# A COMPREHENSIVE FRAMEWORK FOR ANALYZING THE CONVERGENCE OF ADAM: BRIDGING THE GAP WITH SGD

Anonymous authors

Paper under double-blind review

#### Abstract

Adaptive Moment Estimation (Adam) is a cornerstone optimization algorithm in deep learning, widely recognized for its flexibility with adaptive learning rates and efficiency in handling large-scale data. However, despite its practical success, the theoretical understanding of Adam's convergence has been constrained by stringent assumptions, such as almost surely bounded stochastic gradients or uniformly bounded gradients, which are more restrictive than those typically required for analyzing stochastic gradient descent (SGD).

In this paper, we introduce a novel and comprehensive framework for analyzing the convergence properties of Adam. This framework offers a versatile approach to establishing Adam's convergence. Specifically, we prove that Adam achieves asymptotic (last iterate sense) convergence in both the almost sure sense and the  $L_1$  sense under the relaxed assumptions typically used for SGD, namely L-smoothness and the ABC inequality. Meanwhile, under the same assumptions, we show that Adam attains non-asymptotic sample complexity bounds similar to those of SGD.

006

008 009 010

011

013

014

015

016

017

018

019

021

024

025

## 1 INTRODUCTION

Adaptive Moment Estimation (Adam) is one of the most widely used optimization algorithms in
 deep learning due to its adaptive learning rate properties and efficiency in handling large-scale data
 (Kingma & Ba, 2014). Despite its widespread use, the theoretical understanding of Adam's convergence is not as advanced as its practical success. Previous studies have often imposed stringent
 assumptions on the loss function and stochastic gradients, such as uniformly bounded loss functions and almost surely bounded gradients (Reddi et al., 2018b; Zou & Shen, 2019), which are more
 restrictive than those required for analyzing classical stochastic gradient descent (SGD).

In this paper, we introduce a novel and comprehensive framework for analyzing the convergence properties of Adam. Our framework unifies various aspects of convergence analysis, including nonasymptotic (average iterate sense) sample complexity, asymptotic (last iterate sense) almost sure convergence, and asymptotic  $L_1$  convergence. Crucially, we demonstrate that under this framework, Adam can achieve convergence under the same assumptions typically used for SGD—namely, the *L*-smooth condition and the ABC inequality ( $L_2$  sense) (Khaled & Richtárik, 2023; Bottou, 2010; Ghadimi & Lan, 2013).

Several recent works have attempted to relax the stringent conditions required for Adam's convergence, each focusing on different aspects of the stochastic gradient assumptions and convergence guarantees. However, limitations still exist in terms of assumptions and the types of convergence results obtained. Table 1 provides the references and a summary of the works and compares the assumptions on stochastic gradients, the resulting complexities, and the convergence properties achieved.

Our approach builds upon these prior works and seeks to offer a more comprehensive and gen eral framework for analysis. In contrast to these previous works, we study Adam under the ABC
 inequality, which is more general and less restrictive compared to the assumptions made in the previous studies. Our analysis successfully establishes non-asymptotic sample complexity and achieves

054

Table 1: Comparison of Assumptions and Convergence Results. ( $\diamondsuit$ ) The smoothing term  $\mu$  is often set to small values like  $10^{-8}$  in practice. It is difficult and relevant to avoid the  $\mathcal{O}(\text{poly}(\frac{1}{\mu}))$ dependence (Wang et al., 2024a), which our analysis achieves. ( $\bigstar$ ) The work focuses on learning rates and hyperparameters dependent on the total number of iterations *T*, leading to results without a  $\mathcal{O}(\ln T)$  term. As our asymptotic analysis uses *T*-independent parameters, terms regarding  $\mathcal{O}(\ln T)$ inevitably appear, though our method can be easily extended to *T*-dependent settings. ( $\bigstar$ ) These works have weakened the classical *L*-smooth condition, which is different from the focus of this paper.

Reference	Assumptions on Stochastic Gradient	Sample Complexity	A.S. Convergence	L <sub>1</sub> Convergence
Wang et al. (2024a)♠	Bounded Variance (or Coordinate Weak Growth)	$O\left(\frac{1}{\sqrt{T}}\right)$	No	No
He et al. (2023) أ	Almost Surely Bounded Stochastic Gradient	$\mathcal{O}\left(\operatorname{poly}\left(\frac{1}{\mu}\right) \cdot \frac{\ln T}{\sqrt{T}}\right)$	Yes	Yes
Zou et al. (2019)	L2 Bounded Stochastic Gradient	$\mathcal{O}\left(\frac{\ln T}{\sqrt{T}}\right)$	No	No
Zhang et al. (2022)	Randomly Reshuffled Stochastic Gradient	$O\left(\frac{\ln T}{\sqrt{T}}\right)$	No	No
Li et al. (2024)♥◊	Almost Surely Bounded Stochastic Gradient or Sub-Gaussian Variance	$\mathcal{O}\left(\operatorname{poly}\left(\frac{1}{\mu}\right)\cdot\frac{\ln T}{\sqrt{T}}\right)$	No	No
Wang et al. (2024b)♥	Randomly Reshuffled Stochastic Gradient	$O\left(\frac{\ln T}{\sqrt{T}}\right)$	No	No
Xiao et al. (2024)♥	Almost Surely Bounded Stochastic Gradient	No Result	Yes	No
Our Work	ABC Inequality	$O\left(\frac{\ln T}{\sqrt{T}}\right)$	Yes	Yes

asymptotic almost sure convergence and  $L_1$  convergence under conditions that align with those required for SGD. This makes our framework theoretically sound and versatile for analyzing multiple convergence properties of Adam. Our framework might also be of independent interest in analyzing different variants of Adam. In summary, our work presents a novel and general theoretical framework for Adam, unifying various convergence properties. This framework demonstrates that Adam's convergence guarantees can be aligned with those of SGD, which justifies the applicability of Adam across a wide range of machine learning problems.

081 1.1 RELATED WORKS

In recent years, the convergence properties of Adam have been extensively studied, with various works focusing on different assumptions about stochastic gradients and the types of convergence guarantees provided. In the following discussion, we categorize and review key contributions based on the different types of stochastic gradient assumptions they employ, as summarized in Table 1.

**Bounded Variance and Coordinate Weak Growth:** Wang et al. (2024a) considered Adam's convergence under the assumption of bounded variance or coordinate weak growth. The coordinate weak growth condition (Eq. 2) is particularly stringent as it requires that each component of the stochastic gradient satisfies a weak growth inequality, which is stronger than the traditional weak growth condition (Eq. 1) applied to the entire gradient. Under these assumptions, Wang et al. were able to avoid the  $O(1/\mu)$  complexity. However, their work did not focus on analyzing almost sure convergence or  $L_1$  convergence, as the primary emphasis was on the sample complexity of the algorithm's behavior.

095

Almost Surely Bounded Stochastic Gradients: Several works, including He et al. (2023) and 096 Xiao et al. (2024), have explored Adam's convergence under the assumption that the stochastic gradients are almost surely bounded. This is a particularly strong assumption, as it implies several other 098 commonly made assumptions about stochastic gradients, such as bounded variance, weak growth, 099 coordinate weak growth, and sub-Gaussian properties. The assumption is often impractical in non-100 convex settings where gradients can become unbounded. Moreover, studies in Wang et al. (2023) 101 have highlighted that this assumption is unrealistic in many common machine learning frameworks, 102 failing to hold even for simple quadratic functions, let alone for deep neural networks. While these 103 works achieved almost sure convergence and, in some cases,  $L_1$  convergence, they did not address 104 the complexity related to the  $\mathcal{O}(1/\mu)$  term.

105

106  $L_2$  Bounded Stochastic Gradients: Zou et al. (2019) analyzed Adam under the assumption of 107  $L_2$  bounded stochastic gradients. Although this condition is milder than the almost surely bounded gradients assumption, it is still stronger than the traditional weak growth condition and the ABC inequality. In the standard analytical framework, this assumption can at best be weakened to the coordinate weak growth condition, which remains more restrictive than the assumptions typically considered for SGD. At the same time, this work focused on complexity analysis without addressing asymptotic convergence.

112 113

114

115

116

117

118

**Randomly Reshuffled Stochastic Gradients:** In other works, such as those by Zhang et al. (2022) and Wang et al. (2024b), the authors considered the case where the stochastic gradients are randomly reshuffled. Randomly reshuffled stochastic gradients represent a special case where the gradients are typically assumed to satisfy certain inequalities almost surely. This reliance on almost sure properties forms a much stronger and more restrictive analytical framework compared to those based on traditional weak growth conditions or the ABC inequality. While these works successfully avoided  $O(1/\mu)$  complexity, they did not focus on analyzing the asymptotic convergence property.

119 120 121

# 2 PRELIMINARIES

122 123 124

125

126

127

128

129 130 131

137

138

In this section, we introduce the necessary preliminaries and establish the foundational framework for our convergence analysis of the Adam. We begin by recalling the Adam optimization algorithm. We then state the assumptions that will be used throughout our analysis. These assumptions are standard in stochastic optimization and are crucial for deriving our main results. By laying out these assumptions explicitly, we also facilitate a clear comparison with the conditions used in previous works, highlighting the less restrictive nature of our approach.

2.1 Adam

Adam is an extension of SGD that computes adaptive learning rates for each parameter by utilizing
estimates of the first and second moments of the gradients. It combines the advantages of two other
extensions of SGD: AdaGrad, which works well with sparse gradients, and RMSProp, which works
well in online and non-stationary settings.

## Algorithm 1 Adam

139 **Input:** Stochastic oracle  $\mathcal{O}$ , initial learning rate  $\eta_t > 0$ , initial parameters  $w_1 \in \mathbb{R}^d$ , initial exponential moving averages  $m_0 = 0, v_0 = v \cdot \mathbf{1}^\top$  with v > 0, hyperparameters  $\beta_1, \beta_{2,1} \in [0, 1)$ , smoothing 140 141 term  $\mu > 0$ , number of iterations T 142 **Output:** Final parameter  $w_T$ 143 1: for t = 1 to T do 2: Generate conditioner parameter  $\beta_{2,t}$ ; 144 Sample a random data point  $z_t$  and compute the stochastic gradient  $g_t = \mathcal{O}_f(w_t, z_t)$ ; 3: 145

- 4: Update the second moment estimate:  $v_t = \beta_{2,t} v_{t-1} + (1 \beta_{2,t}) g_t^{\circ 2}$ ;
- 5: Update the first moment estimate:  $m_t = \beta_1 m_{t-1} + (1 \beta_1) g_t$ ;
- 6: Compute the adaptive learning rate:  $\eta_{v_t} = \eta_t \circ \frac{1}{\sqrt{v_t + \mu}}$ ;
- 7: Update the parameters:  $w_{t+1} = w_t \eta_{v_t} \circ m_t$ ;
- 149
   7:
   Upc

   150
   8:
   end for

151 152

146

147

148

In Adam, the random variables  $\{z_t\}_{t\geq 1}$  are mutually independent. The stochastic gradient at iteration t is denoted by  $g_t$ . The quantities  $m_t$  and  $v_t$  represent the exponential moving averages of the first and second moments of the gradients, respectively. The hyperparameters  $\beta_1$  and  $\beta_{2,t}$  control the exponential decay rates for the moment estimates. A small smoothing term  $\mu$  is introduced to prevent division by zero, and  $\eta_{v_t}$  represents the adaptive learning rate for each parameter.

In terms of notation, all vectors are column vectors unless specified otherwise, and  $\mathbf{1}^{\top}$  denotes a row vector with all elements equal to 1. For vectors  $\beta, \gamma \in \mathbb{R}^d$ , the Hadamard product (element-wise multiplication) is represented by  $\beta \circ \gamma$ , and the element-wise square root of a vector  $\gamma \in \mathbb{R}^d$  is written as  $\sqrt{\gamma}$ . Operations such as  $\beta + v_0$ ,  $\frac{1}{\beta}$ , and  $\beta^{\circ 2}$  are performed element-wise. Additionally, the *i*-th component of a vector  $\beta_t$  is denoted as  $\beta_{t,i}$ . 162 163 164 165 When analyzing Adam,  $\nabla f(w_t)$  refers to the true gradient of the loss function at iteration t. We define  $\mathscr{F}_t = \sigma(g_1, \dots, g_t)$  as the  $\sigma$ -algebra generated by the stochastic gradients up to iteration t, with  $\mathscr{F}_0 = \{\Omega, \emptyset\}$  and  $\mathscr{F}_\infty = \sigma\left(\bigcup_{t \ge 1} \mathscr{F}_t\right)$ .

167 2.2 ASSUMPTIONS

To establish our convergence results, we make the following standard assumptions, focusing on the stochastic gradient conditions. These assumptions are less restrictive than those imposed in some prior works, as highlighted in Table 1.

Assumption 2.1. (Bounded from Below Loss Function) Let  $f : \mathbb{R}^d \to \mathbb{R}$  be a loss function defined on  $\mathbb{R}^d$ . We assume that there exists a constant  $f^* \in \mathbb{R}$  such that for all  $w \in \mathbb{R}^d$ , the following inequality holds:  $f(w) \ge f^*$ .

This assumption ensures that the loss function f is bounded from below, preventing it from decreasing indefinitely during the optimization process.

Assumption 2.3. (ABC Inequality) We assume that the stochastic gradient  $g_t$  is an unbiased estimate of the true gradient, i.e.,  $\mathbb{E}[g_t | \mathscr{F}_{t-1}] = \nabla f(w_t)$ , and there exist constants  $A, B, C \ge 0$  such that for all iterations t, we have:  $\mathbb{E}[||g_t||^2 | \mathscr{F}_{t-1}] \le A(f(w_t) - f^*) + B||\nabla f(w_t)||^2 + C$ .

The ABC inequality provides a bound on the second moment of the stochastic gradients, which is crucial for analyzing the convergence of stochastic optimization algorithms.

187 188 189

193

194

195 196

166

2.3 COMPARISON WITH PRIOR WORKS ON STOCHASTIC GRADIENT ASSUMPTIONS

Our assumption on the stochastic gradient (Assumption 2.3) is relatively mild compared to those in prior works. Here, we focus on comparing with the traditional weak growth condition, coordinate weak growth assumption, and the almost surely bounded stochastic gradient assumption.

**Traditional Weak Growth Condition** The traditional weak growth condition (e.g., Bottou et al. (2018); Nguyen et al. (2018)) assumes that there exist constants  $B \ge 0$  and  $C \ge 0$  such that:

$$\mathbb{E}[\|g_t\|^2 \mid \mathscr{F}_{t-1}] \le B \|\nabla f(w_t)\|^2 + C.$$

$$\tag{1}$$

This condition bounds the expected squared norm of the stochastic gradient by a linear function of the squared norm of the true gradient plus a constant. It is stronger than our ABC inequality because it does not include the term involving the function value difference  $f(w_t) - f^*$ .

Even under this condition, current methods for analyzing Adam encounter significant difficulties.
 We will explain these challenges in the proof sketch of Lemma 4.1.

Coordinate Weak Growth Assumption Wang et al. (2024a) introduce the *coordinate weak* growth assumption, which requires that each component of the stochastic gradient satisfies a weak growth inequality. Specifically, for each coordinate i, there exist constants  $B, C \ge 0$  such that:

$$\mathbb{E}[g_{t,i}^2 \mid \mathscr{F}_{t-1}] \le B \|\nabla_i f(w_t)\|^2 + C,\tag{2}$$

where  $g_{t,i}$  and  $\nabla_i f(w_t)$  are the *i*-th components of  $g_t$  and  $\nabla f(w_t)$ , respectively.

This assumption is stronger than the traditional weak growth condition because it imposes the inequality on each coordinate individually, rather than on the overall gradient.

212

207

**Almost Surely Bounded Stochastic Gradient Assumption** Some prior works, such as He et al. (2023); Xiao et al. (2024), assume that the stochastic gradients are almost surely bounded. That is, there exists a constant  $M \ge 0$  such that for all iterations  $t: ||g_t|| \le M$  almost surely. This is a strong assumption, as it requires that the stochastic gradient norm is uniformly bounded almost surely

216 at all iterations. In practice, especially in non-convex optimization problems, this assumption is 217 often violated (see Wang et al. 2023). For instance, when optimizing deep neural networks, gradient 218 norms can become unbounded due to the complexity and non-linearity of the models. Moreover, this 219 assumption implies that the true gradient is also bounded by M, because  $\|\nabla f(w_t)\|^2 \leq \mathbb{E}[\|g_t\|^2]$ 220  $\mathscr{F}_{t-1} \leq M^2$ . Our assumption is clearly weaker than the almost surely bounded stochastic gradient assumption, as we only require a bound on the expected squared norm of the stochastic gradient, 221 which can depend on the current function value and gradient norm, rather than a uniform almost 222 sure bound. 223

Moreover, assuming almost surely bounded stochastic gradients is hard to satisfy in practice and may not reflect realistic scenarios. As discussed in Wang et al. (2023); Khaled & Richtárik (2023), such assumptions can be unrealistic and limit the applicability of theoretical results.

Next, we introduce a property. We know that when the loss function is *L*-smooth, the true gradient of the loss function can be controlled by the loss function value  $f(w_t) - f^*$  (as shown in Lemma B.1). Therefore, we can simplify the ABC inequality as follows.

**Property 1.** Under Assumptions 2.2 and 2.3, for all iterations t, we have:

$$\mathbb{E}[\|g_t\|^2 \mid \mathscr{F}_{t-1}] \le (A + 2L_f B)(f(w_t) - f^*) + C.$$

This property demonstrates that the variance of the stochastic gradients can be bounded by the function value difference, which is a key component in our convergence analysis.

#### 2.4 Hyperparameter Settings

In this paper, to avoid overly lengthy proofs, we choose a class of representative parameter settings, as follows:

$$\begin{array}{ll} \textbf{241} \\ \textbf{242} \\ \end{array} \qquad \beta_{2,t} := \left\{ \begin{array}{ll} 1 - \alpha_0, & \text{if } t = 1 \\ 1 - \frac{1}{t^{\gamma}}, & \text{if } t \geq 2 \end{array}, \ \beta_1 \in [0, 1), \ \eta_t = \frac{1}{t^{\frac{1}{2} + \delta}}, \ \left( \alpha_0 \in [0, 1), \ \gamma \in [1, 2\delta + 1], \ \delta \in \left[0, \frac{1}{2}\right] \right) \end{array} \right.$$

Imposing restrictions on Adam's parameters, particularly  $\beta_{2,t}$ , is necessary to ensure convergence. Early studies (Reddi et al., 2018a) have demonstrated that without appropriate constraints on  $\beta_{2,t}$ , counterexamples exist where the algorithm fails to converge. Moreover, for the gradient norm to converge to zero, it is essential that  $\beta_{2,t}$  approaches 1 (Zou et al., 2019; He et al., 2023), as noted in previous works.

Some studies on complexity allow  $\beta_{2,t}$  to be constant. However, these studies typically focus on the algorithm's complexity over a finite number of iterations T. In such cases, the constant value of  $1 - \beta_{2,t}$  is inversely related to T, effectively causing  $\beta_{2,t}$  to approach 1 as T increases. This is another means of ensuring that  $\beta_{2,t}$  asymptotically approaches 1, which is crucial for convergence.

The hyperparameter settings adopted in this paper are representative and have been considered in previous studies (Zou et al., 2019; He et al., 2023). Our configuration includes settings that can achieve near-optimal complexity of  $\mathcal{O}(\ln T/\sqrt{T})$ . The logarithmic factor  $\ln T$  arises because  $\beta_{2,t}$  is chosen independent of the total number of iterations T, which is an unavoidable consequence with this class of parameters.

Our choice of hyperparameters simplifies the analysis while capturing the essential behavior of the
 Adam. Although the proof techniques can be extended to a broader range of parameter settings, this
 paper focuses primarily on the assumptions related to the convergence of the algorithm rather than
 an exhaustive exploration of hyperparameter configurations.

262 263

235

236 237

238

239

240

#### **3** THEORETICAL RESULTS

264

In this section, we establish both non-asymptotic and asymptotic convergence guarantees for the Adam within our smooth non-convex framework, as defined by Assumptions 2.1–2.3. For the nonasymptotic analysis, we derive a sample complexity bound that is independent of  $O(1/\mu)$ , providing an explicit bound on the number of iterations required to achieve a specified accuracy. In the asymptotic analysis, we consider two forms of convergence: almost sure convergence and convergence in the  $L_1$  norm. The almost sure convergence result demonstrates that, the gradient norm of almost every trajectory converges to zero. Meanwhile, the  $L_1$  convergence result reveals that the convergence across different trajectories is uniform with respect to the  $L_1$  norm of the gradient, where the  $L_1$  norm is taken in the sense of the underlying random variable, meaning the expectation of the gradient norm.

274 275

276

278

279

286

287

288 289

290

291

299

300

301

## 3.1 NON-ASYMPTOTIC SAMPLE COMPLEXITY

277 We first establish a non-asymptotic bound on the sample complexity of Adam.

**Theorem 3.1 (Non-Asymptotic Sample Complexity).** Consider the Adam algorithm as specified in Algorithm 2.1, and suppose that Assumptions 2.1–2.3 hold. Then, for any initial point and for  $T \ge 1$ , the following results hold:

$$\frac{1}{T}\sum_{t=1}^{T}\mathbb{E}[\|\nabla f(w_t)\|] \leq \begin{cases} \mathcal{O}\left(\frac{1}{T^{\frac{1}{2}-\delta}}\right), & \text{if } \delta \in (0, \frac{1}{2}], \\ \mathcal{O}\left(\frac{\ln T}{\sqrt{T}}\right), & \text{if } \gamma > 1, \ \delta = 0, \\ \mathcal{O}\left(\frac{\ln^2 T}{\sqrt{T}}\right), & \text{if } \gamma = 1, \ \delta = 0. \end{cases}$$

The constant hidden in the O notation depends on the initial point, the constants in our required assumptions (excluding  $1/\mu$ ), and the parameters  $\delta$  and  $\alpha_0$ .

This theorem provides a non-asymptotic rate of convergence for the expected gradient norm, highlighting how the choice of hyperparameters affects the convergence rate.

# 292 3.2 ASYMPTOTIC CONVERGENCE

We now present our main asymptotic convergence results, demonstrating that the gradients of the Adam converge to zero both almost surely and in the  $L_1$  sense under appropriate conditions.

**Theorem 3.2** (Asymptotic Almost Sure Convergence). Under Assumptions 2.1–2.3, consider the Adam with hyperparameters specified in Subsection 2.4 with  $\gamma > 1$  and  $\delta > 0$ . Then, the gradients of the Adam converge to zero almost surely, i.e.,  $\lim_{t\to\infty} \|\nabla f(w_t)\| = 0$  a.s.

This theorem shows that the gradients evaluated at the iterates converge to zero almost surely, indicating that the algorithm approaches a critical point of the loss function along almost every trajectory.

