
Convergence of Distributed Adaptive Optimization with Local Updates

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We study distributed adaptive algorithms with local updates (intermittent commu-
2 nication). Despite the great empirical success of adaptive methods in distributed
3 training of modern machine learning models, the theoretical benefits of local up-
4 dates within adaptive methods, particularly in terms of reducing communication
5 complexity, have not been fully understood yet. In this paper, we prove that *Lo-*
6 *cal SGD* with momentum (*Local* SGDM) and *Local* Adam can outperform their
7 minibatch counterparts in convex and weakly convex settings, respectively. Our
8 analysis relies on a novel technique to prove contraction during local iterations,
9 which is a crucial yet challenging step to show the advantages of local updates,
10 under generalized smoothness assumption and gradient clipping strategy.

11 1 Introduction

12 Leveraging parallelism is crucial in accelerating the training of modern machine learning models
13 for large scale optimization problems. In distributed environments such as large data-centers or in
14 the federated learning setting, where the devices working together are spread apart, communication
15 between the distributed workers is a key bottleneck. In this work, we consider solving

$$\min_{x \in \mathbb{R}^d} f(x) := \mathbb{E}_{\xi \sim \mathcal{D}} [F(x; \xi)]. \quad (1.1)$$

16 in a distributed setting with M workers. Each worker has access to f via the stochastic gradient
17 oracle $\nabla F(x; \xi)$, where ξ is independently drawn from the distribution \mathcal{D} . In federated learning, this
18 is known as the *homogeneous* setting, since all workers draw from the same data distribution.

19 Perhaps the simplest algorithm for distributed optimization is distributed *minibatch stochastic gradient*
20 *descent* (*SGD*), in which at each iteration, each worker computes a minibatch of gradients, and a
21 gradient step is taken by averaging the gradient computed among the M workers. However, such an
22 algorithm requires communicating at each gradient step, which may be expensive. Thus numerous
23 works have proposed distributed algorithms with less frequent communication. A popular and well-
24 studied algorithm is *Local* SGD, also known as FedAvg [McMahan et al., 2017], where each worker
25 runs SGD independently and periodically synchronizes with others by averaging the iterates.

26 Despite the success of *Local* SGD in federated learning [McMahan et al., 2017], it may not exhibit
27 good performance when training Transformer-based large language models (LLMs). Many empirical
28 studies suggest that adaptive methods (e.g., Adam [Kingma and Ba, 2014]) are much better suited for
29 natural language processing than vanilla SGD [Goodfellow et al., 2016, Zhang et al., 2020, Kunstner
30 et al., 2023, Pan and Li, 2023]. Furthermore, as shown in Zhang et al. [2019, 2020], language models
31 tend to have unbounded global smoothness and heavy-tailed noise, which may also contribute to
32 the worse performance of SGD. Parallelizing adaptive methods requires an even more expensive
33 communication cost since additional terms, such as the momentum or the Adam denominator, need

34 to be synchronized. Previous works on distributed adaptive optimization have utilized compression
 35 and quantization techniques to address this issue [Bernstein et al., 2018, Wangni et al., 2018, Wang
 36 et al., 2023]. While Douillard et al. [2023] has shown the great empirical success of *Local* Adam, to
 37 the best of our knowledge, there are no theoretical results trying to improve training efficiency or
 38 adaptive methods from the perspective of intermittent communication.

39 In this paper, we investigate **distributed adaptive optimization algorithms in the homogeneous**
 40 **regime**, in order to establish theoretical guarantees for the benefits of local iterations in reducing
 41 communication complexity. We focus on the convex or weakly convex setting, because in the
 42 non-convex setting, without non-standard strong smoothness assumptions, we are not aware of any
 43 theoretical-proven advantages of local iterations, even for non-adaptive methods¹. Further, in the case
 44 of Adam, we consider the weakly convex setting (as opposed to the standard convex setting), since
 45 we are not aware of any results on the convergence rate of Adam which take advantage of convexity.
 46 To handle unbounded global smoothness and heavy-tailed noise, we use the coordinate-wise gradient
 47 clipping mechanism.

