times \frac{2\sqrt{T}}{\sqrt{b}} \left(D_1 + \frac{D_2}{d} \sum_{l=1}^d \mathbb{E} \ln \nu_{T,l} \right).
\end{aligned}$$

439 Solving the above inequality with respect to $\sum_{t=1}^T \mathbb{E} \|\mathbf{G}_t\|_1$ completes the proof. \square