**Remark 1.** (Almost sure vs  $L_1$  convergence) As stated in the introduction, it is important to note that the almost sure convergence does not imply  $L_1$  convergence. To illustrate this concept, let us consider a sequence of random variables  $\{\zeta_n\}_{n\geq 1}$ , where  $\mathbb{P}(\zeta_n = 0) = 1 - 1/n^2$  and  $\mathbb{P}(\zeta_n = n^2) = 1/n^2$ . According to the Borel-Cantelli lemma, it follows that  $\lim_{n\to+\infty} \zeta_n = 0$  almost surely. However, it can be shown that  $\mathbb{E}(\zeta_n) = 1$  for all n > 0 by simple calculations.

Theorem 3.3 (Asymptotic  $L_1$ -Convergence). Under Assumptions 2.1–2.3, consider the Adam with hyperparameters specified in Subsection 2.4 with  $\gamma > 1$  and  $\delta > 0$ . Then, the gradients of the Adam converge to zero in the  $L_1$  sense, i.e.,  $\lim_{t\to\infty} \mathbb{E}[||\nabla f(w_t)||] = 0$ .

This result establishes convergence in the mean sense, showing that the expected gradient norm approaches zero as the number of iterations increases. It indicates that the convergence of gradient norms across different trajectories is uniform in the  $L_1$  norm of the random variables.

In previous works (He et al. (2023); Xiao et al. (2024)), the assumption that the stochastic gradients are uniformly bounded, i.e.,  $||g_t|| \le M$  a.s.  $(\forall t \ge 1)$ , or that the gradients themselves are uniformly bounded, i.e.,  $||\nabla f(w_t)|| \le M$  ( $\forall t \ge 1$ ), allows almost sure convergence to directly imply  $L_1$ convergence via the *Lebesgue's Dominated Convergence* theorem. However, in our framework, which deals with potentially unbounded stochastic gradients or gradients, proving  $L_1$  convergence is much more challenging. We will elaborate on this in the next section.

320 321

322

## 4 FRAMEWORK FOR ANALYZING ADAM

In this section, we present the analytical framework that underpins our convergence analysis for the Adam. Our approach is built upon the insights provided by existing methods, while introducing

new techniques to address the limitations of previous analyses and provide a more comprehensive
 understanding of Adam's behavior under weaker assumptions. Our core innovations are detailed in
 Section 4.3.1, Section 4.4, and Section 4.5.

# 4.1 KEY PROPERTIES OF ADAPTIVE LEARNING RATES

We begin by characterizing the fundamental properties of the adaptive learning rate sequence  $\eta_{v_t}$ . These properties are critical as they directly influence the behavior of the algorithm and are foundational to our subsequent analysis. By understanding how these properties interact with the algorithm's dynamics, we obtain more insights on the conditions under which Adam converges.

**Property 2.** Each element  $\eta_{v_t,i}$  of the sequence  $\{\eta_{v_t}\}_{t\geq 1} = \{[\eta_{v_t,1}, \eta_{v_t,2}, \dots, \eta_{v_t,d}]^\top\}_{t\geq 1}$  is monotonically decreasing with respect to t.

This property ensures that the learning rate becomes progressively smaller as the algorithm progresses, which is a crucial factor in the stability and convergence of Adam.

**Property 3.** Each element  $\eta_{v_t,i}$  of the sequence  $\{\eta_{v_t}\}_{t\geq 1} = \{[\eta_{v_t,1}, \eta_{v_t,2}, \dots, \eta_{v_t,d}]^\top\}_{t\geq 1}$  satisfies the inequality  $t^{\gamma}v_{t,i} \geq \alpha_1 S_{t,i}$ , where we define  $\alpha_1 := \min\{1 - \alpha_0, \alpha_0\}, S_{t,i} := v + \sum_{k=1}^t g_{k,i}^2$  for all  $t \geq 1$ , and  $S_{0,i} := v$ .

This property highlights the relationship between the accumulated gradient information  $S_{t,i}$  and the adaptive learning rate, ensuring that the latter appropriately scales with the former as iterations proceed.

**Remark 4.1.** For the purpose of simplifying the proofs of subsequent theorems, we define two auxiliary parameters:  $\Sigma_{v_t} := \sum_{i=1}^d v_{t,i}$  and  $S_t := \sum_{i=1}^d S_{t,i}$ . Additionally, for convenience in the subsequent proofs, we define a new initial parameter based on  $S_{0,i}$  as  $\eta_{v_0,i} = S_{0,i}/\alpha_1 = v/\alpha_1$ .

These definitions of auxiliary parameters help streamline the analysis, making the mathematical expressions more manageable and the proofs more concise.

With the key properties of the adaptive learning rates established, we now turn our attention to analyzing the momentum term, which plays a crucial role in the Adam.

353 354 355

361 362

371 372

373

377

#### 4.2 HANDLING THE MOMENTUM TERM

To effectively analyze the momentum term in the Adam, we adopt a classical method introduced by Liu et al. (2020). The momentum term introduces additional complexity in the analysis due to its recursive nature, which can complicate the convergence proofs. To address this, we construct an auxiliary variable  $u_t$  that simplifies the analysis by decoupling the momentum term from the update process. This auxiliary variable is defined as follows:

$$u_t := \frac{w_t - \beta_1 w_{t-1}}{1 - \beta_1} = w_t + \frac{\beta_1}{1 - \beta_1} (w_t - w_{t-1}) = w_t - \frac{\beta_1}{1 - \beta_1} \eta_{v_{t-1}} \circ m_{t-1}.$$
 (3)

The introduction of  $u_t$  allows us to handle the momentum term more effectively by transforming the recursive nature of the updates into a more tractable form. Specifically, we can express the relationship between successive iterations of  $u_t$  as follows:

$$u_{t+1} - u_t = -\eta_{v_t} \circ g_t + \frac{\beta_1}{1 - \beta_1} (\underbrace{\eta_{v_{t-1}} - \eta_{v_t}}_{\Delta_t}) \circ m_{t-1}.$$
(4)

This recursive relation is instrumental in breaking down the complex dependencies introduced by the momentum term, which will facilitate the convergence analysis.

#### 4.3 ESTABLISHING THE APPROXIMATE DESCENT INEQUALITY

In the convergence analysis of stochastic gradient descent (SGD), a fundamental tool is the *approximate descent inequality*, which quantifies the expected decrease in the objective function at each iteration. Specifically, for SGD, the approximate descent inequality is given by:

 $f(w_{t+1}) \le f(w_t) - \frac{\eta_t}{2} \|\nabla f(w_t)\|^2 + \frac{\eta_t^2 L}{2} \mathbb{E}[\|g_t\|^2 \mid \mathscr{F}_t] + \eta_t \nabla f(w_t)^\top (\nabla f(w_t) - g_t), \quad (5)$ 

where  $\eta_t$  is the learning rate, L is the Lipschitz constant, and  $g_t$  is the stochastic gradient.

Motivated by the success of this approach in analyzing SGD, we aim to establish a similar approximate descent inequality for the Adam. The goal is to develop a descent inequality that captures the adaptive nature of Adam's learning rates while maintaining the essential structure seen in the analysis of SGD.

To this end, we present the following key result, which forms the cornerstone of our convergence analysis for Adam.

**Lemma 4.1 (Approximate Descent Inequality).** Consider the sequences  $\{w_t\}_{t\geq 1}$ ,  $\{v_t\}_{t\geq 1}$ , and  $\{u_t\}_{t\geq 1}$  generated by Algorithm 2.1 and Eq. 4. Under Assumptions 2.1–2.3, the following sufficient decrease inequality holds:

$$\Pi_{\Delta,t+1}\hat{f}(u_{t+1}) - \Pi_{\Delta,t}\hat{f}(u_t) \le -\frac{5}{16}\Pi_{\Delta,t+1}\sum_{i=1}^d \zeta_i(t) + (L_f+1)\Pi_{\Delta,t+1}\sum_{i=1}^d \eta_{v_t,i}^2 g_{t,i}^2 + \Pi_{\Delta,t+1}M_t.$$
(6)

Here,

394

396 397

399 400 401

$$\hat{f}(u_t) := f(u_t) + C \sum_{i=1}^d \eta_{v_{t-1},i} + \frac{C_2}{(1-\beta_1)\sqrt{v}} \|m_{t-1}\|^2, \ \zeta_i(t) := \eta_{v_t,i} (\nabla_i f(w_t))^2,$$
$$\Pi_{\Delta,t} := \prod_{k=1}^{t-1} (1+C_1 \overline{\Delta}_k)^{-1} \quad (\forall \ t \ge 2), \quad \Pi_{\Delta,1} := 1, \ \overline{\Delta}_t := \sum_{i=1}^d \mathbb{E}(\Delta_{t,i} \mid \mathscr{F}_{t-1}),$$
$$M_t := M_{t,1} + M_{t,2} + M_{t,3}. \tag{7}$$

402 Constants  $C_1$  and  $C_2$  are defined in Eq. 21;  $M_{t,1}$  is defined in Eq. 15;  $M_{t,2}$  and  $M_{t,3}$  are defined in Eq. 16. 404

This lemma introduces  $\Pi_{\Delta,t} \hat{f}(u_t)$  as a new Lyapunov function for the Adam, which plays a crucial role in our analysis. In inequality 6, the term  $-\frac{5}{16}\Pi_{\Delta,t+1}\sum_{i=1}^{d}\zeta_i(t)$  can be interpreted as the descent term, representing the expected decrease in the Lyapunov function. The second term,  $(L_f + 1)\Pi_{\Delta,t+1}\sum_{i=1}^{d}\eta_{v_t,i}^2g_{t,i}^2$ , accounts for the squared error due to the stochastic nature of the gradients. The third term on the right-hand side,  $\Pi_{\Delta,t+1}M_t$ , is a martingale difference sequence with respect to the filtration  $\{\mathscr{F}_t\}_{t\geq 1}$ , which, due to its zero expectation, can be considered to have no overall impact on the algorithm's iteration process.

This structure closely resembles the approximate descent inequality commonly used in the analysis of SGD. For comparison, the approximate descent inequality for SGD is given by Eq. 5.

We now proceed to provide the main idea of proving Lemma 4.1 and highlight the key steps and challenges involved in establishing this result for Adam.

To begin with, we calculate the difference in the loss function values between two consecutive auxiliary variables  $\{u_t\}_{t\geq 1}$  that we introduced. We obtain the following expression (informal):

$$f(u_{t+1}) - f(u_t) \leq -\sum_{i=1}^d \mathbb{E}\left[\eta_{v_t,i} \nabla f(w_t) g_{t,i} | \mathscr{F}_{t-1}\right] + \underbrace{\mathcal{O}\left(\sum_{i=1}^d \eta_{v_t,i}^2 g_{t,i}^2\right)}_{Term_{t,2}} + \underbrace{\sum_{i=1}^d \mathbb{E}\left[\eta_{v_t,i} \nabla f(w_t) g_{t,i} | \mathscr{F}_{t-1}\right] - \sum_{i=1}^d \eta_{v_t,i} \nabla f(w_t) g_{t,i} + R_t.$$
(8)

428 Term<sub>t,3</sub> 429 It can be observed that the above equation is simply a second-order Taylor expansion of  $f(u_{t+1}) - f(u_t)$  (since an L-smooth function is almost everywhere twice differentiable). Term<sub>t,1</sub> represents 431 the first-order term, which in general serves as the descent term. Term<sub>t,2</sub> is the quadratic error, and 431 Term<sub>t,3</sub> is a martingale difference sequence. The remaining term  $R_t$  is negligible and can be ignored.  $-Term_{t,1} = -\sum_{i=1}^{n} \mathbb{E}\left[\eta_{v_{t,i}} \nabla f(w_t) g_{t,i} \middle| \mathscr{F}_{t-1}\right]$ 

432 In the informal explanation provided in the sketch, these were collectively referred to as remainder 433 terms. For the exact formulation, refer to the detailed proof in Appendix D.1. 434

While handling the quadratic error term  $Term_{t,2}$  is relatively straightforward using standard scaling 435 techniques, addressing the first-order term  $Term_{t,1}$  is more challenging due to the adaptive nature of 436 Adam's learning rates. Specifically,  $\eta_{v_{t,i}}$  and  $g_{t,i}$  are both  $\mathscr{F}_t$ -measurable, which necessitates the 437 introduction of an auxiliary random variable  $\tilde{\eta}_{v_t,i} \in \mathscr{F}_{t-1}$  to facilitate the extraction of the learning 438 rate from the conditional expectation. In this paper, we choose the auxiliary random variable  $\eta_{v_{t-1},i}$ 439 to approximate  $\eta_{v_t,i}$ . There are also other forms of this approximation, as discussed by Wang et al. 440 (2023; 2024a). This allows us to rewrite the first-order term as: 441

445 446

447 448

450

451

453

The presence of  $Term_{t,4}$  introduces an additional layer of complexity in the analysis, as it reflects the 449 difference between successive adaptive learning rates. Addressing this extra error term is crucial for establishing robust convergence guarantees under the ABC inequality or weak growth conditions. Existing approaches to handling such terms, which often rely on the cancellation of errors through 452 preceding descent terms, fall short in this context. This necessitates a more innovative strategy, which we present in the following section.

 $= -\underbrace{\sum_{i=1}^{d} \mathbb{E}\left[\eta_{v_{t-1},i} \nabla f(w_t) g_{t,i} | \mathscr{F}_{t-1}\right]}_{\text{Descent-Term}} + \underbrace{\sum_{i=1}^{d} \mathbb{E}\left[(\eta_{v_{t-1},i} - \eta_{v_t,i}) \nabla f(w_t) g_{t,i} | \mathscr{F}_{t-1}\right]}_{\text{Terms}}.$ 

#### 454 455 456

#### 4.3.1 ADDRESSING THE EXTRA ERROR TERM: OUR INNOVATIVE APPROACH

The term  $Term_{t,4}$ , introduced by the difference between  $\eta_{v_{t-1},i}$  and  $\eta_{v_t,i}$ , presents a significant chal-457 lenge in the convergence analysis of Adam under the ABC inequality or weak growth conditions. In 458 existing methods, it is common to attempt to cancel out such error terms by leveraging the preceding 459 descent term Descent- $Term_t$ . However, this approach might not work within the ABC framework. 460 Recent works such as Wang et al. (2023; 2024a) have shown that, under existing techniques, the best 461 one can achieve is a weakened form of the stochastic gradient assumption, namely the coordinate 462 weak growth condition. 463

To overcome these limitations, we introduce a novel approach to handle  $Term_{t,4}$ . We scale it as 464 follows: 465

$$\operatorname{Term}_{t,4} \leq \frac{1}{2} \sum_{i=1}^{d} \mathbb{E} \left[ \eta_{v_{t-1},i} \nabla f(w_t) g_{t,i} | \mathscr{F}_{t-1} \right] + C_1 f(u_t) \cdot \sum_{i=1}^{d} \mathbb{E} [\Delta_{t,i} | \mathscr{F}_{t-1}]$$

$$+C\sum_{i=1}^{a}\Delta_{t,i}+C\sum_{i=1}^{a}\left(\mathbb{E}[\Delta_{t,i} \mid \mathscr{F}_{t-1}]-\Delta_{t,i}\right),$$

$$\underbrace{+C\sum_{i=1}^{a}\Delta_{t,i}+C\sum_{i=1}^{a}\left(\mathbb{E}[\Delta_{t,i} \mid \mathscr{F}_{t-1}]-\Delta_{t,i}\right)}_{Terms},$$

472 473

474

475

479 480

484 485

466 467 468

> where  $\Delta_{t,i} := \eta_{v_{t-1},i} - \eta_{v_t,i}$ , and  $C_1 := \frac{A+2L_fB}{2}(L_f+1)$ . The key term in the inequality is  $C_1 f(u_t) \sum_{i=1}^d \mathbb{E}[\Delta_{t,i} \mid \mathscr{F}_{t-1}]$ , which cannot be easily canceled out by existing methods.

To handle this issue, we move the term  $C_1 f(u_t) \sum_{i=1}^d \mathbb{E}[\Delta_{t,i} | \mathscr{F}_{t-1}]$  to the left-hand side of inequality 8 and combine it with the existing  $f(u_t)$  term. This leads to a new iteration inequality of 476 477 the form: 478

$$f(u_{t+1}) - (1 + C_1\overline{\Delta}_t)f(u_t) \le -\frac{1}{2}Descent \cdot Term_t + M \cdot Term_t + Term_{t,2} + R \cdot Term_t.$$
(9)

481 In the inequality M-Term<sub>t</sub> = Term<sub>t,3</sub> + Term<sub>t,5</sub> is a martingale difference sequence and R-Term<sub>t</sub> is 482 the (neglectable) remainder term by combining all other terms from the inequalities. To express this inequality in a form resembling a Lyapunov function, we introduce an auxiliary product variable: 483

$$\Pi_{\Delta,t} := \prod_{k=1}^{t-1} (1 + C_1 \overline{\Delta}_k)^{-1} \quad (\forall t \ge 2), \quad \Pi_{\Delta,1} := 1.$$

Multiplying both sides of the inequality by  $\Pi_{\Delta,t+1}$ , we obtain the following reformulated inequality:

489 490

494

495

496

498

486

487

 $\Pi_{\Delta,t+1}f(u_{t+1}) - \Pi_{\Delta,t}f(u_t) \le -\frac{1}{2}\Pi_{\Delta,t+1} \cdot \textit{Descent-Term}_t + \Pi_{\Delta,t+1} \cdot \textit{M-Term}_t$  $+ \Pi_{\Delta,t+1} \cdot \textit{Term}_{t,2} + \Pi_{\Delta,t+1} \cdot \textit{R-Term}_t.$ (10)

This reformulation introduces  $\Pi_{\Delta,t}$  as a scaling factor, which, along with the original Lyapunov function, captures the impact of  $Term_{t,4}$ . The resulting inequality closely parallels the approximate descent inequality for SGD, with additional terms accounting for Adam's adaptive nature.

The handling of  $Term_{t,4}$  in our analysis framework is a significant advancement over existing methods. It allows us to establish stronger convergence guarantees under more general conditions.

# 497 4.4 DERIVING SAMPLE COMPLEXITY AND ALMOST SURE CONVERGENCE

After establishing the *Approximate Descent Inequality*, the next step is to derive the sample complexity and almost sure convergence results for Adam. The methodology for obtaining these results largely mirrors the approaches traditionally used in the analysis of SGD. Specifically, the inequality provides a foundation for bounding the expected decrease in the loss function, which can then be used to establish both sample complexity and almost sure Convergence.

However, a key difference in our analysis lies in the introduction of the term  $\Pi_{\Delta,t+1}$  within the *Approximate Descent Inequality*. This term introduces a new layer of complexity not present in the standard SGD analysis. In particular, we are required to bound the *p*-th moment of the reciprocal of this term, i.e.,  $\mathbb{E}[\Pi_{\Delta,t+1}^{-p}]$ ,  $(p \ge 1)$ . Due to the unique structure of  $\Pi_{\Delta,t+1}$ , determining a bound for this *p*-th moment is a non-trivial task.

To address this challenge, we leverage tools from discrete martingale theory, particularly the *Burkholder's* inequality. It allows us to establish a recursive relationship between the *p*-th moment  $\mathbb{E}[\Pi_{\Delta,t+1}^{-p}]$  and the *p*/2-th moment  $\mathbb{E}[\Pi_{\Delta,t+1}^{-p/2}]$ . This recursive structure is crucial as it enables us to iteratively bound the higher moments of  $\Pi_{\Delta,t+1}^{-1}$ .

514 Once the recursive relationship is established, we apply fundamental theorems from measure theory, 515 such as the *Lebesgue's Monotone Convergence* theorem or the *Lebesgue's Dominated Convergence* 516 theorem, to obtain the final bound on the *p*-th moment.

- The detailed process for bounding  $\mathbb{E}[\Pi_{\Delta,t+1}^{-p}]$  can be found in Lemma B.2 and Lemma B.3.
- 519 4.5 ESTABLISHING ASYMPTOTIC  $L_1$  CONVERGENCE

521 Since we have already proved almost sure convergence in Theorem 3.2, it is natural to attempt to 522 prove  $L_1$  convergence via the *Lebesgue's Dominated Convergence* theorem. To achieve this, we 523 need to find a function h that is  $\mathscr{F}_{\infty}$ -measurable and satisfies  $\mathbb{E}|h| < +\infty$ , and such that for all 524  $t \ge 1$ , we have  $\|\nabla f(w_t)\| \le |h|$ . Since for all t we naturally have  $\|\nabla f(w_t)\| \le \sup_{k\ge 1} \|\nabla f(w_k)\|$ , 525 we only need to prove that  $\mathbb{E}[\sup_{k>1} \|\nabla f(w_k)\|] < +\infty$ .

This task presents a significant challenge because, within our analytical framework, we cannot assume that the gradients are uniformly bounded, which means we cannot directly apply the *Lebesgue's Dominated Convergence* theorem. Instead, we need to utilize advanced techniques from discrete martingale theory, specifically the first hitting time decomposition method, to obtain a bound on this maximal expectation. The detailed process can be found in Appendix E.4.

531 532

533

5 CONCLUSION

<sup>534</sup> We have introduced a novel and comprehensive framework for analyzing the convergence properties <sup>535</sup> of Adam. Our frame starts with weak assumptions such as the ABC inequality. By identifying the <sup>536</sup> key properties of the learning rate, handling the momentum term, and establishing the approximate <sup>537</sup> descent inequality, the frame concludes the sample complexity, almost surely convergence, and <sup>538</sup> asymptotic  $L_1$  convergence results of Adam. Our techniques overcome the limitations of existing <sup>539</sup> analyses, and show that Adam's convergence guarantees can be aligned with those of SGD, which <sup>539</sup> justifies the applicability of Adam across a wide range of machine learning problems.