48 We propose a distributed version of Adam, namely, *Local* Adam, with gradient clipping. Our
 49 algorithm also reduces to *Local* SGD with momentum (*Local* SGDM), with some specific hyper-
 50 parameter choices.

- 51 • In Theorem 2, we establish the first convergence guarantee for *Local* SGDM in the convex
 52 setting, which outperforms the convergence rate of *Minibatch* SGDM. The rate we obtain is
 53 in line with the rate of *Local* SGD [Woodworth et al., 2020a].
- 54 • In Theorem 3, we establish a convergence rate for *Local* Adam in the weakly convex setting.
 55 We show that *Local* Adam can provably improve communication efficiency compared to its
 56 minibatch baseline.

57 For the first time, we are able to show the benefits of local iterations for the two commonly used
 58 algorithms, SGDM and Adam. This suggests that we may be able to improve the training efficiency
 59 of LLMs by using intermittent communication.

60 Additionally, our results hold under generalized smoothness and heavy-tailed noise. Our result is the
 61 first high probability bound for distributed optimization algorithms with local updates, to the best of
 62 our knowledge. The conventional in-expectation rate seems fail to capture some important properties
 63 like heavy/light tailed noise distribution. The high probability convergence guarantee can sometimes
 64 be more informative and useful in practice [Gorbunov et al., 2020].

65 As for technical contribution, we use a **novel technique to prove contraction for adaptive methods**,
 66 which bounds the consensus error between the iterates at different workers. This is a key step in
 67 proving benefits of local updates. Different from *Local* SGD, our update direction involves momentum
 68 or even distorted momentum due to the denominator in *Local* Adam, making it challenging to
 69 disentangle these accumulated stochastic gradients. To address this issue, we define and analyze an
 70 auxiliary sequence which is conditionally independent of the latest stochastic gradient and thus can
 71 construct a martingale. We will introduce the technique in more details in Section 4.

72 **Notation** Let $\|\cdot\|$ be the standard Euclidean norm of a vector or the spectral norm of a matrix. For
 73 any $x, y \in \mathbb{R}^d$, the expressions $x + y, x \odot y, \frac{x}{y}$ stand for coordinate-wise sum, product and division,
 74 respectively. And $x \preceq y$ means each coordinate of $x - y$ is no greater than 0. Furthermore, we use
 75 $x^2, \sqrt{x}, |x|$ to denote the coordinate-wise square, square root and absolute value. We use $\mathbb{E}_m[X_m]$
 76 to denote the average $\frac{1}{M} \sum_{m=1}^M X_m$. The coordinate-wise clipping operator $\mathbf{clip}(\cdot, \rho) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is
 77 defined as $[\mathbf{clip}(X, \rho)]_i = \text{sgn}([X]_i) \cdot \min\{|X_i|, \rho\}$. We use $[N]$ to denote the set $\{1, 2, \dots, N\}$. For
 78 a subset $\Omega_0 \subset \mathbb{R}^d$, let $\mathbf{conv}(\cdot)$ denote the convex hull of Ω_0 and $\mathbf{B}_{R_0}(\Omega_0)$ denote the neighborhood

¹Under the stronger assumptions of 3rd-order smoothness [Glasgow et al., 2022] and mean smoothness [Patel et al., 2022], there are demonstrated advantages of local iterations in the non-convex setting. While our theoretical results are for the convex or weakly convex setting, it is likely that local iterations are advantageous in practice for non-convex objectives, just in the same way *Local* SGD has been shown to be advantageous in practice for non-convex objectives [McMahan et al., 2017].

79 of Ω_0 with radius R_0 . Finally, we use standard $\mathcal{O}(\cdot)$, $\Omega(\cdot)$, $\Theta(\cdot)$ to omit constant factors and $\tilde{\mathcal{O}}(\cdot)$ to
80 omit logarithmic factors.

81 2 Problem Setup

82 Consider the distributed optimization problem

$$\min_{x \in \mathbb{R}^d} f(x) := \mathbb{E}_{\xi \sim \mathcal{D}} [F(x; \xi)]. \quad (2.1)$$

83 Here \mathcal{D} is the data distribution and f is the population loss function. We consider a setting with M
84 parallel workers, and a budget of R total communication rounds, and T total gradient computations at
85 each worker. We will describe the implementation of the *local* and *minibatch* versions of a centralized
86 algorithm \mathcal{A} , which uses a single stochastic gradient in each iteration, as illustrated in Figure 1.