# 540 REFERENCES

542 543 544	Léon Bottou. Large-scale machine learning with stochastic gradient descent. In Proceedings of COMPSTAT'2010: 19th International Conference on Computational StatisticsParis France, August 22-27, 2010 Keynote, Invited and Contributed Papers, pp. 177–186. Springer, 2010.
545 546	Léon Bottou, Frank E Curtis, and Jorge Nocedal. Optimization methods for large-scale machine learning. <i>SIAM review</i> , 60(2):223–311, 2018.
548 549	Saeed Ghadimi and Guanghui Lan. Stochastic first-and zeroth-order methods for nonconvex stochas- tic programming. <i>SIAM journal on optimization</i> , 23(4):2341–2368, 2013.
550 551 552	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.
553 554 555	Ahmed Khaled and Peter Richtárik. Better theory for SGD in the nonconvex world. <i>Trans. Mach. Learn. Res.</i> , 2023, 2023.
556 557	Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. <i>arXiv preprint arXiv:1412.6980</i> , 2014.
558 559 560	Haochuan Li, Alexander Rakhlin, and Ali Jadbabaie. Convergence of adam under relaxed assumptions. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.
561 562	Yanli Liu, Yuan Gao, and Wotao Yin. An improved analysis of stochastic gradient descent with momentum. <i>Advances in Neural Information Processing Systems</i> , 33:18261–18271, 2020.
563 564 565	Lam Nguyen, Phuong Ha Nguyen, Marten Dijk, Peter Richtárik, Katya Scheinberg, and Martin Takác. Sgd and hogwild! convergence without the bounded gradients assumption. In <i>International Conference on Machine Learning</i> , pp. 3750–3758. PMLR, 2018.
566 567 568	Sashank J. Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. In 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings. OpenReview.net, 2018a.
569 570 571	Sashank J Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. In <i>International Conference on Learning Representations (ICLR)</i> , 2018b.
572 573 574	Bohan Wang, Huishuai Zhang, Zhiming Ma, and Wei Chen. Convergence of adagrad for non-convex objectives: Simple proofs and relaxed assumptions. In <i>The Thirty Sixth Annual Conference on Learning Theory</i> , pp. 161–190. PMLR, 2023.
575 576 577	Bohan Wang, Jingwen Fu, Huishuai Zhang, Nanning Zheng, and Wei Chen. Closing the gap between the upper bound and lower bound of adam's iteration complexity. <i>Advances in Neural Information Processing Systems</i> , 36, 2024a.
578 579 580 581	Bohan Wang, Yushun Zhang, Huishuai Zhang, Qi Meng, Ruoyu Sun, Zhi-Ming Ma, Tie-Yan Liu, Zhi-Quan Luo, and Wei Chen. Provable adaptivity of adam under non-uniform smoothness. In <i>Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining</i> , pp. 2960–2969, 2024b.
582 583 584 585	Nachuan Xiao, Xiaoyin Hu, Xin Liu, and Kim-Chuan Toh. Adam-family methods for nonsmooth optimization with convergence guarantees. <i>Journal of Machine Learning Research</i> , 25(48):1–53, 2024.
586 587 588	Yushun Zhang, Congliang Chen, Naichen Shi, Ruoyu Sun, and Zhi-Quan Luo. Adam can converge without any modification on update rules. <i>Advances in neural information processing systems</i> , 35:28386–28399, 2022.
589 590 591	Difan Zou and Li Shen. Improved convergence analysis of stochastic optimization algorithms for nonconvex optimization. In Advances in Neural Information Processing Systems (NeurIPS), 2019.
592 593	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.

1	Intr	oduction		
1	1 1	Palatad Works		
	1.1			
2	Prel	iminaries		
	2.1	Adam		
	2.2	Assumptions		
	2.3	Comparison with Prior Works on Stochastic Gradient Assumptions		
	2.4	Hyperparameter Settings		
3	The	oretical Results		
	3.1	Non-Asymptotic Sample Complexity		
	3.2	Asymptotic Convergence		
4	Fra	mework for Analyzing Adam		
	4.1	Key Properties of Adaptive Learning Rates		
	4.2	Handling the Momentum Term		
	4.3	Establishing the Approximate Descent Inequality		
		4.3.1 Addressing the Extra Error Term: Our Innovative Approach		
	4.4	Deriving Sample Complexity and Almost Sure Convergence		
	4.5	Establishing Asymptotic $L_1$ Convergence		
5	Con	clusion		
A	Esta	blishing Key Properties between $w_t$ and $u_t$		
B	Sup	porting Lemmas		
	B.1	Theorem Dependency Graph		
	B.2	The Basic Form of Additional Important Lemmas		
С	Pro	ofs of Vital Properties		
	C.1	The Proof of Property 2		
	C.2	The Proof of Property 3		
	C.3	The Proof of Property 5		
	C.4	The Proof of Property 5		
D	Pro	ofs of Theorems and Lemmas		
	D.1	Proofs of Lemma 4.1		
	D.2	The Proof of Lemma B.1		
	D.3	The Proof of Lemma B.2		
	<b>D</b> 4	The Dreef of Lemma D 2		

648	D.5	The Proof of Lemma B.4	26
649 650	D.6	The Proof of Lemma B.5	27
651	D.7	The Proof of Lemma B.6	27
652	D.8	The Proof of Lemma B.7	29
653 654	D.9	The Proof of Lemma B.9	30
655	D 10	The Proof of Lemma B 10 $\sim$	31
656	D.11	The Proof of Lemma D 11 $($	27
657 658	D.11		32 22
659	D.12	. The Proof of Lemma B.12	33
660 E	Pro	ofs of Theorems	34
662	E.1	The Proof of Theorem 3.1	34
663	E.2	The Proof of Lemma B.8	35
664 665	E.3	The Proof of Theorem 3.2	35
666	E 4	The Proof of Theorem 3.3	38
667	1.1		20
668			
669			

#### 702 A ESTABLISHING KEY PROPERTIES BETWEEN $w_t$ and $u_t$

Here, we establish two key properties that connect the original variable  $w_t$  and the auxiliary variable  $u_t$  which defined in Section 4.2. These properties are crucial for bounding the changes in the momentum term and linking the function values at different points in the iteration process.

**Property 4.** For any iteration step t, the following inequality holds:

 $||m_t||^2 - ||m_{t-1}||^2 \le -(1-\beta_1)||m_{t-1}||^2 + (1-\beta_1)||g_t||^2.$ 

This property establishes a bound on the change in the momentum term, which is critical for ensuring
that the momentum does not increase indefinitely during the optimization process. Controlling the
momentum in this manner is a key aspect of proving convergence.

**Property 5.** For any iteration step t, the following inequality holds:

$$f(w_t) \le (L_f + 1)f(u_t) + \frac{(L_f + 1)\beta_1^2}{2(1 - \beta_1)^2} \left\| \eta_{v_{t-1}} \circ m_{t-1} \right\|^2$$

This property links the function values at  $w_t$  and  $u_t$ , providing a foundation for analyzing the convergence of  $f(w_t)$ . By establishing this relationship, we can relate the behavior of the original variable  $w_t$  to the more manageable auxiliary variable  $u_t$ , thereby simplifying the overall convergence analysis.

# 756 B SUPPORTING LEMMAS

# B.1 THEOREM DEPENDENCY GRAPH

In this section, we will supplement several additional supporting lemmas that are crucial to the overall proof. Due to the large number of lemmas, we have combined these lemmas with those in the main text and theorems to create a lemma-theorem dependency graph. Readers can refer to this graph while following the proofs.



# 810 B.2 THE BASIC FORM OF ADDITIONAL IMPORTANT LEMMAS

**Lemma B.1.** Suppose that f(x) is differentiable and lower bounded  $f^* = \inf_{x \in \mathbb{R}^d} f(x) > -\infty$ and  $\nabla f(x)$  is Lipschitz continuous with parameter  $\mathcal{L} > 0$ , then  $\forall x \in \mathbb{R}^d$ , we have

$$\left\|\nabla f(x)\right\|^2 \le 2\mathcal{L}\big(f(x) - f^*\big)$$

**Lemma B.2.** Let  $\{(X_n, \mathscr{F}_n)\}_{n \ge 1}$  be a non-negative adapted process such that  $\sum_{n=1}^{+\infty} X_n = M < +\infty$  almost surely, where M is a finite constant. Define the partial sum of conditional expectations as  $\Lambda_T := \sum_{n=1}^T \mathbb{E}[X_n | \mathscr{F}_{n-1}]$ . Then:

- (i) The sequence  $\{\Lambda_T\}_{T\geq 1}$  converges almost surely, i.e.,  $\Lambda_T \xrightarrow{a.s.} \Lambda$ , where  $\Lambda := \sum_{n=1}^{+\infty} \mathbb{E}[X_n \mid \mathscr{F}_{n-1}].$
- (ii) For any  $p \ge 1$ , the sequence  $\{\Lambda_T\}_{T\ge 1}$  converges in  $L_p$ , i.e.,  $\lim_{T\to\infty} \mathbb{E}\left[|\Lambda_T \Lambda|^p\right] = 0$ . Meanwhile, the p-th moment of the limit  $\Lambda$  is bounded by a constant  $C_{\Lambda}(p) > 0$ , where  $C_{\Lambda}(p) = o(p^{\sqrt{p}})$ .
- (iii) The arbitrary p-th moment of the random variable  $e^{\Lambda}$  also exists, and its upper bound depends only on p and M. We denote this upper bound by  $C_{e^{\Lambda}}(p, M)$ .

**Lemma B.3.** For  $\Pi_{\Delta,T+1}$  as defined in Eq. 7, for any  $T \ge 0$  and any  $p \ge 1$ , the p-th moment of its reciprocal is bounded, i.e.,

$$\mathbb{E}\left[\Pi_{\Delta,T+1}^{-p}\right] < C_{v,d,p} < +\infty,$$

where  $C_{v,d,p}$  is a constant that depends only on v, d, and p.

*Moreover, we have that*  $\Pi_{\Delta,\infty}^{-1} := \lim_{t \to +\infty} \Pi_{\Delta,t}^{-1} < +\infty$  *a.s., and for any*  $p \ge 1$ *, the p-th moment of*  $\Pi_{\Delta,\infty}^{-1}$  *exists, with* 

$$\mathbb{E}\left[\Pi_{\Delta,\infty}^{-p}\right] \le C_{v,d,p} < +\infty$$

**Lemma B.4.** Consider the Adam in Algorithm 2.1 and suppose that Assumption 2.1~2.3 hold, then for any initial point, and  $T \ge 1$ , we have:

$$\frac{\Pi_{\Delta,t+1}\hat{f}(u_{t+1})}{\sqrt{S_t}} \le \frac{\hat{f}(u_1)}{\sqrt{dv}} + \frac{3(L_f+1)d}{\alpha_1\sqrt{v}} + \sum_{t=1}^T \frac{\Pi_{\Delta,t+1}}{\sqrt{S_{t-1}}} M_t.$$
 (11)

**Lemma B.5.** Consider the Adam in Algorithm 2.1 and suppose that Assumption 2.1~2.3 hold, then for any initial point and  $\forall \phi > 0$ , we have for any  $T \ge 1$ , the following inequality:

$$\frac{\Pi_{\Delta,T}\sqrt{S_T}}{(T+1)^{\phi}} \le \sqrt{dv} + \sum_{t=1}^T \Pi_{\Delta,t}\Lambda_{\phi,t}.$$
(12)

where

$$\Lambda_{\phi,t} := \frac{\|g_t\|^2}{(t+1)^{\phi}\sqrt{S_{t-1}}}$$

and  $S_T$  is defined in Eq. 4.1.

**Lemma B.6.** Consider the Adam as defined in Algorithm 2.1, and suppose that Assumptions 2.1 through 2.3 hold. Then, for any initial point and for all  $T \ge 1$ , there exists a random variable  $\zeta$  such that the following results hold:

(a)  $0 \le \zeta < +\infty$  almost surely, and  $\mathbb{E}(\zeta)$  is uniformly bounded above by a constant  $C_{\zeta}$ , which depends on the initial point and the constants in the required assumptions (excluding  $1/\mu$ ). The explicit form of this upper bound is provided in Eq. 30.

(b) 
$$\sqrt{S_T} \leq (T+1)^2 \zeta$$
, and  $\ln(S_T) \leq \zeta' \ln(T+1)$ , where  $\zeta' := 8 \ln\left(\max\left\{e, \frac{\sqrt{2}\Pi_{\Delta,\infty}^{-1} \zeta}{\sqrt{v}}\right\}\right)$ 

Lemma B.7. Consider the Adam in Algorithm 2.1 and suppose that Assumption 2.1 2.3 hold, then for any initial point, and  $T \ge 1$ , the following results hold: 

where  $C_5$  and  $C_6$  are two constants that depend on the initial point and the constants in our required assumptions (excluding  $1/\mu$ ), and  $C_{4,\delta}$  is a constant that depends on the initial point,  $\delta$ , and the constants in our required assumptions (excluding  $1/\mu$ ).

 $\sum_{t=1}^{T} \mathbb{E}\left[\Pi_{\Delta,t+1} \sum_{i=1}^{d} \zeta_i(t)\right] \leq \begin{cases} C_{4,\delta}, & \text{if } \delta \in (0,1]\\ C_5 + C_6 \mathbb{E}\left[\ln(S_T)\right], & \text{if } \delta = 0 \end{cases},$ 

Lemma B.8 (Subsequence Convergence). Under Assumptions 2.1–2.3, consider the Adam (Al-gorithm 2.1) with hyperparameters as specified in Subsection 2.4, where  $\delta > 0$ . Then, there exists a subsequence  $\{w_{c_t}\}_{t\geq 1}$  such that its gradients converge to zero almost surely, i.e.,  $\lim_{t \to \infty} \|\nabla f(w_{c_t})\| = 0 \quad a.s.$ 

Lemma B.9. Consider the Adam as defined in Algorithm 2.1, and assume that Assumptions 2.1 through 2.3 hold. Then, for any initial point and for all  $T \ge 1$ , the following results hold: 

- When 
$$\delta = 0$$
, we have

$$\sup_{t \ge 1} \frac{\Pi_{\Delta,t+1}(f(w_t) - f^*)}{\ln^2(t+1)} < +\infty \text{ a.s., } \sup_{T \ge 1} \mathbb{E}\left[\frac{\Pi_{\Delta,t+1}(f(w_t) - f^*)}{\ln^2(t+1)}\right] < M_0 < +\infty,$$

- When  $\delta > 0$ , we have

$$\sup_{t \ge 1} \Pi_{\Delta, t+1}(f(w_t) - f^*) < +\infty \text{ a.s., } \sup_{T \ge 1} \mathbb{E}\left[\Pi_{\Delta, t+1}(f(w_t) - f^*)\right] < M_{\delta} < +\infty$$

and

$$\sup_{t \ge 1} \prod_{\Delta, t} ||m_{t-1}|| < +\infty \ a.s.,,$$

where  $M_0$  and  $M_{\delta}$  are two constants that depend on the initial point and the constants in our assumptions (excluding  $1/\mu$ ). 

Lemma B.10. Consider the Adam in Algorithm 2.1 and suppose that Assumption 2.1 2.3 hold, then for any initial point,  $T \ge 1$ ,  $i \in [1, d]$ , there is

$$\mathbb{E}(S_T^{3/4}) = \begin{cases} \mathcal{O}(T^{3/4}), & \text{if } \delta \in (0,1] \\ \mathcal{O}(T^{3/4} \ln^{3/2} T), & \text{if } \delta = 0 \end{cases}.$$

where constant hidden in  $\mathcal O$  depends only on initial point, and the constants in our required assump-tions (not includes  $1/\mu$ ).

**Lemma B.11.** Under Assumptions 2.1–2.3, consider the Adam (Algorithm 2.1) with the hyperparameters specified in Subsection 2.4. Then, for any  $t \ge 1$ , the following inequality holds:

$$\sup_{t\geq 1} \mathbb{E}[\Pi_{\Delta,t+1}\Sigma_{v_t}] < \begin{cases} \mathbb{E}[\Pi_{\Delta,2}\Sigma_{v_1}] + (A+2L_fB)M_\delta + C, & \text{if } \delta \in (0,1]\\ \mathbb{E}[\Pi_{\Delta,2}\Sigma_{v_1}] + (A+2L_fB)M_0 \ln^2 t + C, & \text{if } \delta = 0 \end{cases}$$

Furthermore, if  $\lambda > 1$ , then we have

$$\sup_{t\geq 1} \mathbb{E}[\Pi_{\Delta,t+1}\Sigma_{v_t}] < \begin{cases} \left( (A+2L_f B)M_{\delta} + C \right) \sum_{t=1}^{+\infty} \frac{1}{(t+1)^{\lambda}}, & \text{if } \delta \in (0,1] \\ \left( (A+2L_f B)M_0 + C \right) \sum_{t=1}^{+\infty} \frac{\ln^2 t}{(t+1)^{\lambda}}, & \text{if } \delta = 0 \end{cases} < +\infty,$$

and the following almost sure bound:

$$\sup_{t\geq 1}\Sigma_{v_t}<+\infty\quad a.s.$$

**Lemma B.12.** Under Assumption 2.1-2.3, consider the Adam with hyperparameters in Subsection 2.4 with  $\gamma > 1$ ,  $\delta > 0$ . Then for any initial point, the following results hold: 

916  
917 
$$\sum_{t=1}^{+\infty} \eta_t \|\nabla f(w_t)\|^2 < +\infty \text{ a.s. and } \sum_{t=1}^{+\infty} \eta_t \|\nabla f(u_t)\|^2 < +\infty \text{ a.s.}$$

#### С **PROOFS OF VITAL PROPERTIES**

#### C.1 THE PROOF OF PROPERTY 2

*Proof.* Due to Algorithm 2.1, we know that

$$v_{t+1} = \beta_{2,t+1}v_t + (1 - \beta_{2,t+1})g_{t+1}^{\circ 2} = \left(1 - \frac{1}{(t+1)^{\gamma}}\right)v_t + \frac{1}{(t+1)^{\gamma}}g_{t+1}^{\circ 2}, \ (\forall t \ge 1).$$

which means

$$(13)^{\gamma} v_{t+1,i} = \left( (t+1)^{\gamma} - 1 \right) v_{t,i} + g_{t+1,i}^2 \ge t^{\gamma} v_{t,i}.$$

This implies that  $tv_{t,i}$  is monotonically non-decreasing. Subsequently, we can obtain:

$$\eta_{v_t,i} = \frac{\eta_t}{\sqrt{v_{t,i}} + \mu} = \frac{\sqrt{t^\gamma}\eta_t}{\sqrt{t^\gamma v_{t,i}} + \sqrt{t^\gamma}\mu} = \frac{\frac{1}{t^{\delta - \frac{\gamma - 1}{2}}}}{\sqrt{tv_{t,i}} + \sqrt{t^\gamma}\mu}$$

It can be seen that the numerator is monotonically decreasing and greater than 0, while the denominator is monotonically non-increasing and greater than 0. Therefore, overall, we can deduce the monotonic non-increasing property of  $\eta_{v_t}$ . 

#### C.2 THE PROOF OF PROPERTY 3

*Proof.* For  $v_{1,i}$ , we can derive the following estimate:

$$v_{1,i} = \beta_{2,1}v_{0,i} + (1 - \beta_{2,1})g_{1,i}^2 = (1 - \alpha_0)v + \alpha_0 g_{1,i}^2 = (1 - \alpha_0)v + g_{1,i}^2 - (1 - \alpha_0)g_{1,i}^2.$$

It is easy to find that  $\alpha_1 S_{1,i} \leq v_{1,i} \leq S_{1,i}$ . For  $\forall k \geq 2$ , we back to Eq. 13, acquiring  $k^{\gamma} v_{k,i} \geq 1$  $(k-1)^{\gamma}v_{k-1,i} + g_{k,i}^2$ . Next, by summing up the above iterative equations, we obtain  $\forall t \geq 2$ ,

$$t^{\gamma} v_{t,i} \ge v_{1,i} + \sum_{k=2}^{t} g_{k,i}^2$$

Next, combining the estimate for  $v_{1,i}$ , we obtain  $\forall t \ge 2$ :

$$t^{\gamma} v_{t,i} \ge (1 - \alpha_0) v + \alpha_0 g_{1,i}^2 + \sum_{k=2}^t g_{k,i}^2.$$

It is easy to find that  $t^{\gamma}v_{t,i} \geq \alpha_1 S_{t,i}$ . With this, we complete the proof.

#### C.3 The Proof of Property 5

*Proof.* According to Algorithm 2.1, we have the following iterative equations:

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$$

We take the square of the 2-norm on both sides, yielding

$$\begin{split} \|m_t\|^2 &= \|\beta_1 m_{t-1} + (1-\beta_1)g_t\|^2 \\ &= \beta_1^2 \|m_{t-1}\|^2 + 2\beta_1 (1-\beta_1) m_{t-1}^\top g_t + (1-\beta_1)^2 \|g_t\|^2 \\ &\stackrel{(a)}{\leq} \beta_1 \|m_{t-1}\|^2 + (1-\beta_1) \|g_t\|^2 \\ &\stackrel{\text{Eq. 13}}{\leq} \beta_1 \|m_{t-1}\|^2 + (1-\beta_1) \|g_t\|^2. \end{split}$$

In step (a), we used the AM-GM inequality, i.e.,

$$2\beta_1(1-\beta_1)m_{t-1}^{\top}g_t \leq \beta_1(1-\beta_1)||m_{t-1}||^2 + \beta_1(1-\beta_1)||g_t||^2,$$

that is,

$$||m_t||^2 - ||m_{t-1}||^2 \le -(1-\beta_1)||m_{t-1}||^2 + (1-\beta_1)||g_t||^2.$$

With this, we complete the proof.

# 972 C.4 THE PROOF OF PROPERTY 5

974 Proof. Due to

we have

$$f(w_t) \le f(u_t) + |f(w_t) - f(u_t)| \le (L_f + 1)f(u_t) + \frac{(L_f + 1)\beta_1^2}{2(1 - \beta_1)^2} \|\eta_{v_{t-1}} \circ m_{t-1}\|^2.$$

 $|f(w_t) - f(u_t)| = \left| \nabla f(u_t)^\top (w_t - u_t) + \frac{L_f}{2} \|w_t - u_t\|^2 \right| \le \|\nabla f(u_t)\| \|w_t - u_t\| + \frac{L_f}{2} \|w_t - u_t\|^2$ 

 $\leq \frac{1}{2} \|\nabla f(u_t)\|^2 + \frac{L_f + 1}{2} \|w_t - u_t\|^2$ 

 $= Lf(u_t) + \frac{(L_f + 1)\beta_1^2}{2(1 - \beta_1)^2} \|\eta_{v_{t-1}} \circ m_{t-1}\|^2,$ 

#### D PROOFS OF THEOREMS AND LEMMAS

#### D.1 PROOFS OF LEMMA 4.1

*Proof.* By L-smooth in Assumption 2.2, we have:

$$f(u_{t+1}) - f(u_t) \le \nabla f(u_t)^\top (u_{t+1} - u_t) + \frac{L_f}{2} \|u_{t+1} - u_t\|^2$$

Then, by substituting the iterative formula for  $u_t$  from Eq. 4 into the above inequality, we obtain:

$$f(u_{t+1}) - f(u_t) \leq -\sum_{i=1}^d \eta_{v_t,i} \nabla_i f(u_t) g_{t,i} + \frac{\beta_1}{1 - \beta_1} \sum_{i=1}^d \Delta_{t,i} \nabla_i f(u_t) m_{t-1,i} + L_f \sum_{i=1}^d \eta_{v_t,i}^2 g_{t,i}^2 + L_f \left(\frac{\beta_1}{1 - \beta_1}\right)^2 \sum_{i=1}^d \Delta_{t,i}^2 m_{t-1,i}^2$$

$$\stackrel{(a)}{=} \underbrace{-\sum_{i=1}^d \eta_{v_t,i} \nabla_i f(w_t) g_{t,i}}_{\Theta_{t,1}} + \underbrace{\sum_{i=1}^d (\eta_{v_t,i} (\nabla_i f(w_t) - \nabla_i f(u_t)) g_{t,i})}_{\Theta_{t,2}}}_{\Theta_{t,2}} + \frac{\beta_1}{1 - \beta_1} \underbrace{\sum_{i=1}^d \Delta_{t,i} \nabla_i f(u_t) m_{t-1,i}}_{\Theta_{t,3}} + L_f \underbrace{\sum_{i=1}^d \eta_{v_t,i}^2 g_{t,i}^2}_{\Theta_{t,3}} + L_f \left(\frac{\beta_1}{1 - \beta_1}\right)^2 \underbrace{\sum_{i=1}^d \Delta_{t,i}^2 m_{t-1,i}^2}_{\Theta_{t,4}}.$$
(14)

1018 Step (a) employs the identity  $\nabla_i f(u_t) = \nabla_i f(w_t) + \nabla_i f(u_t) - \nabla_i f(w_t)$ . Next, we handle  $\Theta_{t,1}$ ,  $\Theta_{t,2}, \Theta_{t,3}$  and  $\Theta_{t,4}$  separately. First, for  $\Theta_{t,1}$ , we can perform the following identity transformation:

$$\begin{array}{l} 1019\\ 1020\\ 1021\\ 1022\\ 1023\\ 1024\\ 1025 \end{array} \qquad \Theta_{t,1} = -\sum_{i=1}^{d} \eta_{v_{t,i}} \nabla_i f(w_t) g_{t,i} = -\sum_{i=1}^{d} \eta_{v_{t-1,i}} \nabla_i f(w_t) g_{t,i} + \sum_{i=1}^{d} \Delta_{t,i} \nabla_i f(w_t) g_{t,i} + \sum_{i=1}^{d} \Delta_{t,i} \nabla_i f(w_t) (\nabla_i f(w_t) - g_{t,i}), \\ \Theta_{t,1,1} \qquad \Theta_{t,1,1} \qquad \Theta_{t,1,1} \end{array}$$

$$\begin{array}{l} \Theta_{t,1} = -\sum_{i=1}^{d} \eta_{v_{t-1,i}} (\nabla_i f(w_t))^2 + \sum_{i=1}^{d} \Delta_{t,i} \nabla_i f(w_t) g_{t,i} + \sum_{i=1}^{d} \eta_{v_{t-1,i}} \nabla_i f(w_t) (\nabla_i f(w_t) - g_{t,i}), \\ \Theta_{t,1,1} \qquad \Theta_{t,1,1} \qquad \Theta_{t,1,1} \end{array}$$

$$\begin{array}{l} \Theta_{t,1} = -\sum_{i=1}^{d} \eta_{v_{t-1,i}} (\nabla_i f(w_t))^2 + \sum_{i=1}^{d} \Delta_{t,i} \nabla_i f(w_t) g_{t,i} + \sum_{i=1}^{d} \eta_{v_{t-1,i}} (\nabla_i f(w_t) - g_{t,i}), \\ \Theta_{t,1,1} \qquad \Theta_{t,1,1} \qquad \Theta_{t,1,1} \end{array}$$

$$\begin{array}{l} \Theta_{t,1} = -\sum_{i=1}^{d} \eta_{v_{t-1,i}} (\nabla_i f(w_t))^2 + \sum_{i=1}^{d} \Delta_{t,i} \nabla_i f(w_t) g_{t,i} + \sum_{i=1}^{d} \eta_{v_{t-1,i}} (\nabla_i f(w_t) - g_{t,i}), \\ \Theta_{t,1,1} \qquad \Theta_{t$$

 $\Theta_{t,1,1} = \sum_{i=1}^{d} \mathbb{E} \left( \Delta_{t,i} \nabla_i f(w_t) g_{t,i} \mid \mathscr{F}_{t-1} \right)$ 

where  $\Delta_{t,i}$  in the above inequality represents the *i*-th component of the vector  $\Delta_t$ , which is defined in Eq. 4. It can be observed that we decompose  $\Theta_1$  into a descent term  $-\sum_{i=1}^d \zeta_i(t)$ , an error term  $\Theta_{t,1,1}$ , and a martingale difference term  $M_{t,1}$ . Next, we will further scale and control the error term  $\Theta_{t,1,1}$ . Specifically, we have:

$$+\underbrace{\sum_{i=1}^{d} \left(\Delta_{t,i} \nabla_{i} f(w_{t}) g_{t,i} - \mathbb{E} \left[\Delta_{t,i} \nabla_{i} f(w_{t}) g_{t,i} \mid \mathscr{F}_{t-1}\right]\right)}_{M_{t,2}}$$

$$\stackrel{(a)}{<} \underbrace{\sum_{i=1}^{d} \sqrt{\eta_{v_{t-1},i}} \nabla_{i} f(w_{t}) \mathbb{E} \left[\sqrt{\Delta_{t,i}} g_{t,i} \mid \mathscr{F}_{t-1}\right] + M_{t,2}}_{(b)}$$

$$\stackrel{(b)}{=} 1 \underbrace{\sum_{i=1}^{d} p_{i,i} (\nabla_{i} f(w_{t}))^{2}}_{i} + \frac{1}{2} \underbrace{\sum_{i=1}^{d} \mathbb{E}^{2} \left[\sqrt{\Delta_{i,i}} g_{i,i} \mid \mathscr{F}_{t-1}\right]}_{i} + M_{t,2}$$

$$\leq \frac{1}{2} \sum_{i=1}^{d} \eta_{v_{t-1},i} (\nabla_i f(w_t))^2 + \frac{1}{2} \sum_{i=1}^{d} \mathbb{E}^2 \left[ \sqrt{\Delta_{t,i}} g_{t,i} \mid \mathscr{F}_{t-1} \right] + M_{t,2}$$

$$(c) \ 1 \ \frac{d}{2} \qquad 1 \ \frac{d}{2} \qquad c$$

1044  
1045  
1046  
1046  
1047  
(c) 
$$\frac{1}{2} \sum_{i=1}^{d} \zeta_i(t) + \frac{1}{2} \sum_{i=1}^{d} \mathbb{E}[g_{t,i}^2 \mid \mathscr{F}_{t-1}] \cdot \mathbb{E}[\Delta_{t,i} \mid \mathscr{F}_{t-1}] + M_{t,2}$$

$$\begin{aligned} & \leq \frac{1}{2} \sum_{i=1}^{d} \zeta_{i}(t) + \frac{1}{2} \sum_{i=1}^{d} \mathbb{E}[g_{t,i}^{2} \mid \mathscr{F}_{t-1}] \cdot \mathbb{E}[\Delta_{t,i} \mid \mathscr{F}_{t-1}] + M_{t,2} \\ & 1049 \\ 1050 \\ 1051 \\ 1052 \\ 1052 \\ 1053 \\ 1054 \end{aligned} \\ & \leq \frac{1}{2} \sum_{i=1}^{d} \zeta_{i}(t) + \frac{1}{2} \left( \sum_{i=1}^{d} \mathbb{E}[g_{t,i}^{2} \mid \mathscr{F}_{t-1}] \right) \cdot \left( \sum_{i=1}^{d} \mathbb{E}[\Delta_{t,i} \mid \mathscr{F}_{t-1}] \right) + M_{t,2} \\ & \leq \frac{1}{2} \sum_{i=1}^{d} \zeta_{i}(t) + \frac{1}{2} \left( (A + 2L_{f}B)f(w_{t}) + C \right) \cdot \left( \sum_{i=1}^{d} \mathbb{E}[\Delta_{t,i} \mid \mathscr{F}_{t-1}] \right) + M_{t,2} \end{aligned}$$

$$\leq \overline{2} \sum_{i=1}^{d} \zeta_i(t) + \overline{2} \left( (A + 2L_f D) f(w_t) + C \right)^* \left( \sum_{i=1}^{d} \mathbb{E}[\Delta_{t,i} \mid \mathscr{F}_{t-1}] \right) + M$$

$$= \frac{1}{2} \sum_{i=1}^{d} \zeta_i(t) + \frac{1}{2} (A + 2L_f B) f(w_t) \cdot \left( \sum_{i=1}^{d} \mathbb{E}[\Delta_{t,i} \mid \mathscr{F}_{t-1}] \right)$$

$$+ C \left( \sum_{i=1}^{d} \mathbb{E}[\Delta_{t,i} \mid \mathscr{F}_{t-1}] \right) + M_{t,2}$$

$$= \frac{1}{2} \sum_{i=1}^{d} \zeta_i(t) + \frac{1}{2} (A + 2L_f B) f(w_t) \cdot \left( \underbrace{\sum_{i=1}^{d} \mathbb{E}[\Delta_{t,i} \mid \mathscr{F}_{t-1}]}_{\overline{\Delta}_t} \right) + C \sum_{i=1}^{d} \Delta_{t,i}$$

$$+\underbrace{C\left(\sum_{i=1}^{d} \left(\mathbb{E}[\Delta_{t,i} \mid \mathscr{F}_{t-1}] - \Delta_{t,i}\right)\right)}_{M_{t,3}} + M_{t,2}$$
(16)

In the above derivation, in step (a), we utilized the property of conditional expectation, which states that if random variables  $X \in \mathscr{F}_{n-1}$  and  $Y \in \mathscr{F}_n$ , then  $\mathbb{E}[XY|\mathscr{F}_{n-1}] = X \mathbb{E}[Y|\mathscr{F}_{n-1}]$ . Addi-tionally, we need to note that  $\Delta_{t,i} = \sqrt{\Delta_{t,i}} \sqrt{\Delta_{t,i}} < \sqrt{\eta_{v_{t-1}}} \sqrt{\Delta_{t,i}}$  (due to Property 2, we know  $\Delta_{t,i} \geq 0$ , so taking the square root is well-defined). In step (b), we employed the AM-GM inequal-ity, which states  $ab \leq \frac{a^2+b^2}{2}$ . In step (c), we used the *Cauchy-Schwarz* inequality for conditional expectations:  $\mathbb{E}[XY|\mathscr{F}_{n-1}] \leq \sqrt{\mathbb{E}[X^2|\mathscr{F}_{n-1}]\mathbb{E}[Y^2|\mathscr{F}_{n-1}]}$ . For step (d), we used Property 1. Specifically, we have: 

1078  
1079 
$$\sum_{i=1}^{d} \mathbb{E}[g_{t,i}^2 | \mathscr{F}_{t-1}] = \mathbb{E}[||g_t||^2 | \mathscr{F}_{t-1}] \le (A + 2L_f B)f(w_t) + C.$$

Substituting the estimate of  $\Theta_{t,1,1}$  back into Eq. 16, we obtain: 

$$\Theta_{t,1} = -\frac{1}{2} \sum_{i=1}^{d} \zeta_i(t) + \frac{A + 2L_f B}{2} \overline{\Delta}_t \cdot f(w_t) + C \sum_{i=1}^{d} \Delta_{t,i} + \underbrace{M_{t,1} + M_{t,2} + M_{t,3}}_{M_t},$$

Then, we use Property 5 to replace  $f(w_t)$  with  $f(u_t)$  to obtain:

$$\Theta_{t,1} = -\frac{1}{2} \sum_{i=1}^{d} \zeta_i(t) + \frac{(A + 2L_f B)(L_f + 1)}{2} \overline{\Delta}_t \cdot f(u_t) + C \sum_{i=1}^{d} \Delta_{t,i} + \frac{(L_f + 1)\beta_1^2}{2(1 - \beta_1)^2} \|\eta_{v_{t-1}} \circ m_{t-1}\|^2 + M_t,$$
(17)

Next, we deal with  $\Theta_{t,2}$ . Specifically, we have the following derivation: 

$$\begin{aligned}
& \Theta_{t,2} = \frac{1}{2} \sum_{i=1}^{d} (\eta_{v_{t,i}} (\nabla_i f(w_t) - \nabla_i f(u_t)) g_{t,i}) \\
& \Theta_{t,2} = \frac{1}{2} \sum_{i=1}^{d} (\eta_{v_{t,i}} g_{t,i}^2 + \frac{1}{2} \sum_{i=1}^{d} (\nabla_i f(w_t) - \nabla_i f(u_t))^2 \\
& \leq \sum_{i=1}^{d} \eta_{v_{t,i}}^2 g_{t,i}^2 + \frac{1}{2} \|\nabla f(w_t) - \nabla f(u_t)\|^2 \\
& = \sum_{i=1}^{d} \eta_{v_{t,i}}^2 g_{t,i}^2 + \frac{L_f^2}{2} \|w_t - u_t\|^2 \\
& \leq \sum_{i=1}^{d} \eta_{v_{t,i}}^2 g_{t,i}^2 + \frac{\beta_1^2 L_f^2}{2(1 - \beta_1)^2} \|\eta_{v_{t-1}} \circ m_{t-1}\|^2.
\end{aligned}$$
(18)

Next, we deal with  $\Theta_{t,3}$ , and we obtain: 

$$\begin{array}{ll} & \begin{array}{l} 1110\\ 1111\\ 1112\\ 1112\\ 1112\\ 1113\\ 1114\\ 1115\\ 1115\\ 1115\\ 1115\\ 1116\\ 1116\\ 1115\\ 1116\\ 1117\\ 1116\\ 1117\\ 1116\\ 1117\\ 11$$

$$\overset{(a)}{\leq} \frac{3}{16} \sum_{i=1}^{d} \Delta_{t,i} (\nabla_i f(w_t))^2 + \frac{3}{8} \sum_{i=1}^{d} (\nabla_i f(u_t) - \nabla_i f(w_t))^2 + 2 \sum_{i=1}^{d} \Delta_{t,i} m_{t-1,i}^2$$

$$\stackrel{(b)}{\leq} \frac{3}{16} \sum_{i=1}^{d} \zeta_i(t) + \frac{3\beta_1^2 L_f^2}{8(1-\beta_1)^2} \|\eta_{v_{t-1}} \circ m_{t-1}\|^2 + 2\sum_{i=1}^{d} \Delta_{t,i} m_{t-1,i}^2.$$

$$(19)$$

In step (a), we used the following substitution:

In step (b), we used the L-Smooth condition (Assumption 2.2), i.e., 

1129  
1130  
1131  
1132  
1133  

$$\sum_{i=1}^{d} (\nabla_i f(u_t) - \nabla_i f(w_t))^2 = \|\nabla f(w_t) - \nabla f(u_t)\|^2$$

$$\leq L_f^2 \|w_t - u_t\|^2 \stackrel{\text{Eq. 4}}{=} \frac{\beta_1^2 L_f}{(1 - \beta_1)^2} \|\eta_{v_{t-1}} \circ m_{t-1}\|^2.$$

For  $\Theta_{t,4}$ , we have that 

$$\Theta_{t,4} = \sum_{i=1}^{d} \Delta_{t,i}^2 m_{t-1,i}^2 \overset{\Delta_{t,i} \le \eta_{v_{t-1},i}}{<} \sum_{i=1}^{d} \eta_{t,i}^2 m_{t-1,i}^2 = \|\eta_{v_{t-1}} \circ m_{t-1}\|^2.$$
(20)

Finally, substituting the estimates of  $\Theta_{t,1}$  from Eq. 17,  $\Theta_{t,2}$  from Eq. 18,  $\Theta_{t,3}$  from Eq. 19, and  $\Theta_{t,4}$  from Eq. 20 back into Eq. 14, we obtain: 

$$\begin{array}{l} 1141\\ 1142\\ 1143\\ 1143\\ 1144 \end{array} \qquad f(u_{t+1}) - f(u_t) \leq -\frac{5}{16} \sum_{i=1}^d \zeta_i(t) + \frac{(A + 2L_f B)(L_f + 1)}{2} \overline{\Delta}_t \cdot f(u_t) + C \sum_{i=1}^d \eta_{v_{t-1},i} - C \sum_{i=1}^d \eta_{v_t,i} + \frac{1}{2} \sum_{i=1}^d \eta_{v_t,i}$$

$$+\frac{(L_f+1)\beta_1^2}{2(1-\beta_1)^2}\|\eta_{v_{t-1}}\circ m_{t-1}\|^2 + L_f\left(\frac{\beta_1}{1-\beta_1}\right)^2\sum_{i=1}^{\infty}\eta_{v_t,i}^2g_{t,i}^2$$

1147  
1148  
1149  

$$+ \frac{\beta_1^2 L_f^2}{2(1-\beta_1)^2} \|\eta_{v_{t-1}} \circ m_{t-1}\|^2 + L_f \left(\frac{\beta_1}{1-\beta_1}\right)^2 \|\eta_{v_{t-1}} \circ m_{t-1}\|^2$$

1149  
1150  
1151  
1152  

$$+ \frac{3\beta_1^2 L_f^2}{8(1-\beta_1)^2} \|\eta_{v_{t-1}} \circ m_{t-1}\|^2 + 2\sum_{i=1}^d \Delta_{t,i} m_{t-1,i}^2 + L_f \sum_{i=1}^d \eta_{v_t,i}^2 g_{t,i}^2$$

 $+ M_t$ . 

Then we note that  $\Delta_{t,i} < \eta_{v_{t-1},i}$  and  $\eta_{v_{t-1},i}^2 < \frac{1}{\sqrt{v}}\eta_{v_{t-1},i}$ . After simplification, we obtain: 

$$\left( f(u_{t+1}) + C \sum_{i=1}^{d} \eta_{v_{t},i} \right) - \left( f(u_{t}) + C \sum_{i=1}^{d} \eta_{v_{t-1},i} \right) \le -\frac{5}{16} \sum_{i=1}^{d} \zeta_{i}(t) + C_{1} \overline{\Delta}_{t} \cdot f(u_{t})$$
$$+ C_{2} \| \sqrt{\eta_{v_{t-1}}} \circ m_{t-1} \|^{2} + (L_{f} + 1) \sum_{i=1}^{d} \eta_{v_{t-1},i}^{2} g_{t,i}^{2} + M_{t},$$

$$+ C_2 \|\sqrt{\eta_{v_{t-1}}} \circ m_{t-1}\|^2 + (L_f + 1) \sum_{i=1}^d \eta_{v_t,i}^2 g_{t,i}^2 +$$

where 

$$C_{1} := \frac{(A + 2L_{f}B)(L_{f} + 1)}{2},$$

$$C_{2} := \frac{1}{\sqrt{v}} \cdot \left(\frac{7L_{f}^{2}}{8} + \frac{L_{f} + 1}{2}\right) \cdot \frac{\beta_{1}^{2}}{(1 - \beta_{1})^{2}} + 2 + \frac{L_{f}}{\sqrt{v}} \cdot \left(\frac{\beta_{1}}{1 - \beta_{1}}\right)^{2}.$$
(21)

We add the following term to both sides of the above inequality: 

$$\frac{C_2}{(1-\beta_1)\sqrt{v}}(\|m_t\|^2 - \|m_{t-1}\|^2),$$

to obtain:

1188  

$$\leq -C_2 \|\sqrt{\eta_{v_{t-1}}} \circ m_{t-1}\|^2 + C_2 \|\sqrt{\eta_{v_{t-1}}} \circ m_{t-1}\|^2 = 0.$$
1189

1190 We apply an obvious inequality to the second term on the right side of Eq. 22:

$$f(u_t) < f(u_t) + C \sum_{i=1}^d \eta_{v_{t-1},i} + \frac{C_2}{(1-\beta_1)\sqrt{v}} ||m_{t-1}||^2 = \hat{f}(u_t),$$

and then move the expanded term to the left side of the inequality and combine like terms to obtain:

$$\hat{f}(u_{t+1}) - (1 + C_1 \overline{\Delta}_t) \hat{f}(u_t) \le -\frac{5}{16} \sum_{i=1}^d \zeta_i(t) + (L_f + 1) \sum_{i=1}^d \eta_{v_t,i}^2 g_{t,i}^2 + M_t.$$

1199 Next, we construct an auxiliary variable

$$\Pi_{\Delta,t} := \prod_{k=1}^{t-1} (1 + C_1 \overline{\Delta}_k)^{-1} \quad (\forall t \ge 2), \ \Pi_{\Delta,1} := 1.$$

1204 Multiplying both sides of the above inequality by  $\Pi_{\Delta,t+1}$ , we obtain

$$\Pi_{\Delta,t+1}\hat{f}(u_{t+1}) - \Pi_{\Delta,t}\hat{f}(u_t) \le -\frac{5}{16}\Pi_{\Delta,t+1}\sum_{i=1}^d \zeta_i(t) + (L_f+1)\Pi_{\Delta,t+1}\sum_{i=1}^d \eta_{v_t,i}^2 g_{t,i}^2 + \Pi_{\Delta,t+1}M_t.$$

1210 With this, we complete the proof.

#### 1213 D.2 THE PROOF OF LEMMA B.1

*Proof.* For  $\forall x \in \mathbb{R}^N$ , we define function

$$g(t) = f\left(x + t\frac{x' - x}{\|x' - x\|}\right),$$

1219 where x' is a constant point such that x' - x is parallel to  $\nabla f(x)$ . By taking the derivative, we obtain

$$g'(t) = \nabla_{x+t\frac{x'-x}{\|x'-x\|}} f\left(x+t\frac{x'-x}{\|x'-x\|}\right)^T \frac{x'-x}{\|x'-x\|}.$$
(23)

1223 Through the Lipschitz condition of  $\nabla f(x)$ , we get  $\forall t_1, t_2$ 

$$|g'(t_1) - g'(t_2)| = \left| \left( \nabla_{x+t \frac{x'-x}{\|x'-x\|}} f\left(x+t_1 \frac{x'-x}{\|x'-x\|}\right) - \nabla_{x+t \frac{x'-x}{\|x'-x\|}} f\left(x+t_2 \frac{x'-x}{\|x'-x\|}\right) \right)^T \frac{x'-x}{\|x'-x\|}$$

$$\leq \left\| \nabla_{x+t\frac{x'-x}{\|x'-x\|}} f\left(x+t_1\frac{x'-x}{\|x'-x\|}\right) - \nabla_{x+t\frac{x'-x}{\|x'-x\|}} f\left(x+t_2\frac{x'-x}{\|x'-x\|}\right) \right\| \left\| \frac{x'-x}{\|x'-x\|} \right\| \leq \mathcal{L}|t_1-t_2|.$$

So g'(t) satisfies the Lipschitz condition, and we have  $\inf_{t \in \mathbb{R}} g(t) \ge \inf_{x \in \mathbb{R}^N} f(x) > -\infty$ . Let  $g^* = \inf x \in \mathbb{R} g(x)$ , then it holds that for  $\forall t_0 \in \mathbb{R}$ ,

$$g(0) - g^* \ge g(0) - g(t_0).$$
 (24)

By using the Newton-Leibniz's formula, we get that

$$g(0) - g(t_0) = \int_{t_0}^0 g'(\alpha) d\alpha = \int_{t_0}^0 \left( g'(\alpha) - g'(0) \right) d\alpha + \int_{t_0}^0 g'(0) d\alpha.$$

1239 Through the Lipschitz condition of g', we get that 1240

$$g(0) - g(t_0) \ge \int_{t_0}^0 -\mathcal{L} |\alpha - 0| d\alpha + \int_{t_0}^0 g'(0) d\alpha = \frac{1}{2\mathcal{L}} (g'(0))^2.$$

Then we take a special value of  $t_0$ . Let  $t_0 = -g'(0)/\mathcal{L}$ , then we get 

$$g(0) - g(t_0) \ge -\int_{t_0}^0 \mathcal{L}|\alpha| d\alpha + \int_{t_0}^0 g(0) dt = -\frac{\mathcal{L}}{2} (0 - t_0)^2 + g'(0)(-t_0)$$
  
$$= -\frac{1}{2\mathcal{L}} (g'(0))^2 + \frac{1}{\mathcal{L}} (g'(0))^2 = \frac{1}{2\mathcal{L}} (g'(0))^2.$$
(25)

Substituting Eq. 25 into Eq. 24, we get

$$g(0) - g^* \ge \frac{1}{2\mathcal{L}} (g'(0))^2.$$

Due to  $g^* \ge f^*$  and  $(g'(0))^2 = \|\nabla f(x)\|^2$ , it follows that 

$$\left\|\nabla f(x)\right\|^2 \le 2\mathcal{L}\big(f(x) - f^*\big)$$

 $]^{p}$ 

#### D.3 THE PROOF OF LEMMA B.2

*Proof.* (i) Consider the non-negative adapted process  $\{X_n, \mathscr{F}_n\}_{n \ge 1}$  and define the partial sum of conditional expectations as  $\Lambda_T := \sum_{n=1}^T \mathbb{E}[X_n \mid \mathscr{F}_{n-1}].$ 

First, we compute the expectation of  $\Lambda_T$ : 

$$\mathbb{E}[\Lambda_T] = \mathbb{E}\left[\sum_{n=1}^T \mathbb{E}[X_n \mid \mathscr{F}_{n-1}]\right] = \sum_{n=1}^T \mathbb{E}[X_n] \le M.$$

Since  $X_n$  are non-negative, we know that  $\Lambda_T$  is a non-decreasing sequence, and considering that  $\mathbb{E}(\Lambda_T)$   $(\forall T \ge 1)$  is also bounded by M, we can apply the Lebesgue's Monotone Convergence theorem.