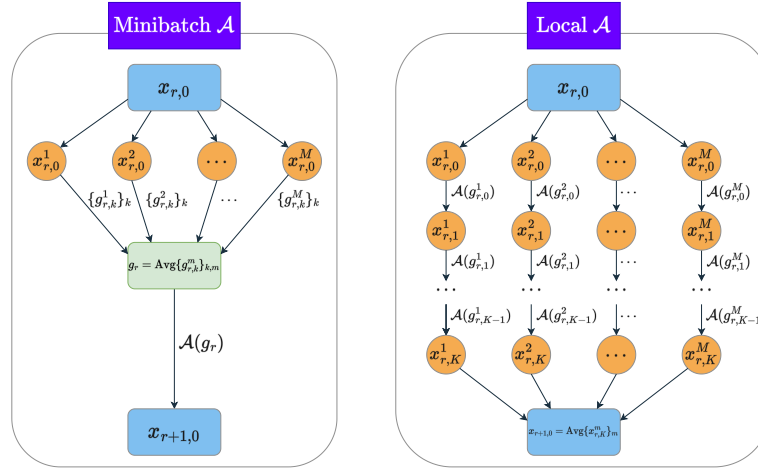


Figure 1: *Minibatch A* v.s. *Local A* in one communication round. Minibatch version computes the average of all KM gradients and then executes one step of \mathcal{A} , while local version runs \mathcal{A} independently for K steps at each worker.

87 In the *local* version of algorithm \mathcal{A} , in each round r of the R total communication rounds, each
88 worker m independently executes $K = T/R$ steps of local updates (according to the algorithm \mathcal{A}).
89 For a worker m , we denote the k th gradient computed in round r by $g_{r,k}^m$. Then the M workers
90 synchronize the iterates and related momentum state. We use *Minibatch A* to denote a distributed
91 implementation of \mathcal{A} run for R rounds, where KM stochastic gradients are computed and averaged
92 at each step. This is a fair baseline to compare the local update algorithms to, since the number of
93 gradient calls and communication rounds are the same.

94 *Local Adam* is shown in Algorithm 1, which is a natural extension of centralized Adam [Kingma and
95 Ba, 2014]. The stochastic gradient is clipped by a coordinate-wise clipping operator with threshold
96 ρ . After K steps of local updates, all the workers average their current iterates x_t^m , their first order
97 momentum u_t^m , and their second order momentum v_t^m . These averaged quantities become the values
98 used at the beginning of the next local round. Note that there are two slight differences from original
99 Adam. First, we do not involve bias correction here, i.e., u_t^m and v_t^m are not divided by $1 - \beta_1^t$ or
100 $1 - \beta_2^t$, respectively. Second, λ in the denominator is in the square root, while it is outside of the
101 denominator in original Adam. These modifications do not harm the spirit of Adam and are made for
102 the convenience of analysis.

103 2.1 Assumptions

104 Throughout this work, we will use the following assumptions.

105 **Assumption 1** (Lower-boundedness). f is closed, twice continuously differentiable and $\inf_{x \in \mathbb{R}^d} f(x) =$
106 $f(x_*) = f_* > -\infty$.