Thus,  $\Lambda_T$  converges almost surely to a limit  $\Lambda$ : 

$$\Lambda := \lim_{T \to \infty} \Lambda_T = \sum_{n=1}^{\infty} \mathbb{E}[X_n \mid \mathscr{F}_{n-1}]$$
 a.s.

This concludes that the sequence of conditional expectation sums converges almost surely. 

(ii) We begin by normalizing  $X_n$  by considering the expression  $Y_n = \frac{X_n}{2M}$ . According to the Lebesgue's monotone convergence theorem, we only need to prove that 

$$\forall p \ge 1, \ \mathbb{E}\left[\sum_{n=1}^{\infty} \mathbb{E}[Y_n | \mathscr{F}_{n-1}]\right]^p := M(p) < +\infty$$

Next, we proceed with the calculation, and we obtain  $\forall p \ge 2$ , there is: 

$$\begin{array}{ll} 1285\\ 1286\\ 1287\\ 1287\\ 1287\\ 1289\\ 1289 \end{array} \qquad M(p) = \mathbb{E}\left[\sum_{n=1}^{\infty} \mathbb{E}[Y_n|\mathscr{F}_{n-1}]\right]^p = \mathbb{E}\left[\sum_{n=1}^{\infty} Y_n + \sum_{n=1}^{\infty} (\mathbb{E}[Y_n|\mathscr{F}_{n-1}] - Y_n)\right]^p \\ \overset{(a)}{\leq} \mathbb{E}\left[\frac{1}{2} + \sum_{n=1}^{\infty} (\mathbb{E}[Y_n|\mathscr{F}_{n-1}] - Y_n)\right]^p \overset{(b)}{\leq} 2^{p-1} \left(\frac{1}{2^p} + \mathbb{E}\left[\sum_{n=1}^{\infty} (\mathbb{E}[Y_n|\mathscr{F}_{n-1}] - Y_n)\right]^p \right]^p \\ \overset{(b)}{\leq} 2^{p-1} \left(\frac{1}{2^p} + \mathbb{E}\left[\sum_{n=1}^{\infty} (\mathbb{E}[Y_n|\mathscr{F}_{n-1}] - Y_n)\right]^p \right]^p \\ \overset{(b)}{\leq} 2^{p-1} \left(\frac{1}{2^p} + \mathbb{E}\left[\sum_{n=1}^{\infty} (\mathbb{E}[Y_n|\mathscr{F}_{n-1}] - Y_n)\right]^p \right]^p \\ \overset{(b)}{\leq} 2^{p-1} \left(\frac{1}{2^p} + \mathbb{E}\left[\sum_{n=1}^{\infty} (\mathbb{E}[Y_n|\mathscr{F}_{n-1}] - Y_n)\right]^p \right]^p \\ \overset{(b)}{\leq} 2^{p-1} \left(\frac{1}{2^p} + \mathbb{E}\left[\sum_{n=1}^{\infty} (\mathbb{E}[Y_n|\mathscr{F}_{n-1}] - Y_n)\right]^p \right]^p \\ \overset{(b)}{\leq} 2^{p-1} \left(\frac{1}{2^p} + \mathbb{E}\left[\sum_{n=1}^{\infty} (\mathbb{E}[Y_n|\mathscr{F}_{n-1}] - Y_n)\right]^p \right]^p \\ \overset{(b)}{\leq} 2^{p-1} \left(\frac{1}{2^p} + \mathbb{E}\left[\sum_{n=1}^{\infty} (\mathbb{E}[Y_n|\mathscr{F}_{n-1}] - Y_n)\right]^p \right]^p \\ \overset{(b)}{\leq} 2^{p-1} \left(\frac{1}{2^p} + \mathbb{E}\left[\sum_{n=1}^{\infty} (\mathbb{E}[Y_n|\mathscr{F}_{n-1}] - Y_n)\right]^p \right]^p \\ \overset{(b)}{\leq} 2^{p-1} \left(\frac{1}{2^p} + \mathbb{E}\left[\sum_{n=1}^{\infty} (\mathbb{E}[Y_n|\mathscr{F}_{n-1}] - Y_n\right)\right]^p \\ \overset{(b)}{\leq} 2^{p-1} \left(\frac{1}{2^p} + \mathbb{E}\left[\sum_{n=1}^{\infty} (\mathbb{E}[Y_n|\mathscr{F}_{n-1}] - Y_n\right)\right]^p \\ \overset{(b)}{\leq} 2^{p-1} \left(\frac{1}{2^p} + \mathbb{E}\left[\sum_{n=1}^{\infty} (\mathbb{E}[Y_n|\mathscr{F}_{n-1}] - Y_n\right)\right]^p \\ \overset{(b)}{\leq} 2^{p-1} \left(\frac{1}{2^p} + \mathbb{E}\left[\sum_{n=1}^{\infty} (\mathbb{E}[Y_n|\mathscr{F}_{n-1}] - Y_n\right)\right]^p \\ \overset{(b)}{\leq} 2^{p-1} \left(\frac{1}{2^p} + \mathbb{E}\left[\sum_{n=1}^{\infty} (\mathbb{E}[Y_n|\mathscr{F}_{n-1}] - Y_n\right)\right]^p \\ \overset{(b)}{\leq} 2^{p-1} \left(\frac{1}{2^p} + \mathbb{E}\left[\sum_{n=1}^{\infty} (\mathbb{E}[Y_n|\mathscr{F}_{n-1}] - Y_n\right)\right]^p \\ \overset{(b)}{\leq} 2^{p-1} \left(\frac{1}{2^p} + \mathbb{E}\left[\sum_{n=1}^{\infty} (\mathbb{E}[Y_n|\mathscr{F}_{n-1}] - Y_n\right)\right]^p \\ \overset{(b)}{\leq} 2^{p-1} \left(\frac{1}{2^p} + \mathbb{E}\left[\sum_{n=1}^{\infty} (\mathbb{E}[Y_n|\mathscr{F}_{n-1}] - Y_n\right)\right]^p \\ \overset{(b)}{\leq} 2^{p-1} \left(\frac{1}{2^p} + \mathbb{E}\left[\sum_{n=1}^{\infty} (\mathbb{E}[Y_n|Y_n|Y_n] - Y_n\right)\right]^p \\ \overset{(b)}{\leq} 2^{p-1} \left(\frac{1}{2^p} + \mathbb{E}\left[\sum_{n=1}^{\infty} (\mathbb{E}[Y_n|Y_n|Y_n] - Y_n\right)\right]^p \\ \overset{(b)}{\leq} 2^{p-1} \left(\frac{1}{2^p} + \mathbb{E}\left[\sum_{n=1}^{\infty} (\mathbb{E}[Y_n|Y_n|Y_n] - Y_n\right)\right]^p \\ \overset{(b)}{\leq} 2^{p-1} \left(\frac{1}{2^p} + \mathbb{E}\left[\sum_{n=1}^{\infty} (\mathbb{E}[Y_n|Y_n] - Y_n\right)\right]^p \\ \overset{(b)}{\leq} 2^{p-1} \left(\frac{1}{2^p} + \mathbb{E}\left[\sum_{n=1}^{\infty} (\mathbb{E}[Y_n|Y_$$

$$\leq \mathbb{E}\left[\frac{1}{2} + \sum_{n=1}^{\infty} (\mathbb{E}[Y_n|\mathscr{F}_{n-1}] - Y_n)\right] \leq 2^{p-1} \left(\frac{1}{2^p} + \mathbb{E}\left[\sum_{n=1}^{\infty} (\mathbb{E}[Y_n|\mathscr{F}_{n-1}] - Y_n)\right]\right)$$

$$\leq \frac{1}{2} + 2^{p-1}C_p \mathbb{E}\left[\sum_{n=1}^{\infty} |\mathbb{E}[Y_n|\mathscr{F}_{n-1}] - Y_n|^2\right]^{p/2} \leq \frac{1}{2} + 2^{p-1}C_p \mathbb{E}\left[\sum_{n=1}^{\infty} |\mathbb{E}[Y_n|\mathscr{F}_{n-1}] - Y_n|\right]^{p/2}$$

$$\stackrel{(c)}{\leq} \frac{1}{2} + 2^{p-1}C_p \mathbb{E}\left[\sum_{n=1}^{\infty} |\mathbb{E}[Y_n|\mathscr{F}_{n-1}] - Y_n|^2\right] \stackrel{(d)}{\leq} \frac{1}{2} + 2^{p-1}C_p \mathbb{E}\left[\sum_{n=1}^{\infty} |\mathbb{E}[Y_n|\mathscr{F}_{n-1}] - Y_n|^2\right]$$

$$\stackrel{(f)}{\leq} \frac{1}{2} + 2^{p-2}C_p + 2^{\frac{3}{2}p-2}C_p \mathbb{E}\left[\sum_{n=1}^{\infty} \mathbb{E}[Y_n|\mathscr{F}_{n-1}]\right]^{p/2}$$

1296  
1297 
$$= \frac{1}{2} + 2^{p-2}C_p + 2^{\frac{3}{2}p-2}C_p M(p/2).$$
(26)

In the above derivation, inequality (a) requires noting that 

Inequality (b) uses the AM-GM inequality, specifically,

$$\left(\frac{a+b}{2}\right)^p \le \frac{a^p + b^p}{2}$$

Inequality (c) involves using Burkholder's inequality <sup>1</sup>, where  $C_p$  is a constant depending only on p, and its order with respect to p is  $\mathcal{O}(p)$ . Inequality (d) requires noting that 

$$|\mathbb{E}[Y_n|\mathscr{F}_{n-1}] - Y_n|^2 \le |\mathbb{E}[Y_n|\mathscr{F}_{n-1}] - Y_n|.$$

-.

 $\sum_{n=1}^{+\infty} Y_n = \frac{1}{2}.$ 

By repeatedly iterating Eq. 26 and using the fact that  $C_p = \mathcal{O}(p)$ , we can finally obtain the following estimate: 

$$M(p) = o(p^{\sqrt{p}}).$$

 $\mathbb{E}[\Lambda^p] = o((2M)^p \cdot p^{\sqrt{p}})$ 

that is, 

(iii) Similarly, we first consider the partial sum sequence  $\Lambda_T$ . Since there exists an obvious upper bound for this sequence that depends on T, given by  $\Lambda_T < M \cdot T \ a.s.$ , and because the Taylor series of the function  $h_p(x) = e^{px}$  at x = 0 converges uniformly on every compact subset of  $\mathbb{R}$ , we can expand  $e^{\Lambda_T}$  using its Taylor series, yielding the following expression:

1321  
1322  
1323 
$$e^{\Lambda_T} = \sum_{n=0}^{\infty} \frac{\Lambda_T^n}{n!}$$

Since the Taylor series converges uniformly on every compact subset, we can interchange the sum-mation and the expectation operators when taking the expectation on both sides of the above equa-tion. Specifically, we have: 

$$\mathbb{E}[e^{p\Lambda_T}] = \mathbb{E}\left[\sum_{n=0}^{\infty} \frac{p^n \Lambda_T^n}{n!}\right] = \sum_{n=0}^{\infty} \frac{p^n \mathbb{E}[\Lambda_T^n]}{n!}$$

Noting that  $\mathbb{E}[\Lambda_T^n] < \mathbb{E}[\Lambda^n] < o((2M)^n \cdot n^{\sqrt{n}})$ , we have: 

$$\mathbb{E}[e^{p\Lambda_T}] = \sum_{n=0}^{\infty} \frac{p^n \mathbb{E}[\Lambda_T^n]}{n!} \le \sum_{n=0}^{\infty} \frac{o((2pM)^n \cdot n^{\sqrt{n}})}{n!} = \mathcal{O}\bigg(\sum_{n=0}^{\infty} \frac{(2pM)^n \cdot n^{\sqrt{n}}}{n!}\bigg).$$

We use *Stirling's* approximation  $^2$  to substitute the factorial in the denominator. It is evident that the series inside the  $\mathcal{O}$  notation converges and depends only on p and M. Then, by applying the Lebesgue's Monotone Convergence theorem, we can prove that  $\mathbb{E}[e^{p\Lambda}]$  exists, and its upper bound depends only on p and M. 

#### With this, we complete the proof.

<sup>1</sup>Burkholder's inequality: For any martingale  $(M_n, \mathscr{F}_n)$  with  $M_0 = 0$  almost surely, and for any  $1 \le p < \infty$  $\infty$ , there exist constants  $c_p > 0$  and  $C_p > 0$  depending only on p such that:

$$c_p \mathbb{E}[(S(M))^p] \le \mathbb{E}[(M^*)^p] \le C_p \mathbb{E}[(S(M))^p],$$

where  $M^* = \sup_{n \ge 0} |M_n|$  and  $S(M) = \left(\sum_{i \ge 1} (M_i - M_{i-1})^2\right)^{1/2}$ . 

1348 
$$n! = \sqrt{2\pi n} \left(\frac{n}{e}\right)^n e^{\frac{\theta_n}{12n}}$$
1349

where  $0 < \theta_n < 1$ .

#### D.4 THE PROOF OF LEMMA B.3

*Proof.* First, using the expression in Eq. 7, we express  $\Pi_{\Delta,T+1}^{-p}$  as follows:

$$\Pi_{\Delta,T+1}^{-p} = \prod_{k=1}^{T} (1 + C_1 \overline{\Delta}_k)^p.$$

By applying the logarithmic transformation, we obtain: 

$$\Pi_{\Delta,T+1}^{-p} = \exp\left\{p\sum_{k=1}^{T}\ln(1+C_1\overline{\Delta}_k)\right\} \stackrel{\ln(1+x)\leq x \ (\forall \ x>0)}{\leq} \exp\left\{pC_1\sum_{k=1}^{T}\overline{\Delta}_k\right\}.$$

Based on the definition of  $\overline{\Delta}_k$ , we have 

$$\overline{\Delta}_{k} = \sum_{i=1}^{d} \mathbb{E}[\Delta_{k,i} | \mathscr{F}_{k-1}] = \mathbb{E}\left[\sum_{i=1}^{d} \Delta_{k,i} \middle| \mathscr{F}_{k-1}\right].$$

 $\sum_{k=1}^{+\infty} \sum_{i=1}^{d} \Delta_{k,i} < \frac{d}{\sqrt{v}},$ 

Since  $\sum_{i=1}^{d} \Delta_{k,i} > 0$  and 

we can use Lemma B.2 to prove that: 

$$\mathbb{E}\left[\exp\left\{pC_1\sum_{k=1}^T\overline{\Delta}_k\right\}\right] < C_{v,d,p}$$

which implies 

$$\mathbb{E}[\Pi_{\Delta,T+1}^{-p}] < C_{v,d,p}$$

Next, noting the monotonicity of the sequence  $\{\Pi_{\Delta,T+1}^{-1}\}_{T\geq 1}$ , we take  $T \to +\infty$  and apply the Lebesgue's Monotone Convergence theorem to obtain the result regarding  $\Pi_{\Delta,\infty}$ . 

With this, the proof is complete. 

#### D.5 THE PROOF OF LEMMA B.4

*Proof.* Recalling the approximate descent inequality 4.1, we have: 

$$\Pi_{\Delta,t+1}\hat{f}(u_{t+1}) - \Pi_{\Delta,t}\hat{f}(u_t) \le -\frac{5}{16}\Pi_{\Delta,t+1}\sum_{i=1}^d \zeta_i(t) + (L_f+1)\Pi_{\Delta,t+1}\sum_{i=1}^d \eta_{v_t,i}^2 g_{t,i}^2 + \Pi_{\Delta,t+1}M_t.$$

We multiply both sides of the above equation by  $\frac{1}{\sqrt{S_{t-1}}}$ , we obtain:

$$\frac{\Pi_{\Delta,t+1}\hat{f}(u_{t+1})}{\sqrt{S_t}} - \frac{\Pi_{\Delta,t}\hat{f}(u_t)}{\sqrt{S_{t-1}}} \le (L_f+1)\frac{\Pi_{\Delta,t+1}}{\sqrt{S_{t-1}}}\sum_{i=1}^d \eta_{v_t,i}^2 g_{t,i}^2 + \frac{\Pi_{\Delta,t+1}}{\sqrt{S_{t-1}}}M_t$$
$$\stackrel{\Pi_{\Delta,t+1}\le 1}{\le} (L_f+1)\frac{1}{\sqrt{S_{t-1}}}\sum_{i=1}^d \eta_{v_t,i}^2 g_{t,i}^2 + \frac{\Pi_{\Delta,t+1}}{\sqrt{S_{t-1}}}M_t$$

We sum the above inequality with respect to the index t from 1 to T, and we obtain: 

$$\frac{\Pi_{\Delta,t+1}\hat{f}(u_{t+1})}{\sqrt{S_t}} \le \frac{\hat{f}(u_1)}{\sqrt{S_0}} + (L_f+1)\sum_{t=1}^T \underbrace{\frac{1}{\sqrt{S_{t-1}}}\sum_{i=1}^d \eta_{v_t,i}^2 g_{t,i}^2}_{\Gamma_{t,3}} + \sum_{t=1}^T \frac{\Pi_{\Delta,t+1}}{\sqrt{S_{t-1}}} M_t.$$
(27)

For  $\Gamma_{t,3}$ , we can perform the following simplification, and we obtain: 

$$\frac{1}{\sqrt{S_{t-1}}} \sum_{i=1}^d \eta_{v_t,i}^2 g_{t,i}^2 = \frac{1}{\sqrt{S_t}} \sum_{i=1}^d \eta_{v_t,i}^2 g_{t,i}^2 + \left(\frac{1}{\sqrt{S_{t-1}}} - \frac{1}{\sqrt{S_t}}\right) \sum_{i=1}^d \eta_{v_t,i}^2 g_{t,i}^2$$

$$\begin{array}{c} 1408 \\ 1409 \\ 1410 \\ 1$$

By using the series-integral inequality, we obtain:

$$\begin{split} \sum_{t=1}^{T} \Gamma_{t,3} &\leq \frac{1}{\alpha_1} \sum_{i=1}^{d} \sum_{t=1}^{T} \frac{g_{t,i}^2}{S_{t,i}^{3/2}} + \frac{d}{\alpha_1} \sum_{t=1}^{T} \left( \frac{1}{\sqrt{S_{t-1}}} - \frac{1}{\sqrt{S_t}} \right) \\ &\leq \frac{1}{\alpha_1} \sum_{i=1}^{d} \int_{S_{0,i}}^{+\infty} \frac{1}{x^{3/2}} \mathrm{d}x + \frac{d}{\alpha_1 \sqrt{S_0}} = \frac{2d}{\alpha_1 \sqrt{v}} + \frac{1}{\alpha_1} \sqrt{\frac{d}{v}} \\ &< \frac{3d}{\alpha_1 \sqrt{v}}. \end{split}$$

Substituting above inequality into Eq. 27, we acquire

$$\frac{\prod_{\Delta,t+1}\hat{f}(u_{t+1})}{\sqrt{S_t}} \le \frac{\hat{f}(u_1)}{\sqrt{dv}} + \frac{3(L_f+1)d}{\alpha_1\sqrt{v}} + \sum_{t=1}^T \frac{\prod_{\Delta,t+1}}{\sqrt{S_{t-1}}} M_t.$$

With this, we complete the proof.

## D.6 THE PROOF OF LEMMA B.5

*Proof.* For any  $\phi \in \mathbb{R}$ , we consider  $\frac{\sqrt{S_T}}{(T+1)^{\phi}}$ , and we obtain: 

$$\frac{\sqrt{S_T}}{(T+1)^{\phi}} = \frac{S_T}{(T+1)^{\phi}\sqrt{S_T}} = \frac{S_0 + \sum_{t=1}^T \|g_t\|^2}{(T+1)^{\phi}\sqrt{S_T}} = \frac{S_0}{(T+1)^{\phi}\sqrt{S_T}} + \sum_{t=1}^T \frac{\|g_t\|^2}{(T+1)^{\phi}\sqrt{S_T}}$$

$$\leq \frac{S_0}{(T+1)^{\phi}\sqrt{S_T}} + \sum_{t=1}^T \frac{\|g_t\|^2}{(T+1)^{\phi}\sqrt{S_T}} \leq \sqrt{S_0} + \sum_{t=1}^T \frac{\|g_t\|^2}{(t+1)^{\phi}\sqrt{S_{t-1}}}$$
$$= \sqrt{dv} + \sum_{t=1}^T \frac{\|g_t\|^2}{(t+1)^{\phi}\sqrt{S_{t-1}}}.$$

Next, by multiplying both sides of the above inequality by  $\Pi_{\Delta,T}$  and noting the monotonicity of  ${\Pi_{\Delta,t}}_{t\geq 1}$  as well as the fact that  $\Pi_{\Delta,T} \leq 1$  for all  $T \geq 1$ , we immediately obtain the result.

#### D.7 THE PROOF OF LEMMA B.6

*Proof.* We take  $\phi = 2$  in Lemma B.5 and bound the expectation of the partial sum  $\sum_{t=1}^{T} \Lambda_{2,t}$ . We have:

1456  
1457 
$$\mathbb{E}\left[\sum_{t=1}^{T}\Pi_{\Delta,t}\Lambda_{2,t}\right] = \sum_{t=1}^{T}\mathbb{E}[\Pi_{\Delta,t}\Lambda_{2,t}] = \sum_{t=1}^{T}\mathbb{E}\left[\frac{\Pi_{\Delta,t}\|g_t\|^2}{(t+1)^2\sqrt{S_{t-1}}}\right] = \sum_{t=1}^{T}\mathbb{E}\left[\frac{\Pi_{\Delta,t}\mathbb{E}[\|g_t\|^2|\mathscr{F}_{t-1}]}{(t+1)^2\sqrt{S_{t-1}}}\right]$$

1458  
1459  
1460  
1461  
Property 1 
$$\sum_{t=1}^{T} \mathbb{E}\left[\frac{(A+2L_fB)\Pi_{\Delta,t}(f(w_t)-f^*)+C}{(t+1)^2\sqrt{S_{t-1}}}\right]$$
  
1460  
1461  
Property 1  $\sum_{t=1}^{T} \mathbb{E}\left[\frac{(A+2L_fB)\Pi_{\Delta,t}(f(w_t)-f^*)+C}{(t+1)^2\sqrt{S_{t-1}}}\right]$ 

$$\frac{\Pr_{\text{property 5}} \sum_{t=1}^{T} \mathbb{E} \left[ \frac{(A+2L_f B) \Pi_{\Delta,t} ((L_f+1)f(u_t) + \frac{(L_f+1)\beta_1}{2(1-\beta_1)^2} \|\eta_{v_{t-1}} \circ m_{t-1}\|^2 - f^*) + C}{(t+1)^2 \sqrt{S_{t-1}}} \right] \\
\leq C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right] + C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right]} \right] + C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right]} + C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right]} + C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right]} \right] + C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right]} + C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right]} + C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right]} + C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right]} + C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right]} + C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right]} + C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right]} + C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right]} + C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right]} + C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right]} + C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right]} + C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right]} + C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right]} + C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right]} + C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right]} + C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right]} + C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right]} + C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right]} + C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right]} + C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right]} + C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right]} + C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right]} + C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right]} + C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right]} + C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right]} + C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right]} + C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right]} + C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right]} + C \sum_{t=1}^{T} \frac{1}{1 - \mathbb{E} \left[ \Pi_{\Delta,t} \hat{f}(u_t) \right]} +$$

$$\leq C_3 \sum_{t=1}^{1} \frac{1}{(t+1)^2} \mathbb{E} \left[ \frac{\Pi_{\Delta,t} f(u_t)}{\sqrt{S_{t-1}}} \right] + C_4 \sum_{t=1}^{1} \frac{1}{(t+1)^2},$$
(28)

where

$$C_3 := (A + 2L_f B) \max\left\{L_f + 1, \frac{(L_f + 1)\beta_1^2}{2C_2(1 - \beta_1)}\right\}, \ C_4 := \frac{(A + 2L_f B)|f^*| + C}{\sqrt{S_0}}.$$

Based on the results in Lemma B.4, we can compute:

$$\mathbb{E}\left[\frac{\Pi_{\Delta,t}\hat{f}(u_t)}{\sqrt{S_{t-1}}}\right] \le \frac{\mathbb{E}[\hat{f}(u_1)]}{\sqrt{dv}} + \frac{3(L_f+1)d}{\alpha_1\sqrt{v}} + 0.$$

1476 Substitute above result into Eq. 28, and combine

$$\sum_{t=1}^{T} \frac{1}{(t+1)^2} \le \sum_{t=1}^{+\infty} \frac{1}{t^2} = \frac{\pi^2}{6}$$

1480 we get: 

$$\mathbb{E}\left[\sum_{t=1}^{T} \Pi_{\Delta,t} \Lambda_{2,t}\right] \leq \frac{C_3 \pi^2 \mathbb{E}[\hat{f}(u_1)]}{6\sqrt{dv}} + \frac{\pi^2 C_3 (L_f + 1)d}{2\alpha_1 \sqrt{v}} + \frac{\pi^2 C_4}{6}.$$

1484 It can be observed that the right side of the above inequality is independent of T. Thus, according 1485 to the *Lebesgue's Monotone Convergence* theorem, we have

$$\sum_{t=1}^{T} \Pi_{\Delta,t} \Lambda_{2,t} \to \sum_{t=1}^{+\infty} \Pi_{\Delta,t} \Lambda_{2,t} \quad \text{a.s.}$$

1490 and

$$\mathbb{E}\left[\sum_{t=1}^{+\infty}\Pi_{\Delta,t}\Lambda_{2,t}\right] = \lim_{T\to\infty}\mathbb{E}\left[\sum_{t=1}^{T}\Pi_{\Delta,t}\Lambda_{2,t}\right] = \lim_{T\to\infty}\sum_{t=1}^{T}\mathbb{E}[\Pi_{\Delta,t}\Lambda_{2,t}] \le \frac{C_3\pi^2\,\mathbb{E}[\hat{f}(u_1)]}{6\sqrt{dv}} + \frac{\pi^2C_3(L_f+1)d}{2\alpha_1\sqrt{v}} + \frac{\pi^2C_4}{6}.$$

Next, we set

$$\zeta := \sqrt{dv} + \sum_{t=1}^{+\infty} \Pi_{\Delta,t} \Lambda_{2,t},$$

and combine Lemma B.5. We get:

$$\sqrt{S_T} \le \Pi_{\Delta,T}^{-1} (T+1)^2 \zeta < \Pi_{\Delta,\infty}^{-1} (T+1)^2 \zeta.$$
(29)

1501 Meanwhile,

$$\mathbb{E}[\zeta] = \sqrt{dv} + \mathbb{E}\left[\sum_{t=1}^{+\infty} \Lambda_{2,t}\right] \le \sqrt{dv} + \frac{C_3 \pi^2 \mathbb{E}[\hat{f}(u_1)]}{6\sqrt{dv}} + \frac{\pi^2 C_3 (L_f + 1)d}{2\alpha_1 \sqrt{v}} + \frac{\pi^2 C_4}{6}.$$
 (30)

Then through Eq. 29, we have

Next, we note that 

$$\ln\left(\max\left\{e, \frac{\sqrt{2}\Pi_{\Delta,\infty}^{-1}\zeta}{\sqrt{v}}\right\}\right) \ge 1, \quad \text{and} \quad \frac{1}{2} \le \ln(T+1),$$

from which we obtain:

$$\frac{1}{2}\ln(S_T) \le 2\ln\left(\max\left\{e, \frac{\sqrt{2}\Pi_{\Delta,\infty}^{-1}\zeta}{\sqrt{v}}\right\}\right) \left(\frac{\ln(T+1)}{\ln\left(\max\left\{e, \frac{\sqrt{2}\Pi_{\Delta,\infty}^{-1}\zeta}{\sqrt{v}}\right\}\right)} + \frac{1}{2}\right)$$

$$\leq 4\ln\left(\max\left\{e,\frac{\sqrt{2}\Pi_{\Delta,\infty}^{-1}\zeta}{\sqrt{v}}\right\}\right)\ln(T+1).$$

With this, we complete the proof.

#### D.8 THE PROOF OF LEMMA B.7

*Proof.* For any T > 0, taking the expectation on both sides of the recursive inequality in the Sufficient Decreasing lemma (Lemma 4.1) and summing the indices from n = 1 to T, noting  $\mathbb{E}[\Pi_{\Delta,t+1}M_t] = \mathbb{E}[\Pi_{\Delta,t+1}\mathbb{E}[M_t|\mathscr{F}_{t-1}]] = 0$ , we obtain: 

To prove the conclusion of this lemma, we actually only need to bound  $\sum_{t=1}^{T} \mathbb{E}[\Gamma_t]$ . Specifically, we perform the following transformation on  $\Gamma_t$ . We have: 

$$\Gamma_{t} = \Pi_{\Delta,t+1} \sum_{i=1}^{d} \eta_{v_{t},i}^{2} g_{t,i}^{2} = \Pi_{\Delta,t+1} \sum_{i=1}^{d} \frac{\eta_{t}^{2} g_{t,i}^{2}}{(\sqrt{v_{t,i}} + \mu)^{2}} \leq \Pi_{\Delta,t+1} \sum_{i=1}^{d} \frac{1}{t^{2\delta}} \frac{g_{t,i}^{2}}{t_{v_{t,i}}}$$
Property 3 (t+1)<sup>2δ</sup>

$$\prod_{\Delta,t+1} \sum_{i=1}^{d} \frac{1}{t^{2\delta}} \prod_{\Delta,t+1} \sum_{i=1}^{d} \frac{1}{t^{2\delta}} \frac{g_{t,i}^{2}}{(\sqrt{v_{t,i}} + \mu)^{2}} \leq \Pi_{\Delta,t+1} \sum_{i=1}^{d} \frac{1}{t^{2\delta}} \frac{g_{t,i}^{2}}{t_{v_{t,i}}}$$

$$\stackrel{\simeq}{=} \frac{1}{t^{2\delta}} \prod_{\Delta,t+1} \sum_{i=1}^{d} \overline{\alpha_1(t+1)^{2\delta}} \overline{S_{t,i}} \stackrel{\simeq}{=} \frac{1}{\alpha_1} \zeta \prod_{\Delta,t+1} \sum_{i=1}^{d} \overline{S_{t,i}^{1+\frac{\delta}{2}}}$$

$$\stackrel{\Pi_{\Delta,t+1}\leq 1}{\leq} \frac{2^{2\delta}}{\alpha_1} \zeta^{\delta} \sum_{i=1}^{d} \frac{g_{t,i}^2}{S_{t,i}^{1+\frac{\delta}{2}}}.$$

$$(32)$$

In step (a) of the above derivation, we need to apply Lemma B.6 to  $(t+1)^{2\delta}$ . Specifically, according to Lemma B.6, we have

$$\sqrt{S_t} \le (t+1)^2 \zeta,$$

which means: 

$$\frac{1}{(t+1)^{2\delta}} \le \frac{\zeta^{\delta}}{(t+1)^{2\delta}S_t^{\frac{\delta}{2}}} \le \frac{\zeta^{\delta}}{(t+1)^{2\delta}S_{t,i}^{\frac{\delta}{2}}} \Pi_{\Delta,t}^{-1}$$

Next, with the estimate for  $\Gamma_t$ , we can estimate  $\sum_{t=1}^T \mathbb{E}[\Gamma_t]$ . Specifically, we have: 

1564  
1565 
$$\sum_{t=1}^{T} \mathbb{E}[\Gamma_t] = \mathbb{E}\left[\frac{2^{2\delta}\zeta^{\delta}}{\alpha_1} \sum_{i=1}^{d} \sum_{t=1}^{T} \frac{g_{t,i}^2}{S_{t,i}^{1+\frac{\delta}{2}}}\right] \le \mathbb{E}\left[\frac{2^{2\delta}\zeta^{\delta}}{\alpha_1} \sum_{i=1}^{d} \int_{S_{0,i}}^{S_{T,i}} \frac{1}{x^{1+\frac{\delta}{2}}} dx\right]$$

 $\leq \begin{cases} \frac{2^{2\delta}d}{\alpha_1 v^{\frac{\delta}{2}}} \mathbb{E}\left[\zeta^{\delta}\right], & \text{ if } \delta \in (0,1] \\ \frac{2^{2\delta}}{\alpha_1} \mathbb{E}\left[\ln\left(\frac{S_T}{dv}\right)\right], & \text{ if } \delta = 0 \end{cases}$  $\stackrel{(a)}{\leq} \begin{cases} \mathcal{O}(1), & \text{if } \delta \in (0,1] \\ \frac{2^{2\delta}}{\alpha_1} \mathbb{E}\left[\ln\left(\frac{S_T}{dv}\right)\right], & \text{if } \delta = 0 \end{cases}.$ 

In step (a), we used the following *Hölder's* inequality to obtain the  $\mathcal{O}(1)$  result 

$$\mathbb{E}\left[\zeta^{\delta}\right] \leq (\mathbb{E}\left[\zeta\right])^{\delta} = C_{\zeta}^{\delta}.$$

Next, we substitute the estimate of  $\sum_{t=1}^{T} \mathbb{E}[\Gamma_t]$  from Eq. 33 back into Eq. 31, and we obtain the result. This completes the proof. 

(33)

#### D.9 THE PROOF OF LEMMA B.9

*Proof.* We only analyze the case where  $\delta = 0$ ; the case where  $\delta > 0$  can be treated using exactly the same analytical approach. Returning to the approximate descent inequality (Lemma 4.1), we have: 

$$\Pi_{\Delta,t+1}\hat{f}(u_{t+1}) - \Pi_{\Delta,t}\hat{f}(u_t) \le -\frac{5}{16}\Pi_{\Delta,t+1}\sum_{i=1}^d \zeta_i(t) + (L_f+1)\Pi_{\Delta,t+1}\sum_{i=1}^d \eta_{v_t,i}^2 g_{t,i}^2 + \Pi_{\Delta,t+1}M_t.$$

We divide both sides of the above inequality by  $\ln^2(t+1)$ , and noting that  $\ln^2(t+1) < \ln^2(t+2)$ , we obtain: 

$$\frac{\Pi_{\Delta,t+1}\hat{f}(u_{t+1})}{\ln^{2}(t+2)} - \frac{\Pi_{\Delta,t}\hat{f}(u_{t})}{\ln^{2}(t+1)} \leq (L_{f}+1)\frac{\Pi_{\Delta,t+1}}{\ln^{2}(t+1)}\sum_{i=1}^{d}\eta_{v_{t},i}^{2}g_{t,i}^{2} + \frac{\Pi_{\Delta,t+1}M_{t}}{\ln^{2}(t+1)}$$

$$\frac{\Pi_{\Delta,t+1}\leq 1}{\leq}(L_{f}+1)\underbrace{\frac{1}{\ln^{2}(t+1)}\sum_{i=1}^{d}\eta_{v_{t},i}^{2}g_{t,i}^{2}}_{\Omega_{t}} + \frac{\Pi_{\Delta,t+1}M_{t}}{\ln^{2}(t+1)}.$$
(34)

For  $\Omega_t$ , we can perform the following transformation, and we have: 

$$\Omega_t = \frac{1}{\ln^2(t+1)} \sum_{i=1}^d \eta_{v_t,i}^2 g_{t,i}^2 = \frac{1}{\ln^2(t+1)} \sum_{i=1}^d \frac{g_{t,i}^2}{tv_{t,i}} \le \frac{1}{\alpha_1} \sum_{i=1}^d \frac{g_{t,i}^2}{\ln^2(t+1)S_{t,i}}$$

$$\stackrel{(a)}{\leq} \frac{1}{\alpha_1} \sum_{i=1}^{a} \frac{(\zeta'^2) \frac{2g_{t,i}}{v}}{\ln^2(\frac{2S_{t,i}}{v}) \frac{2S_{t,i}}{v}}.$$

In step (a), we use the last result from Lemma B.6, which states:

$$\ln\left(\frac{2S_{t,i}}{v}\right) \le \ln\left(\frac{2S_t}{v}\right) \le \zeta' \ln(T+1).$$

Then, using the series-integral inequality, we can bound  $\sum_{t=1}^{+\infty} \mathbb{E}[\Omega_t]$ , and we obtain: 

$$\sum_{t=1}^{+\infty} \mathbb{E}[\Omega_t] \le \frac{1}{\alpha_1} \sum_{i=1}^d \mathbb{E}\left[\sum_{t=1}^{+\infty} \frac{(\zeta'^2) \frac{2g_{t,i}^2}{v}}{\ln^2(\frac{2S_{t,i}}{v}) \frac{2S_{t,i}}{v}}\right] < \frac{1}{\alpha_1} \sum_{i=1}^d \mathbb{E}\left[\int_2^{+\infty} \frac{\zeta'^2}{x \ln^2 x} dx\right] = \frac{d \mathbb{E}[\zeta'^2]}{\alpha_1 \ln 2}$$

Next, we use the explicit expression for  $\zeta'$  given in Lemma B.6 to bound  $\mathbb{E}[\zeta'^2]$ . We have: 

$$\mathbb{E}[\zeta'^2] = \mathbb{E}\left[8\ln\left(\max\left\{e, \frac{\sqrt{2}\Pi_{\Delta,\infty}^{-1}\zeta}{\sqrt{v}}\right\}\right)\right] = 8\mathbb{E}\left[\max\left\{1, \ln\left(\frac{\sqrt{2}\Pi_{\Delta,\infty}^{-1}\zeta}{\sqrt{v}}\right)\right\}\right] \\ \leq 8\mathbb{E}\left[\max\left\{1, \ln\sqrt{2} - \ln\Pi_{\Delta,\infty} + \ln(\zeta) - \ln\sqrt{v}\right\}\right]$$

Lemma B.3 and Lemma B.6 
$$< +\infty$$
.

1622 As a result, we have:

1629 1630 1631

1634 1635 1636

1620

1621

$$\sum_{t=1}^{+\infty} \mathbb{E}[\Omega_t] < +\infty.$$
(35)

According to the *Lebesgue's Monotone Convergence* theorem, we know that the above result implies:

$$\sum_{t=1}^{+\infty} \mathbb{E}[\Omega_t | \mathscr{F}_{t-1}] < +\infty \text{ a.s.}$$
(36)

1632 Next, we take the conditional expectation with respect to  $\mathscr{F}_{t-1}$  on both sides of Eq. 34, and we 1633 obtain:

$$\mathbb{E}\left[\frac{\Pi_{\Delta,t+1}\hat{f}(u_{t+1})}{\ln^2(t+2)}\middle|\mathscr{F}_{t-1}\right] \le \frac{\Pi_{\Delta,t}\hat{f}(u_t)}{\ln^2(t+1)} + (L_f+1)\mathbb{E}[\Omega_t|\mathscr{F}_{t-1}] + 0.$$
(37)

1637 Based on the result from Eq. 36 and the *supermartingale convergence* theorem, we deduce that 1638  $\frac{\prod_{\Delta,t} \hat{f}(u_t)}{\ln^2(t+1)}$  convergence almost surely. Then, according to Property 5, we can bound  $f(w_t) - f^*$ 1640 using  $\hat{f}(u_t)$ , i.e.,

1641 1642

1643 1644

1645

1649 1650

$$f(w_t) - f^* \le (L_f + 1)(f(u_t) - f^*) + \frac{(L_f + 1)\beta_1^2}{2(1 - \beta_1)^2} \|\eta_{v_{t-1}} \circ m_{t-1}\|^2 + L_f f^*$$
  
$$\le \max\left\{L_f + 1, \frac{(L_f + 1)\beta_1^2}{2C_2(1 - \beta_1)}\right\} \hat{f}(u_t) + L_f f^*.$$

Then we can acquire our first result. Next, we take the expectation on both sides of Eq. 34, and we obtain: 1648  $\begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$ 

$$\mathbb{E}\left[\frac{\Pi_{\Delta,t+1}\hat{f}(u_{t+1})}{\ln^2(t+2)}\right] \le \mathbb{E}\left[\frac{\Pi_{\Delta,t}\hat{f}(u_t)}{\ln^2(t+1)}\right] + (L_f+1)\mathbb{E}[\Omega_t] + 0.$$
(38)

1651 1652 Based on the convergence result of the expectation summation in Eq. 35 and a simple summation 1653 formula for a recursive sequence, we obtain our second result. Thus, the case for  $\delta = 0$  has been 1654 fully analyzed. For the case where  $\delta > 0$ , we can reach the conclusion using the same method. This 1655

1658 1659

1660

1661 1662 1663

1666 1667

#### D.10 THE PROOF OF LEMMA B.10

*Proof.* Since the case of  $\delta > 0$  is relatively straightforward, we first analyze the scenario where  $\delta > 0$ . According to the second conclusion for  $\delta > 0$  in Lemma B.9, we easily obtain:

$$\mathbb{E}[S_T^{3/4}] = \mathbb{E}[\Pi_{\Delta,T+1}^{-3/4} \Pi_{\Delta,T+1}^{3/4} S_T^{3/4}] \stackrel{\text{Hölder's inequality}}{\leq} \mathbb{E}^{1/4}[\Pi_{\Delta,T+1}^{-3}] \mathbb{E}^{3/4}[\Pi_{\Delta,T+1} S_T].$$

1664 Then according to Lemma B.3, we have  $\mathbb{E}[\Pi_{\Delta,T+1}^{-3}] \leq C_{v,d,3}$ . For the other term,  $\mathbb{E}[\Pi_{\Delta,T+1}S_T]$ , 1665 we can handle it as follows:

$$\mathbb{E}[\Pi_{\Delta,T+1}S_T] \le S_0 + \mathbb{E}\left[\Pi_{\Delta,T+1}\sum_{t=1}^T \|g_t\|^2\right] \le dv + \mathbb{E}\left[\sum_{t=1}^T \Pi_{\Delta,t+1} \|g_t\|^2\right]$$

1669  
1670 
$$= dv + \sum_{t=1}^{T} \mathbb{E} \left[ \Pi_{\Delta,t+1} \mathbb{E}[\|g_t\|^2 | \mathscr{F}_{t-1}] \right]$$

1673 
$$\stackrel{\text{Property 1}}{\leq} dv + \sum_{t=1}^{I} \mathbb{E} \left[ \Pi_{\Delta,t+1} ((A+2L_f B)(f(w_t) - f^*) + C) \right]$$

$$\stackrel{\text{Lemma B.9}}{\leq} dv + ((A + 2L_f B)M_{\delta} + C)T.$$

1676 This implies that:

$$\mathbb{E}[\sqrt{S_T}] \le C_{v,d,3}^{1/4} (dv + ((A + 2L_f B)M_\delta + C)T^{3/4} = \mathcal{O}(T^{3/4})$$

1679 For 1680 spc

For the case where 
$$\delta = 0$$
, we use the same approach as in the case of  $\delta > 0$  and apply the corresponding conclusion for  $\delta = 0$  from Lemma B.9. Thus, we obtain:

$$\mathbb{E}[S_T^{3/4}] = \mathcal{O}(T^{3/4} \ln^{3/2} T).$$

1685 D.11 THE PROOF OF LEMMA B.11 

*Proof.* We discuss two cases based on the value of  $\lambda$ . In the first case, when  $\lambda = 1$ , we naturally have:

$$v_{t+1} = \left(1 - \frac{1}{t+1}\right)v_t + \frac{1}{t+1}g_t^{\circ 2} \ (\forall t \ge 1),$$

that is

 $(t+1)v_{t+1} = tv_t + g_t^{\circ 2}.$ 

1693 Summing over all coordinates, we obtain:

$$(t+1)\Sigma_{v_{t+1}} = t\Sigma_{v_t} + ||g_t||^2.$$
(39)

Multiplying both sides of the above equation by  $\Pi_{\Delta,t+1}$ , and noting that  $\Pi_{\Delta,t+1} \ge \Pi_{\Delta,t+2}$ , we obtain:

$$(t+1)\Pi_{\Delta,t+2}\Sigma_{v_{t+1}} = t\Pi_{\Delta,t+1}\Sigma_{v_t} + \Pi_{\Delta,t+1}||g_t||^2.$$

1699 Taking the expectation on both sides, we obtain: 1700

$$\begin{array}{ll} \text{1700} \\ \text{(}t+1) \mathbb{E}[\Pi_{\Delta,t+2}\Sigma_{v_{t+1}}] \leq t \mathbb{E}[\Pi_{\Delta,t+1}\Sigma_{v_t}] + \mathbb{E}[\Pi_{\Delta,t+1}\|g_t\|^2] \\ = t \mathbb{E}[\Pi_{\Delta,t+1}\Sigma_{v_t}] + \mathbb{E}[\Pi_{\Delta,t+1}\mathbb{E}[\|g_t\|^2|\mathscr{F}_{t-1}]] \\ \text{1703} \\ \overset{\text{Property 1}}{\leq} t \mathbb{E}[\Pi_{\Delta,t+1}\Sigma_{v_t}] + (A+2L_fB) \mathbb{E}[\Pi_{\Delta,t+1}(f(w_t)-f^*)] + C \\ \leq t \mathbb{E}[\Pi_{\Delta,t+1}\Sigma_{v_t}] + (A+2L_fB) \Big(\sup_{t\geq 1}\mathbb{E}[\Pi_{\Delta,t+1}(f(w_t)-f^*)]\Big) + C \\ \overset{\text{Lemma B.9}}{\leq} t \mathbb{E}[\Pi_{\Delta,t+1}\Sigma_{v_t}] + (A+2L_fB)M_{\delta} + C. \end{array}$$