Algorithm 1 *Local Adam*

Require: initial model x_0 , learning rate η , momentum $\beta_1, \beta_2 \in [0, 1)$
Set $x_{0,0}^m = x_0$, $u_{0,-1}^m = 0$, $v_0 = 0$ for each worker $m \in [M]$
for $r = 0, \dots, R-1$ **do**
 for each worker $m \in [M]$ **in parallel do**
 for $k = 0, \dots, K-1$ **do**
 $g_{r,k}^m = \nabla F(x_{r,k}^m; \xi_{r,k}^m)$, $\widehat{g}_{r,k}^m = \text{clip}(g_{r,k}^m, \rho)$ ▷ Compute clipped stochastic gradient
 $u_{r,k}^m = \beta_1 u_{r,k-1}^m + (1 - \beta_1) \widehat{g}_{r,k}^m$ ▷ Update 1st-order momentum
 $v_{r,k}^m = \beta_2 v_{r,k-1}^m + (1 - \beta_2) \widehat{g}_{r,k}^m \odot \widehat{g}_{r,k}^m$ ▷ Update 2nd-order momentum
 $x_{r,k+1}^m = x_{r,k}^m - \frac{\eta}{\sqrt{v_{r,k}^m + \lambda^2}} \odot u_{r,k}^m$ ▷ Update model
 end for
 end for
 $x_{r+1,0}^m = \mathbb{E}_m[x_{r,K}^m]$, $u_{r+1,-1}^m = \mathbb{E}_m[u_{r,K-1}^m]$, $v_{r+1,-1}^m = v_{r+1} := \mathbb{E}_m[v_{r,K-1}^m]$
 ▷ Communicate and average
end for

107 **Assumption 2** (Smoothness). *There exists some set $\Omega \subset \mathbb{R}^d$ and $L > 0$, such that for any $x, y \in \Omega$,*

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|, \quad (2.2)$$

108

$$\|\nabla f(x)\|^2 \leq 2L(f(x) - f_*). \quad (2.3)$$

109 Similar to [Sadiev et al. \[2023\]](#), we only requires some properties of f on a subset Ω of \mathbb{R}^d , since we
110 can prove that all the iterates will not leave this subset with high probability. In contrast, the typical
111 smoothness assumption requires (2.2) on the entire domain.

112 There are many works [[Zhang et al., 2019](#), [Crawshaw et al., 2022](#), [Faw et al., 2023](#), [Wang et al.,](#)
113 [2022](#), [Li et al., 2024b](#)] that make weaker smoothness assumptions (typically called “generalized
114 smoothness”), most of which are in the form of (L_0, L_1) -smoothness:

$$\|\nabla^2 f(x)\| \leq L_0 + L_1 \|\nabla f(x)\|, \quad \forall x \in \mathbb{R}^d. \quad (2.4)$$

115 [Li et al. \[2024a\]](#) considers an extension called ℓ -smoothness, which replaces the linear function of
116 $\|\nabla f\|$ in the right hand side of (2.4) with a sub-quadratic function $\ell(\cdot)$. As pointed out in [Li et al.](#)
117 [\[2024a, Corollary 3.6\]](#), all of these will induce Assumption 2 if Ω is some level-set of the objective
118 function². Therefore, we directly use this more general assumption to get cleaner results.

119 **Assumption 3** (Bounded α -moment noise). *There exists some set $\Omega \subset \mathbb{R}^d$, $\alpha \geq 4$ and constant*
120 *vector $\sigma \succeq 0$ such that for any $x \in \Omega$,*

$$\mathbb{E}_{\xi \sim \mathcal{D}} |\nabla F(x; \xi) - \nabla f(x)|^\alpha \preceq \sigma^\alpha. \quad (2.5)$$

121 *Let $\sigma_\infty := \|\sigma\|_\infty = \max_i \{\sigma_i\}$, $\sigma := \|\sigma\| = (\sigma_1^2 + \dots + \sigma_d^2)^{1/2}$.*

122 **Remark 1.** *To get a high probability bound under generalized smoothness, the assumption on*
123 *stochastic noise is crucial. Light-tailed noise with bounded exponential moment (e.g., bounded,*
124 *sub-exponential, sub-gaussian) are considered in [Harvey et al. \[2019\]](#), [Li and Orabona \[2020\]](#), [Li](#)
125 *et al. [2024b]*. *There are also attempts for heavy-tailed noise with finite α -moment [[Gorbunov et al.,](#)*
126 *2020, [Cutkosky and Mehta, 2021](#), [Faw et al., 2023](#)]. In the most literature studying heavy-tailed*
127 *noise, they restrict to the case where $1 < \alpha \leq 2$. However, in the matter of getting a logarithmic*
128 *dependence on $1/\delta$, where δ is the confidence level, the essence lies in whether we assume bounded*
129 *exponential moment or just polynomial moment (see [Appendix E](#) for detailed discussions). For*
130 *technical convenience, we only consider $\alpha \geq 4$ in this paper, but our analysis methods can be easily*
131 *extended to the case where $\alpha < 4$.**

²e.g., if $\Omega \subset \{x : f(x) - f_* \leq \Delta\}$, then (L_0, L_1) -smoothness would imply Assumption 2 for $L \asymp L_0 + L_1 \Delta$.
Note that we may not obtain the optimal dependence on L_0, L_1 in this way though.