By iterating the above inequality, we finally obtain:

$$(t+1)\mathbb{E}[\Pi_{\Delta,t+2}\Sigma_{v_{t+1}}] \leq \begin{cases} \mathbb{E}[\Pi_{\Delta,2}\Sigma_{v_1}] + \left((A+2L_fB)M_{\delta}+C\right)t, & \text{if } \delta \in (0,1]\\ \mathbb{E}[\Pi_{\Delta,2}\Sigma_{v_1}] + \left((A+2L_fB)M_0\ln^2 t + C\right)t, & \text{if } \delta = 0 \end{cases}.$$

This implies that for any  $t \ge 1$ , we always have:

$$\mathbb{E}[\Pi_{\Delta,t+2}\Sigma_{v_{t+1}}] \leq \begin{cases} \mathbb{E}[\Pi_{\Delta,2}\Sigma_{v_1}] + (A+2L_fB)M_\delta + C, & \text{if } \delta \in (0,1] \\ \mathbb{E}[\Pi_{\Delta,2}\Sigma_{v_1}] + (A+2L_fB)M_0\ln^2 t + C, & \text{if } \delta = 0 \end{cases}$$

that is

$$\sup_{t\geq 1} \mathbb{E}[\Pi_{\Delta,t+1}\Sigma_{v_t}] < \begin{cases} \mathbb{E}[\Pi_{\Delta,2}\Sigma_{v_1}] + (A+2L_fB)M_\delta + C, & \text{if } \delta \in (0,1] \\ \mathbb{E}[\Pi_{\Delta,2}\Sigma_{v_1}] + (A+2L_fB)M_0 \ln^2 t + C, & \text{if } \delta = 0 \end{cases}$$

<sup>1721</sup> Next, we discuss the scenario when  $\lambda > 1$ . In this case, we have the following inequality:

1722  
1723  
1724 
$$\Sigma_{v_{t+1}} \leq \Sigma_{v_t} + \frac{1}{(t+1)^{\lambda}} \|g_t\|^2.$$

1725 We multiply both sides of the above inequality by  $\Pi_{\Delta,t+1}$  and, noting its monotonicity, we obtain:

$$\Pi_{\Delta,t+2} \Sigma_{v_{t+1}} \le \Pi_{\Delta,t+1} \Sigma_{v_t} + \frac{\Pi_{\Delta,t+1}}{(t+1)^{\lambda}} \|g_t\|^2.$$
(40)

Taking the conditional expectation with respect to  $\mathscr{F}_{t-1}$  on both sides of the inequality, we have:

$$\mathbb{E}[\Pi_{\Delta,t+2}\Sigma_{v_{t+1}}|\mathscr{F}_{t-1}] \leq \Pi_{\Delta,t+1}\Sigma_{v_t} + \frac{\Pi_{\Delta,t+1}}{(t+1)^{\lambda}}\mathbb{E}[\Pi_{\Delta,t+1}||g_t||^2|\mathscr{F}_{t-1}]$$

According to Property 1 and Lemma B.9, we easily obtain:

$$\begin{split} &\sum_{t=1}^{+\infty} \frac{\Pi_{\Delta,t+1}}{(t+1)^{\lambda}} \,\mathbb{E}[\Pi_{\Delta,t+1} \|g_t\|^2 |\mathscr{F}_{t-1}] \\ &\leq \begin{cases} \left( (A+2L_f B) \sup_{t\geq 1} \left( \Pi_{\Delta,t+1}(f(w_t)-f^*) \right) + C \right) \cdot \sum_{t=1}^{+\infty} \frac{1}{(t+1)^{\lambda}}, & \text{if } \delta \in (0,1] \\ \left( (A+2L_f B) \sup_{t\geq 1} \left( \frac{\Pi_{\Delta,t+1}(f(w_t)-f^*)}{\ln^2 t} \right) + C \right) \cdot \sum_{t=1}^{+\infty} \frac{\ln^2 t}{(t+1)^{\lambda}}, & \text{if } \delta = 0 \\ &< +\infty \text{ a.s.} \end{split}$$

By the supermartingale convergence theorem, we easily obtain that  $\Pi_{\Delta,t+1}\Sigma_{v_t}$  converges almost surely, which implies that  $\sup_{t\geq 1}\Pi_{\Delta,t+1}\Sigma_{v_t} < +\infty$  a.s. According to Lemma B.3, where  $\sup_{t\geq 1}\Pi_{\Delta,t+1}^{-1} < +\infty$  a.s., we can immediately deduce that  $\sup_{t\geq 1}\Sigma_{v_t} < +\infty$  a.s.. Next, we prove that the expected supremum is finite. Taking the expectation on both sides of Eq. 40, we obtain:

$$\mathbb{E}\left[\Pi_{\Delta,t+2}\Sigma_{v_{t+1}}\right] \leq \mathbb{E}\left[\Pi_{\Delta,t+1}\Sigma_{v_t}\right] + \frac{1}{(t+1)^{\lambda}} \mathbb{E}\left[\Pi_{\Delta,t+1} \|g_t\|^2\right].$$

By summing the above recursive inequalities and using the results from Property 1 and Lemma B.9, we can easily prove that

$$\sup_{t \ge 1} \mathbb{E}[\Pi_{\Delta, t+1} \Sigma_{v_t}] < \begin{cases} \left( (A + 2L_f B) M_{\delta} + C \right) \sum_{t=1}^{+\infty} \frac{1}{(t+1)^{\lambda}}, & \text{if } \delta \in (0, 1] \\ \left( (A + 2L_f B) M_0 + C \right) \sum_{t=1}^{+\infty} \frac{\ln^2 t}{(t+1)^{\lambda}}, & \text{if } \delta = 0 \end{cases} < +\infty.$$

1754 With this, we complete the proof.

#### 1757 D.12 THE PROOF OF LEMMA B.12

*Proof.* According to the result from Lemma B.7, it is straightforward to see that when  $\delta > 0$ , we have:

$$\sum_{t=2}^{T} \mathbb{E}\left[\Pi_{\Delta,t+1} \frac{\eta_{t-1} \|\nabla f(w_t)\|^2}{\sqrt{\Sigma_{v_{t-1}} + \mu}}\right] \le \sum_{t=2}^{T} \mathbb{E}\left[\Pi_{\Delta,t+1} \frac{\eta_{t-1}}{\sqrt{\Sigma_{v_{t-1}} + \mu}} \sum_{i=1}^{d} (\nabla_i f(w_t))^2\right]$$
$$\le \sum_{t=1}^{T} \mathbb{E}\left[\Pi_{\Delta,t+1} \sum_{i=1}^{d} \zeta_i(t)\right]$$

 $< C_{4,\delta} < +\infty.$ 

1769 Next, we apply the *Lebesgue's Monotone Convergence* theorem:

$$\sum_{t=2}^{+\infty} \Pi_{\Delta,t+1} \frac{\eta_{t-1} \|\nabla f(w_t)\|^2}{\sqrt{\Sigma_{v_{t-1}} + \mu}} < +\infty \ \text{ a.s}$$

Then, by combining the almost surely boundedness of  $\sup_{t\geq 1} \prod_{\Delta,t+1}^{-1}$  and  $\sup_{t\geq 1} \sum_{v_t}$  from Lemma B.3 and Lemma B.11, we immediately obtain:

$$\begin{aligned} & \overset{1776}{1777} & \sum_{t=1}^{+\infty} \eta_t \|\nabla f(w_t)\|^2 \le \|\nabla f(w_1)\|^2 + \sum_{t=2}^{+\infty} \eta_{t-1} \|\nabla f(w_t)\|^2 \\ & \overset{1778}{1779} & < \|\nabla f(w_1)\|^2 + \left(\sup_{t\ge 1} \Pi_{\Delta,t+1}^{-3/2}\right) \cdot \left(\sqrt{\sup_{t\ge 1} \Pi_{\Delta,t+1} \Sigma_{v_t}} + \mu\right) \cdot \sum_{t=2}^{+\infty} \Pi_{\Delta,t+1} \frac{\eta_{t-1} \|\nabla f(w_t)\|^2}{\sqrt{\Sigma_{v_{t-1}}} + \mu} \\ & \overset{1781}{1781} \end{aligned}$$

 $< +\infty$  a.s..

According to the L-smooth assumption (Assumption 2.2), it is easy to see that 

1784  
1785 
$$|\|\nabla f(w_t)\| - \|\nabla f(u_t)\|| \le L_f \|w_t - u_t\| = \frac{L_f \beta_1}{1 - \beta_1} \|\eta_{v_{t-1}} \circ m_{t-1}\| \le \frac{L_f \beta_1}{(1 - \beta_1)t^{\frac{1}{2} + \delta} \mu} \|m_{t-1}\|,$$
1786

 $\|\nabla f(u_t)\|^2 \le \left(\|\nabla f(w_t)\| + \frac{L_f \beta_1}{(1-\beta_1)t^{\frac{1}{2}+\delta}\mu} \|m_{t-1}\|\right)^2$ 

 $\leq 2 \left( \|\nabla f(w_t)\|^2 + \frac{L_f^2 \beta_1^2}{(1-\beta_1)^2 t^{1+2\delta} \mu^2} \|m_{t-1}\|^2 \right)$ 

 $\leq 2 \|\nabla f(w_t)\|^2 + \frac{2L_f^2 \beta_1^2}{(1-\beta_1)^2 t^{1+2\delta} \mu^2} \Big(\sup_{t>1} \|m_t\|\Big)^2.$ 

which implies

Thid implies that 

$$\sum_{t=1}^{+\infty} \eta_t \|\nabla f(u_t)\|^2 \le 2\sum_{t=1}^{+\infty} \eta_t \|\nabla f(w_t)\|^2 + \frac{2L_f^2 \beta_1^2}{(1-\beta_1)^2 \mu^2} \Big(\sup_{t\ge 1} \|m_t\|\Big)^2 \sum_{t=1}^{+\infty} \frac{1}{t^{1+2\delta}}$$
Lemma B.9 and  $\delta > 0$ 

$$\le +\infty \quad \text{a.s.}$$

With this, we complete the proof.

#### Ε **PROOFS OF THEOREMS**

THE PROOF OF THEOREM 3.1 E.1

*Proof.* According to Lemma B.7, we have: 

$$\sum_{t=1}^{T} \mathbb{E} \left[ \Pi_{\Delta,t+1} \sum_{i=1}^{d} \zeta_i(t) \right] \leq \begin{cases} C_{4,\delta}, & \text{if } \delta \in (0,1] \\ C_5 + C_6 \mathbb{E} \left[ \ln(S_T) \right], & \text{if } \delta = 0 \end{cases}.$$

According to the monotonicity of  $\eta_{v_t,i}$  in Property 2 and the monotonicity of  $\Pi_{\Delta,t}$  itself, we obtain the following inequality: 

$$\sum_{t=1}^{T} \mathbb{E} \left[ \Pi_{\Delta,T+1} \frac{\|\nabla f(w_t)\|^2}{T^{\frac{1}{2}+\delta}(\sqrt{v_T}+\mu)} \right] \le \sum_{t=1}^{T} \mathbb{E} \left[ \Pi_{\Delta,t+1} \sum_{i=1}^{d} \zeta_i(t) \right] \le \begin{cases} C_{4,\delta}, & \text{if } \delta \in (0,1] \\ C_5 + C_6 \mathbb{E} \left[ \ln(S_T) \right], & \text{if } \delta = 0 \end{cases}.$$

For the leftmost part of the above inequality, we apply the *Cauchy-Schwarz* inequality and obtain:

$$\mathbb{E}\left[\Pi_{\Delta,T+1}^{-1}T^{\frac{1}{2}+\delta}(\sqrt{v_T}+\mu)\right]\left(\sum_{t=1}^{T}\mathbb{E}\left[\Pi_{\Delta,T+1}\frac{\|\nabla f(w_t)\|^2}{T^{\frac{1}{2}+\delta}(\sqrt{v_T}+\mu)}\right]\right) \ge \sum_{t=1}^{T}\mathbb{E}[\|\nabla f(w_t)\|],$$

which means 

$$\sum_{t=1}^{1826} \mathbb{E}[\|\nabla f(w_t)\|] \leq \begin{cases} C_{4,\delta} \mathbb{E}\left[\Pi_{\Delta,T+1}^{-1} T^{\frac{1}{2}+\delta}(\sqrt{v_T}+\mu)\right], & \text{if } \delta \in (0,1] \\ C_5 \mathbb{E}\left[\Pi_{\Delta,T+1}^{-1} T^{\frac{1}{2}+\delta}(\sqrt{v_T}+\mu)\right] + C_6 \mathbb{E}\left[\Pi_{\Delta,T+1}^{-1} T^{\frac{1}{2}+\delta}(\sqrt{v_T}+\mu)\right] \mathbb{E}[\ln(S_T)], & \text{if } \delta = 0 \end{cases}$$

$$\text{Combining the results from Lemma B 10, Lemma B 11, and Lemma B 3, we obtain:}$$

Combining the results from Lemma B.10, Lemma B.11 and Lemma B.3, we obtain:

$$\begin{split} \mathbb{E}\left[\Pi_{\Delta,T+1}^{-1}T^{\frac{1}{2}+\delta}(\sqrt{v_{T}}+\mu)\right] &\leq 2T^{\frac{1}{2}+\delta}\sqrt{\mathbb{E}[\Pi_{\Delta,T+1}^{-3}]}\sqrt{\mathbb{E}[\Pi_{\Delta,T+1}(v_{T}+\mu^{2})]} \\ &\leq \begin{cases} C_{v,d,3}^{1/2}\mathcal{O}(T^{\frac{1}{2}+\delta}), & \text{if } \gamma > 1\\ C_{v,d,3}^{1/2}\mathcal{O}(T^{\frac{1}{2}+\delta}), & \text{if } \gamma = 1, \ \delta \in (0,1] \ .\\ C_{v,d,3}^{1/2}\mathcal{O}(\sqrt{T}\ln T) & \text{if } \gamma = 1, \ \delta = 0 \end{cases} \end{split}$$

1836 and 1837  $\mathbb{E}[\ln(S_T)] = \frac{4}{3} \mathbb{E}[\ln(S_T^{3/4})] \le \frac{4}{3} \ln(\mathbb{E}[S_T^{3/4}]) = \begin{cases} \mathcal{O}(\ln T), & \text{if } \delta \in (0,1] \\ \mathcal{O}(\ln T) + \mathcal{O}(\ln \ln T), & \text{if } \delta = 0 \end{cases}.$ 1838 1840 Combining the two estimates above, we finally obtain: 1841 1842  $\sum_{t=1}^{T} \mathbb{E}[\|\nabla f(w_t)\|] \le \begin{cases} \mathcal{O}(T^{\frac{1}{2}+\delta}), & \text{if } \delta \in (0,1]\\ \mathcal{O}(\sqrt{T}\ln T), & \text{if } \gamma > 1, \ \delta = 0\\ \mathcal{O}(\sqrt{T}\ln^2 T), & \text{if } \gamma = 1, \ \delta = 0 \end{cases}$ 1843 1844 1845 that is 1847  $\frac{1}{T}\sum_{t=1}^{T}\mathbb{E}[\|\nabla f(w_t)\|] \le \begin{cases} \mathcal{O}\left(\frac{1}{T^{\frac{1}{2}-\delta}}\right), & \text{if } \delta \in (0,1]\\ \mathcal{O}\left(\frac{\ln T}{\sqrt{T}}\right), & \text{if } \gamma > 1, \ \delta = 0 \\ \mathcal{O}\left(\frac{\ln^2 T}{\sqrt{T}}\right), & \text{if } \gamma = 1, \ \delta = 0 \end{cases}$ 1849 1850 1851 With this, we complete the proof. E.2 THE PROOF OF LEMMA B.8 1855 *Proof.* From Theorem 3.1, we have 1857  $\lim_{T \to +\infty} \min_{1 \le t \le T} \mathbb{E} \left\| \nabla f(w_t) \right\| \le \lim_{T \to +\infty} \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left\| \nabla f(w_t) \right\| = 0.$ 1858 1859 1860 This implies that there exists a subsequence  $\{w_{d_t}\}_{t\geq 1}$  of  $\{w_t\}_{t\geq 1}$  such that 1861 1862  $\lim_{t \to +\infty} \mathbb{E} \left\| \nabla f(w_{d_t}) \right\| = 0.$ 1863 1864 By the *Riesz* theorem, we can find a further subsequence  $\{w_{c_t}\}_{t>1}$  from  $\{w_{d_t}\}_{t>1}$  such that 1865  $\lim_{t \to +\infty} \|\nabla f(w_{c_t})\| = 0 \quad \text{a.s.}$ 1866 1867 This completes the proof. 1868 E.3 THE PROOF OF THEOREM 3.2 1870 1871 Proof. According to the L-smooth assumption (Assumption 2.2), it is easy to see that 1872 1873  $||\nabla f(w_t)|| - ||\nabla f(u_t)||| \le L_f ||w_t - u_t|| = \frac{L_f \beta_1}{1 - \beta_1} ||\eta_{v_{t-1}} \circ m_{t-1}|| \le \frac{L_f \beta_1}{(1 - \beta_1)t^{\frac{1}{2} + \delta_H}} ||m_{t-1}||.$ 1874 1875 According to Lemma B.9, we know that  $\sup_{t\geq 1} ||m_{t-1}|| < +\infty$  a.s.. This implies that 1876 1877  $\lim_{t \to +\infty} |\|\nabla f(w_t)\| - \|\nabla f(u_t)\|| \le \left(\sup_{t > 1} \|m_{t-1}\|\right) \cdot \lim_{t \to +\infty} \frac{L_f \beta_1}{(1 - \beta_t) t^{\frac{1}{2} + \delta_H}} = 0 \text{ a.s.}$ 1878 1879 1880 This implies that we only need to prove  $\lim_{t\to+\infty} \|\nabla f(u_t)\| = 0$  a.s.. To achieve this objective, 1881 we proceed as follows. 1882 For any l > 0, we construct the following stopping time <sup>3</sup> sequence  $\{\tau_{l,n}\}_{n \ge 1}$ : 1883  $\tau_{l,1} := \min\{t \ge 1 : \|\nabla f(u_t)\| > l\}, \ \tau_{l,2} := \min\{t > \tau_{l,1} : \|\nabla f(u_t)\| \le l\},$ 1884 1885  $\tau_{l,2k-1} := \min\{t > \tau_{l,2k-2} : \|\nabla f(u_t)\| > l\}, \ \tau_{l,2k} := \min\{t > \tau_{l,2k-1} : \|\nabla f(u_t)\| \le l\}.$ 

<sup>&</sup>lt;sup>3</sup>In this paper, we adopt the following definition of stopping time: Let  $\tau$  be a random variable defined on the filtered probability space  $(\Omega, \mathscr{F}, (\mathscr{F}_n)_{n \in \mathbb{N}}, \mathbb{P})$  with values in  $\mathbb{N} \cup \{+\infty\}$ . Then  $\tau$  is called a stopping time (with respect to the filtration  $(\mathscr{F}_n)_{n \in \mathbb{N}}$ ) if the following condition holds:  $\{\tau = n\} \in \mathscr{F}_n$  for all n.

According to the subsequence convergence result in Lemma B.8, we know that when  $\tau_{2k-1} < +\infty$  ( $\forall k \ge 1$ ), it must hold that  $\tau_{2k} < +\infty$  a.s.. We now discuss two cases:

1893 1. When there exists some  $k_0 \ge 1$  such that  $\tau_{2k_0-1} = +\infty$ , this implies that eventually  $\{\|\nabla f(u_t)\|\}_{t\ge 1}$  will remain below l, i.e.,

$$\limsup_{t \to +\infty} \|\nabla f(u_t)\| < l.$$
(41)

1898 2. Next, we focus on the second case, where for all  $\tau_{2k-1}$ , we have  $\tau_{2k-1} < +\infty$ . In this situation, 1899 we examine the behavior of  $\sup_{\tau_{2k-1} \le t < \tau_{2k}} \|\nabla f(u_t)\|$ . It is easy to see that:

$$\begin{split} \sup_{\tau_{2k-1} \le t < \tau_{2k}} \|\nabla f(u_t)\| \le l + \sup_{\tau_{2k-1} \le t < \tau_{2k}} \|\nabla f(u_t)\| - \|\nabla f(u_{\tau_{2k-1}-1})\| \\ \le l + \left( \sum_{t=\tau_{2k-1}-1}^{\tau_{2k}-1} \left\| \|\nabla f(u_t)\| - \|\nabla f(u_{t-1})\| \right| \right) \\ \overset{\text{l-smooth}}{\le} l + \left( L_f \sum_{t=\tau_{2k-1}-1}^{\tau_{2k}-1} \left\| u_t - u_{t-1} \right\| \right) \\ \overset{\text{Eq. 4}}{\le} l + L_f \underbrace{\left( \sum_{t=\tau_{2k-1}-1}^{\tau_{2k}-1} \left\| \eta_{v_t} \circ g_t \right\| \right)}_{\Upsilon_{k,1}} \\ + \frac{\beta_1 L_f^2}{1 - \beta_1} \underbrace{\left( \sum_{t=\tau_{2k-1}-1}^{\tau_{2k}-1} \left\| \Delta_t \circ m_{t-1} \right\| \right)}_{\Upsilon_{k,2}}. \end{split}$$

Our next goal is to prove separately that  $\limsup_{k\to+\infty} \Upsilon_{k,1} = 0$  a.s. and  $\limsup_{k\to+\infty} \Upsilon_{k,2} = 0$  a.s.. For  $\Upsilon_{k,1}$ , we have:

$$\begin{split} & \begin{array}{l} 1920\\ 1921\\ 1922\\ 1922\\ 1922\\ 1923\\ 1924\\ 1925\\ 1926\\ 1926\\ 1926\\ 1926\\ 1926\\ 1927\\ 1926\\ 1927\\ 1926\\ 1927\\ 1928\\ 1926\\ 1927\\ 1928\\ 1926\\ 1927\\ 1928\\ 1926\\ 1927\\ 1928\\ 1930\\ 1931\\ 1931\\ 1931\\ 1931\\ 1934\\ 1935\\ 1936\\ 1936\\ 1936\\ 1936\\ 1937\\ 1936\\ 1936\\ 1937\\ 1936\\ 1936\\ 1937\\ 1936\\ 1936\\ 1937\\ 1936\\ 1$$

For  $\Upsilon_{k,1,1}$ , we have:

$$\Upsilon_{k,1,1} \stackrel{\text{Property 1}}{\leq} \left( \sum_{t=\tau_{2k-1}-1}^{\tau_{2k}-1} \eta_t \left( (A+2L_f B)(f(w_t)-f^*)+C \right) \right)$$

$$\leq \left( (A + 2L_f B) \sup_{t \geq 1} (f(w_t) - f^*) + C \right) \cdot \left( \sum_{t = \tau_{2k-1} - 1}^{\tau_{2k} - 1} \eta_t \right)$$

1944  
1945 
$$= \left( (A + 2L_f B) \sup_{t \ge 1} (f(w_t) - f^*) + C \right) \cdot \left( \eta_{\tau_{2k-1}} + \left( \sum_{t = \tau_{2k-1} - 1}^{\tau_{2k}} \eta_t \right) \right)$$
  
1946

$$\stackrel{(a)}{\leq} \frac{1}{l^2} \left( (A + 2L_f B) \sup_{t \ge 1} (f(w_t) - f^*) + C \right) \cdot \left( \eta_{\tau_{2k-1}} + \left( \sum_{t = \tau_{2k-1}}^{\tau_{2k}-1} \eta_t \| \nabla f(u_t) \|^2 \right) \right).$$

In step (a), this is due to the fact that, over the interval  $[\tau_{2k-1}, \tau_{2k})$ , we always have  $\|\nabla f(u_t)\|^2 > l^2$ . Based on Lemma B.12, we know that

$$\sum_{t=1}^{+\infty} \eta_t \|\nabla f(u_t)\|^2 < +\infty \quad \text{a.s.}$$

<sup>1956</sup> By applying the *Cauchy's* convergence principle, we can prove that

$$\lim_{k \to +\infty} \sum_{t=\tau_{2k-1}}^{\tau_{2k}-1} \eta_t \|\nabla f(u_t)\|^2 = 0 \quad \text{a.s.}$$

1961 On the other hand, it is evident that  $\lim_{k \to +\infty} \eta_{\tau_{2k-1}} = 0$ . Meanwhile, based on Lemma B.9 and 1962 Lemma B.3, we can easily prove that

$$\sup_{t \ge 1} (f(w_t) - f^*) \le \left( \sup_{t \ge 1} \Pi_{\Delta, t+1}^{-1} \right) \cdot \left( \sup_{t \ge 1} \left( \Pi_{\Delta, t+1} (f(w_t) - f^*) \right) \right) < +\infty \quad \text{a.s.}$$

1966 Therefore, we have proven that

$$\limsup_{k \to +\infty} \Upsilon_{k,1,1} = \lim_{k \to +\infty} \Upsilon_{k,1,1} = 0$$

For  $\Upsilon_{k,1,2}$ , we consider the following martingale difference sequence:

$$\overline{X}_T := \sum_{t=1}^T \sum_{i=1}^d \eta_t (|g_{t,i}| - \mathbb{E}[|g_{t,i}| \mid \mathscr{F}_{t-1}]).$$

We can compute

$$\begin{array}{ll} 1976 & \sum_{t=1}^{+\infty} \mathbb{E} \left[ \left( \sum_{i=1}^{d} \eta_t (|g_{t,i}| - \mathbb{E}[|g_{t,i}| \mid \mathscr{F}_{t-1}]) \right)^2 \middle| \mathscr{F}_{t-1} \right] \leq d \sum_{t=1}^{+\infty} \eta_t^2 \sum_{i=1}^{d} \mathbb{E}[(|g_{t,i}| - \mathbb{E}[|g_{t,i}| \mid \mathscr{F}_{t-1}])^2] \\ 1979 & \leq d \sum_{t=1}^{+\infty} \eta_t^2 \sum_{i=1}^{d} \mathbb{E}[||g_t||^2 \mid \mathscr{F}_{t-1}] \\ 1980 & \leq d \sum_{t=1}^{+\infty} \eta_t^2 \sum_{i=1}^{d} \mathbb{E}[||g_t||^2 \mid \mathscr{F}_{t-1}] \\ 1983 & \leq d \sum_{t=1}^{+\infty} \eta_t^2 \sum_{i=1}^{d} \left( (A + 2L_f B) \sup_{t \geq 1} (f(w_t) - f^*) + C \right) \\ 1984 & \leq d \sum_{i=1}^{d} \left( (A + 2L_f B) \left( \sup_{t \geq 1} \Pi_{\Delta,t} \right) \left( \sup_{t \geq 1} (f(w_t) - f^*) + C \right) \right) \cdot \sum_{t=1}^{+\infty} \eta_t^2 \\ 1987 & = M \\ 1988 & \leq M \\ 1989 & = M \\ 1987 \end{array}$$

1990 By the *Martingale Convergence* theorem, we obtain

$$\lim_{T \to +\infty} \overline{X}_T = \sum_{t=1}^{+\infty} \sum_{i=1}^d \eta_t (|g_{t,i}| - \mathbb{E}[|g_{t,i}| \mid \mathscr{F}_{t-1}]) < +\infty \text{ a.s.}$$

Using the *Cauchy's Convergence* principle, we can easily prove that

1996  
1997 
$$\limsup_{k \to +\infty} \Upsilon_{k,1,2} = \lim_{k \to +\infty} \sum_{t=\tau_{2k-1}-1}^{\tau_{2k}-1} \sum_{i=1}^{d} \eta_t (|g_{t,i}| - \mathbb{E}[|g_{t,i}||\mathscr{F}_{t-1}]) = 0 \quad \text{a.s.}$$

Combining the above two limit proofs for  $\Upsilon_{t,1,1}$  and  $\Upsilon_{t,1,2}$ , we can conclude that

$$\limsup_{k \to +\infty} \Upsilon_{t,1} = 0 \quad \text{a.s}$$

Next, we will demonstrate that  $\lim_{k\to+\infty} \Upsilon_{k,2} = 0$  a.s.. This is relatively easy; we only need to examine the series sum

$$\sum_{t=1}^{+\infty} \|\Delta_t \circ m_{t-1}\|,$$

2007 which can easily be proven to satisfy

$$\sum_{t=1}^{2008} \|\Delta_t \circ m_{t-1}\| = \sum_{t=1}^{+\infty} \sum_{i=1}^d \Delta_{t,i}^{-1} |m_{t-1,i}| < \sqrt{d} \Big( \sup_{t \ge 1} \Pi_{\Delta,t} \Big) \cdot \Big( \sup_{t \ge 1} \Pi_{\Delta,t} \|m_{t-1}\| \Big) \cdot \sum_{i=1}^d \Delta_{t,i} \overset{\text{Lemma B.3 and B.9}}{<} + \infty \quad \text{a.s.}$$

Thus, by the *Cauchy's Convergence* principle, we obtain

2012 2013 2014

2015

2018 2019

2022

2023 2024

2026

2027

2029

2000 2001

2004 2005 2006

$$\limsup_{k \to +\infty} \Upsilon_{k,2} = \lim_{k \to +\infty} \sum_{t=\tau_{2k-1}-1}^{\tau_{2k}-1} \sum_{i=1}^{d} \eta_t (|g_{t,i}| - \mathbb{E}[|g_{t,i}||\mathscr{F}_{t-1}]) = 0 \quad \text{a.s.}$$

2016 Combining the limit results for  $\Upsilon_{k,1}$  and  $\Upsilon_{k,2}$ , we conclude that

$$\limsup_{k \to +\infty} \sup_{\tau_{2k-1} \le t < \tau_{2k}} \|\nabla f(u_t)\| \le l+0 = l.$$

Moreover, combining  $\sup_{\tau_{2k} \le t < \tau_{2k+1}} \|\nabla f(u_t)\| < l$ , we can deduce that

$$\limsup_{t \to +\infty} \|\nabla f(u_t)\| \le l \quad \text{a.s}$$

Then, due to the arbitrariness of l, we conclude that 2025

$$\limsup_{t \to +\infty} \|\nabla f(u_t)\| = 0 \quad \text{a.s.}$$

2028 This implies that

$$\lim_{t\to+\infty} \|\nabla f(u_t)\| = 0 \quad \text{a.s.}$$

2030 2031 Thus, we complete the proof.

2032 2033

#### 2034 E.4 THE PROOF OF THEOREM 3.3 2035

**Proof.** Since we have already proven almost sure convergence in Theorem 3.2, it is natural to attempt to prove  $L_1$  convergence via the *Lebesgue's Dominated Convergence* theorem. To achieve this, we need to find a function h that is  $\mathscr{F}_{\infty}$ -measurable and satisfies  $\mathbb{E} |h| < +\infty$ , and such that for all  $t \ge 1$ , we have  $\|\nabla f(w_t)\| \le |h|$ . Since for all t, we naturally have  $\|\nabla f(w_t)\| \le \sup_{k\ge 1} \|\nabla f(w_k)\|$ , we only need to prove that  $\mathbb{E}[\sup_{k\ge 1} \|\nabla f(w_k)\|] < +\infty$ . We proceed to achieve this goal.

Returning to the *Approximate Descent Inequality* (Lemma 4.1), we have:

$$\Pi_{\Delta,t+1}\hat{f}(u_{t+1}) - \Pi_{\Delta,t}\hat{f}(u_t) \le -\frac{5}{16}\Pi_{\Delta,t+1}\sum_{i=1}^d \zeta_i(t) + (L_f+1)\Pi_{\Delta,t+1}\sum_{i=1}^d \eta_{v_t,i}^2 g_{t,i}^2 + \Pi_{\Delta,t+1}M_t.$$
(42)

2045 2046

2050

2043 2044

For any  $\lambda > 0$ , define the stopping time  $\tau_{\lambda}$  as the first time the sequence  $\{\Pi_{\Delta,t}\hat{f}(u_t)\}_{t\geq 1}$  exceeds  $\lambda$ , i.e.,

$$\tau_{\lambda} := \min\{t \ge 2 : \Pi_{\Delta, t} f(u_t) > \lambda\}.$$

It can be rigorously verified that  $\tau_{\lambda}$  is a stopping time with respect to the filtration  $\{\mathscr{F}_t\}_{t\geq 1}$ , and satisfies a special property  $[\tau_{\lambda} = n] \in \mathscr{F}_{n-1}$  for all  $n \geq 1$ . This implies that the preceding time

 $\tau_{\lambda} - 1$  is also a stopping time. Next, for any deterministic time  $T \ge 3$ , we define  $\tau_{\lambda,T} := \tau_{\lambda} \wedge T$ . 2053 We then sum the indices of Eq. 42 from 1 to  $\tau_{\lambda,T} - 1$ . Specifically, we have:

$$\Pi_{\Delta,\tau_{\lambda,T}}\hat{f}(u_{\tau_{\lambda,T}}) \le \Pi_{\Delta,1}\hat{f}(u_1) + (L_f+1)\sum_{t=1}^{\tau_{\lambda,T}-1}\Pi_{\Delta,t+1}\sum_{i=1}^d \eta_{v_t,i}^2 g_{t,i}^2 + \sum_{t=1}^{\tau_{\lambda,T}-1}\Pi_{\Delta,t+1}M_t.$$

Taking the expectation on both sides, we obtain:

$$\mathbb{E}\left[\Pi_{\Delta,\tau_{\lambda,T}}\hat{f}(u_{\tau_{\lambda,T}})\right] \leq \mathbb{E}\left[\Pi_{\Delta,1}\hat{f}(u_{1})\right] + (L_{f}+1)\mathbb{E}\left[\sum_{t=1}^{\tau_{\lambda,T}-1}\Pi_{\Delta,t+1}\sum_{i=1}^{d}\eta_{v_{t},i}^{2}g_{t,i}^{2}\right] \\ + \mathbb{E}\left[\sum_{t=1}^{\tau_{\lambda,T}-1}\Pi_{\Delta,t+1}M_{t}\right].$$

$$+ \mathbb{E}\left[\sum_{t=1}^{N,1} \Pi_{\Delta,t+1}M_t\right].$$

Since  $\{\prod_{\Delta,t+1}M_t, \mathscr{F}_{t-1}\}_{t\geq 1}$  is a martingale difference sequence and  $\tau_{\lambda,T} \leq T < +\infty$ , by *Doob's Stopped* theorem, we know that:

$$\mathbb{E}\left[\sum_{t=1}^{\tau_{\lambda,T}-1} \Pi_{\Delta,t+1} M_t\right] = 0.$$

2071 This implies that:

$$\mathbb{E}\left[\Pi_{\Delta,\tau_{\lambda,T}}\hat{f}(u_{\tau_{\lambda,T}})\right] \leq \mathbb{E}\left[\Pi_{\Delta,1}\hat{f}(u_{1})\right] + (L_{f}+1)\mathbb{E}\left[\sum_{t=1}^{\tau_{\lambda,T}-1}\Pi_{\Delta,t+1}\sum_{i=1}^{d}\eta_{v_{t,i}}^{2}g_{t,i}^{2}\right]$$
$$< \mathbb{E}\left[\Pi_{\Delta,1}\hat{f}(u_{1})\right] + (L_{f}+1)\mathbb{E}\left[\sum_{t=1}^{T}\Pi_{\Delta,t+1}\sum_{i=1}^{d}\eta_{v_{t,i}}^{2}g_{t,i}^{2}\right].$$

According to the estimates of  $\sum_{t=1}^{T} \mathbb{E}[\Gamma_t]$  from Eq. 32 to Eq. 33, we easily obtain that when  $\delta > 1$ , the following holds:

$$\mathbb{E}\left[\Pi_{\Delta,\tau_{\lambda,T}}\hat{f}(u_{\tau_{\lambda,T}})\right] \leq \overline{M} < +\infty,$$

2082 where

$$\overline{M} := \mathbb{E}\left[\Pi_{\Delta,1}\hat{f}(u_1)\right] + \frac{2^{2\delta}d(L_f+1)}{\alpha_1 v^{\frac{\delta}{2}}}C_{\zeta}^{\delta}$$

2085 On the other hand, we easily observe the following event decomposition:

$$\left[\sup_{2 \le t < T} \prod_{\Delta, t} \hat{f}(u_t) > \lambda\right] = \bigcup_{k=2}^{T-1} [\tau_\lambda = k] = \bigcup_{k=2}^{T-1} [\tau_{\lambda, T} = k].$$

Moreover, since for any  $j \neq k$ , we have  $[\tau_{\lambda,T} = j] \cap [\tau_{\lambda,T} = k] = \emptyset$ , it follows that:

$$\mathbb{P}\left[\sup_{2\leq t< T}\Pi_{\Delta,t}\hat{f}(u_t) > \lambda\right] = \sum_{k=2}^{T-1} \mathbb{P}[\tau_{\lambda,T} = k] \overset{\text{Markov's inequality}}{\leq} \frac{1}{\lambda} \sum_{k=2}^{T-1} \mathbb{E}\left[\Pi_{\Delta,k}\hat{f}(u_k)\mathbb{I}_{[\tau_{\lambda,T}=k]}\right] \\ < \frac{1}{\lambda} \mathbb{E}\left[\Pi_{\Delta,\tau_{\lambda,T}}\hat{f}(u_{\tau_{\lambda,T}})\right] \leq \frac{\overline{M}}{\lambda}.$$
(43)

Next, for any  $K \ge 1$ , we compute  $\mathbb{E}\left[\left(\sup_{2\le t < T} \prod_{\Delta, t} \hat{f}(u_t)\right)^{3/4} \land K\right]$ . We have

$$\mathbb{E}\left[\left(\sup_{2\leq t< T}\Pi_{\Delta,t}\hat{f}(u_{t})\right)^{3/4}\wedge K\right] = -\int_{0}^{+\infty} x \, d\left(\mathbb{P}\left[\left(\sup_{2\leq t< T}\Pi_{\Delta,t}\hat{f}(u_{t})\right)^{3/4}\wedge K > x\right]\right)\right)$$

$$= -\int_{0}^{+\infty}\left(\int_{0}^{x}1d\lambda\right)d\left(\mathbb{P}\left[\left(\sup_{2\leq t< T}\Pi_{\Delta,t}\hat{f}(u_{t})\right)^{3/4}\wedge K > x\right]\right)$$

$$= -\int_{0}^{+\infty}\left(\int_{0}^{+\infty}1d\lambda\right)d\left(\mathbb{P}\left[\left(\sup_{2\leq t< T}\Pi_{\Delta,t}\hat{f}(u_{t})\right)^{3/4}\wedge K > x\right]\right)\right)$$

$$Fubini's theorem = -\int_{0}^{+\infty}\left(\int_{\lambda}^{+\infty}1d\left(\mathbb{P}\left[\left(\sup_{2\leq t< T}\Pi_{\Delta,t}\hat{f}(u_{t})\right)^{3/4}\wedge K > x\right]\right)\right)d\lambda$$

$$\begin{aligned} &= \int_{0}^{+\infty} \mathbb{P}\left[\left(\sup_{2 \le t < T} \Pi_{\Delta,t} \hat{f}(u_t)\right)^{3/4} \wedge K > \lambda\right] d\lambda \\ &= \int_{0}^{+\infty} \mathbb{P}\left[\left(\sup_{2 \le t < T} \Pi_{\Delta,t} \hat{f}(u_t)\right)^{3/4} \wedge K > \lambda\right] d\lambda \\ &\leq 1 + \int_{1}^{+\infty} \mathbb{P}\left[\left(\sup_{2 \le t < T} \Pi_{\Delta,t} \hat{f}(u_t)\right) \wedge K^{4/3} > \lambda^{4/3}\right] d\lambda \\ &= 1 + \int_{1}^{+\infty} \mathbb{P}\left[\left(\sup_{2 \le t < T} \Pi_{\Delta,t} \hat{f}(u_t)\right) \wedge K^{4/3} > \lambda^{4/3}\right] d\lambda \\ &< 1 + \int_{1}^{+\infty} \mathbb{P}\left[\left(\sup_{2 \le t < T} \Pi_{\Delta,t} \hat{f}(u_t)\right) > \lambda^{4/3}\right] d\lambda \\ &\leq 1 + \int_{1}^{+\infty} \mathbb{P}\left[\left(\sup_{2 \le t < T} \Pi_{\Delta,t} \hat{f}(u_t)\right) > \lambda^{4/3}\right] d\lambda \\ &= 1 + 3\overline{M}. \end{aligned}$$

2121 Next, we take  $K \to +\infty$  and apply the *Lebesgue's Monotone Convergence* theorem. We get:

$$\mathbb{E}\left[\left(\sup_{2\leq t< T} \Pi_{\Delta,t} \hat{f}(u_t)\right)^{3/4}\right] \leq 1 + 3\overline{M}.$$

2125 Next, by taking  $T \to +\infty$  and applying the *Lebesgue's Monotone Convergence* theorem once again, 2126 we obtain:

$$\mathbb{E}\left[\left(\sup_{t\geq 2}\Pi_{\Delta,t}\hat{f}(u_t)\right)^{3/4}\right] \le 1+3\overline{M}$$

Note that for any finite t, we have  $\Pi_{\Delta,t+1} \ge \Pi_{\Delta,\infty}$  (where  $\Pi_{\Delta,\infty}$  is defined in Lemma B.3). Thus, we have:

$$\mathbb{E}\left[\Pi_{\Delta,\infty}^{3/4}\left(\sup_{t\geq 2}\hat{f}(u_t)\right)^{3/4}\right] \leq \mathbb{E}\left[\left(\sup_{t\geq 2}\Pi_{\Delta,t}\hat{f}(u_t)\right)^{3/4}\right] \leq 1+3\overline{M}.$$

2133 Next, by applying *Hölder's* inequality, we obtain: 2134

$$\mathbb{E}\left[\left(\sup_{t\geq 2}\hat{f}(u_t)\right)^{1/2}\right] \le \mathbb{E}^{1/3}\left[\Pi_{\Delta,\infty}^{-3/2}\right] \mathbb{E}^{2/3}\left[\Pi_{\Delta,\infty}^{3/4}\left(\sup_{t\geq 2}\hat{f}(u_t)\right)^{3/4}\right] \stackrel{\text{Lemma B.3}}{\le} C_{v,d,3/2}^{1/3} (1+3\overline{M})^{2/3}.$$

Then, according to Property 5, we can bound  $f(w_t) - f^*$  using  $\hat{f}(u_t)$ , i.e.,

$$f(w_t) - f^* \le (L_f + 1)(f(u_t) - f^*) + \frac{(L_f + 1)\beta_1^2}{2(1 - \beta_1)^2} \|\eta_{v_{t-1}} \circ m_{t-1}\|^2 + L_f f^*$$
  
$$\le \max\left\{L_f + 1, \frac{(L_f + 1)\beta_1^2}{2C_2(1 - \beta_1)}\right\} \hat{f}(u_t) + L_f f^*.$$

That means

2122 2123 2124

2127 2128

2131 2132

2135 2136

2145 2146 2147

$$\mathbb{E}\left[\left(\sup_{t\geq 2} \left(f(w_t) - f^*\right)\right)^{1/2}\right] \le \left(\max\left\{L_f + 1, \frac{(L_f + 1)\beta_1^2}{2C_2(1 - \beta_1)}\right\}\right)^{1/2} C_{v,d,3/2}^{1/3} (1 + 3\overline{M})^{2/3} + \sqrt{L_f|f^*|}.$$

Finally, according to Lemma B.1, we obtain:

$$\mathbb{E}\left[\sup_{t\geq 2} \|\nabla f(w_t)\|\right] \leq \sqrt{2L_f} \mathbb{E}\left[\left(\sup_{t\geq 2} \left(f(w_t) - f^*\right)\right)^{1/2}\right] \\ < \sqrt{2L_f} \left(\max\left\{L_f + 1, \frac{(L_f + 1)\beta_1^2}{2C_2(1 - \beta_1)}\right\}\right)^{1/2} C_{v,d,3/2}^{1/3} (1 + 3\overline{M})^{2/3} \\ + \sqrt{2L_f^2}|f^*|.$$

2156

2157 2158 2159 By adding the first term, we obtain:

$$\mathbb{E}\left[\sup_{t\geq 1} \|\nabla f(w_t)\|\right] < \|\nabla f(w_1)\| + \sqrt{2L_f} \left(\max\left\{L_f + 1, \frac{(L_f + 1)\beta_1^2}{2C_2(1 - \beta_1)}\right\}\right)^{1/2} C_{v,d,3/2}^{1/3} (1 + 3\overline{M})^{2/3}$$

2160 
$$+\sqrt{2L_f^2|f^*|} < +\infty.$$

Finally, combining the almost sure convergence result from Theorem 3.2 with the *Lebesgue's Dominated Convergence* theorem, we obtain the  $L_1$  convergence result, namely:

$$\lim_{t \to +\infty} \mathbb{E}[\|\nabla f(w_t)\|] = 0.$$

With this, we complete the proof.

2165	
2166	
2167	
2168	
2169	
2170	
2171	
2172	
2173	
2174	
2175	
2176	
2177	
2178	
2179	
2180	
2181	
2182	
2183	
2184	
2185	
2186	
2187	
2188	
2189	
2190	
2191	
2192	
2193	
2194	
2195	
2196	
2197	
2198	
2199	
2200	
2201	
2202	
2203	
2204	
2205	
2206	
2207	
2208	
2209	
2210	
2211	
2212	
2213	