132 **Remark 2** (Noise of minibatch). *It follows from Petrov [1992] that if the gradient is estimated by a*
 133 *batch of i.i.d samples with batch size N , the α -moment of noise has upper bound of:*

$$\mathbb{E}_{\{\xi_i\} \stackrel{i.i.d.}{\sim} \mathcal{D}} \left| \frac{1}{N} \sum_{i=1}^N \nabla F(x; \xi_i) - \nabla f(x) \right|^\alpha \preceq c(\alpha) (\sigma / \sqrt{N})^\alpha, \quad (2.6)$$

134 *where $c(\alpha)$ is a problem-independent constant. It is easy to see that this bound is tight when the noise*
 135 *is Gaussian. Therefore, to get the rate for batch size N , we can just simply replace σ with σ / \sqrt{N}*
 136 *(up to a constant depending on α) in the original convergence guarantee for batch size 1.*

137 3 Main Results

138 In this section, we provide our main results for *Local Adam* and its simplified version: *Local SGDM*.
 139 For the first time, we will be able to show the benefits of local iterations for the two algorithms,
 140 compared with their minibatch baselines in certain regime of M, K, R .

141 3.1 Local SGDM

142 Before getting into *Local Adam*, we start with a simpler yet also important algorithm: *Local SGD*
 143 *with momentum. Note that when $\beta_2 = 1, \lambda = 1$, Algorithm 1 will reduce to *Local SGDM*. We restate*
 144 *the complete version of *Local SGDM* in Algorithm 2 in Appendix C.*

145 **Assumption 4** (Convexity). *There exists some set $\Omega \subset \mathbb{R}^d$ and constant $\mu \geq 0$ such that f is*
 146 *μ -strongly convex on Ω , i.e., for any $x, y \in \Omega$,*

$$\langle \nabla f(x) - \nabla f(y), x - y \rangle \geq \mu \|x - y\|^2, \quad (3.1)$$

147

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\mu}{2} \|x - y\|^2. \quad (3.2)$$

148 Let $D_0 := \|x_0 - x_*\|$. Now we state the results for *Local SGDM* below. Notably, our results are the
 149 first convergence guarantee for distributed SGDM with local updates in (strongly) convex setting.

150 **Theorem 1** (Strongly convex, full version see Theorem C.4). *Let Assumption 1, 2, 3, 4 hold for*
 151 *$\Omega := \{\|x - x_*\| \leq \sqrt{3}D_0\}$ and $\mu > 0$. Further assume that $K \gtrsim \log \frac{MKR}{\delta}, 1 - \beta_1 = \Omega(1)$ and*
 152 *$\|\sigma\|_{2\alpha} d^{\frac{1}{2} - \frac{1}{2\alpha}} = \mathcal{O}(\sigma)$. Then with probability no less than $1 - \delta$, *Local SGDM* yields*

$$f(\hat{x}) - f_* \leq \exp\left(-\Theta\left(\frac{\mu KR}{L}\right)\right) + \tilde{\mathcal{O}}\left(\frac{\sigma^2}{\mu MKR} + \frac{L\sigma^2}{\mu^2 KR^2} + \frac{\sigma^2}{\mu} \left(\frac{L^{\frac{1}{2}}}{\mu^{\frac{1}{2}} KR}\right)^{\frac{2(\alpha-1)}{\alpha}}\right). \quad (3.3)$$

153 **Theorem 2** (Convex, full version see Theorem C.5). *Let Assumption 1, 2, 3, 4 hold for $\Omega :=$*
 154 *$\{\|x - x_*\| \leq \sqrt{3}D_0\}$ and $\mu = 0$. Further assume that $K \gtrsim \log \frac{MKR}{\delta}, 1 - \beta_1 = \Omega(1)$ and*
 155 *$\|\sigma\|_{2\alpha} d^{\frac{1}{2} - \frac{1}{2\alpha}} = \mathcal{O}(\sigma)$. Then with probability no less than $1 - \delta$, *Local SGDM* yields*

$$f(\hat{x}) - f_* \leq \tilde{\mathcal{O}}\left(\frac{LD_0^2}{KR} + \frac{\sigma D_0}{\sqrt{MKR}} + \frac{L^{\frac{1}{3}} \sigma^{\frac{2}{3}} D_0^{\frac{4}{3}}}{K^{\frac{1}{3}} R^{\frac{2}{3}}} + D_0 \left(\frac{(LD_0)^{\frac{1}{2}} \sigma^{\frac{\alpha}{\alpha-1}}}{KR}\right)^{\frac{2(\alpha-1)}{3\alpha-1}}\right). \quad (3.4)$$

156 **Remark 3** (Confidence level δ). *δ does not appear in the bound since we have $\log \frac{1}{\delta}$ dependence.*

157 Our method can also be applied to *Minibath SGDM* (by substituting M, K with 1 and σ with $\frac{\sigma}{\sqrt{MK}}$;
 158 see Remark 2), whose convergence guarantee is

$$f(\hat{x}) - f_* \lesssim \begin{cases} \exp\left(-\Theta\left(\frac{\mu R}{L}\right)\right) + \tilde{\mathcal{O}}\left(\frac{\sigma^2}{\mu MKR}\right), & \text{if } \mu > 0, \\ \tilde{\mathcal{O}}\left(\frac{LD_0^2}{R} + \frac{\sigma D_0}{\sqrt{MKR}}\right), & \text{otherwise.} \end{cases} \quad (3.5)$$

159 This rate matches the well-known in-expectation lower bound on the convergence rate of *Minibatch*
160 SGD (up to logarithmic factors). In fact, our analysis improves the state-of-the-art rate for strongly-
161 convex SGDM (given in Liu et al. [2020b]), which has a stochastic term as $\tilde{\mathcal{O}}\left(\frac{L\sigma^2}{\mu^2 MKR}\right)$. In the
162 convex setting, our rate is consistent with the state-of-the-art centralized in-expectation bound of
163 SGDM in Sebbouh et al. [2021]. Further notice that the last term in both (3.3) and (3.4) is due to the
164 bias of gradient clipping and would be negligible as long as $K^{\alpha-2} \gtrsim \frac{\mu R^2}{L}$ or $K^{\frac{3\alpha-5}{2}} \gtrsim \frac{\sigma R^2}{LD_0}$. In
165 this case, our guarantee for *Local* SGDM is aligned with the rate of *Local* SGD in Woodworth et al.
166 [2020a], Khaled et al. [2020] up to logarithmic factor. Therefore, we can see the benefits of local
167 iterations in the large M and large K regime compared to minibatch baseline.
168 We defer the detailed proof to Appendix C.

169 3.2 Local Adam

170 The convergence of Adam is much more difficult to prove. Reddi et al. [2019] pointed out that the
171 original proof in Kingma and Ba [2014] in centralized convex setting was incorrect. Therefore, the
172 convergence of Adam in for convex function is of independent interest and beyond our scope. Instead,
173 we turn to consider Adam in the weakly convex setting.

174 **Assumption 5** (Weak convexity). *There exists constant $\tau > 0$ such that f is τ -weakly convex, i.e.,*
175 *for any $x, y \in \mathbb{R}^d$,*

$$176 \langle \nabla f(x) - \nabla f(y), x - y \rangle \geq -\tau \|x - y\|^2, \quad (3.6)$$

$$176 f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle - \frac{\tau}{2} \|x - y\|^2, \quad \nabla^2 f(x) \succeq -\tau I_d. \quad (3.7)$$

177 Note that L -smoothness implies that Assumption 5 always holds with $\tau = L$. Also note that here we
178 assume the weak convexity holds in \mathbb{R}^d for technical simplicity. Let $H_r = \mathbf{diag}(\sqrt{v_r + \lambda^2}) \succeq \lambda I_d$
179 and $\Delta := f(x_0) - f_*$. Furthermore, inspired by Liu et al. [2020b], define an auxiliary sequence
180 $\{z_{r,k}^m\}$ as:

$$180 z_{r,k+1}^m = \begin{cases} (x_{r,k+1}^m - \beta_1 x_{r,k}^m)/(1 - \beta_1) & \text{if } k \neq K - 1, \\ (x_{r,k+1}^m - \beta_1 \bar{x}_{r,k})/(1 - \beta_1) & \text{otherwise.} \end{cases} \quad (3.8)$$

181 Let $\bar{z}_{r,k} := \mathbb{E}_m[z_{r,k}^m]$. Now we state the main result of *Local* Adam below (see Theorem D.2 for
182 more general results on Moreau envelope).

183 **Theorem 3** (Full version see Theorem D.3). *Let Assumption 1, 2, 3, 5 hold for $\Omega = \mathbf{conv}(\mathbf{B}_{R_0}(\Omega_0))$,*
184 *where $\Omega_0 := \{f(x) - f_* \leq 4\Delta\}$ and $R_0 = \sqrt{\Delta}/(80L)$. Further assume $K \gtrsim \log(MKR/\delta)$,*
185 *$1 - \beta_1 = \Omega(1)$, $\|\sigma\|_{2\alpha} d^{\frac{1}{2} - \frac{1}{2\alpha}} = \mathcal{O}(\sigma)$ and $1 - \beta_2 = \tilde{\mathcal{O}}(K^{-3/2}R^{-1/2})$. Then with probability no*
186 *less than $1 - \delta$, Local Adam yields*

$$186 \frac{\lambda}{KR} \sum_{r=0}^{R-1} \sum_{k=0}^{K-1} \|\nabla f(\bar{z}_{r,k})\|_{H_r^{-1}}^2 \quad (3.9)$$

$$= \tilde{\mathcal{O}}\left(\frac{\tau\Delta}{R} + \frac{L\Delta}{KR} + \sqrt{\frac{L\Delta\sigma^2}{MKR}} + \frac{(L\Delta\sigma)^{\frac{2}{3}}}{K^{\frac{1}{3}}R^{\frac{2}{3}}} + \left(\frac{L\Delta\sigma^{\frac{\alpha-1}{\alpha-1}}}{KR}\right)^{\frac{2(\alpha-1)}{3\alpha-2}}\right).$$

187 The RHS of (3.9) consists of four parts. The first part is $\frac{\tau\Delta}{R} + \frac{L\Delta}{KR}$, which is the optimization term
188 and determined by the upper bound of learning rate η . The second term is $\sqrt{\frac{L\Delta\sigma^2}{MKR}}$, corresponding to
189 the standard statistical lower bound from MKR stochastic gradients [Arjevani et al., 2023]. The third
190 component is $\frac{(L\Delta\sigma)^{\frac{2}{3}}}{K^{\frac{1}{3}}R^{\frac{2}{3}}}$, which is sourced from the discrepancy overhead of doing local iterations.
191 And the last one, $\left(\frac{L\Delta\sigma^{\frac{\alpha-1}{\alpha-1}}}{KR}\right)^{\frac{2(\alpha-1)}{3\alpha-2}}$, is induced by the bias of clipped stochastic gradient and can be
192 dominated when $K^{\frac{3\alpha-4}{2}} \gtrsim \sigma^2 R/(L\Delta)$.

193 Our analysis method can also be applied to *Minibatch* Adam (by substituting M, K with 1 and σ
 194 with $\frac{\sigma}{\sqrt{MK}}$; see Remark 2), and the convergence rate is

$$\tilde{\mathcal{O}}\left(\frac{L\Delta}{R} + \sqrt{\frac{L\Delta\sigma^2}{MKR}}\right), \quad (3.10)$$

195 aligned with (up to logarithmic factor) the state-of-the-art convergence guarantees for smooth weakly
 196 convex functions [Davis and Drusvyatskiy, 2019, Deng and Gao, 2021]. Suppose $K^{\frac{3\alpha-4}{2}} \gtrsim$
 197 $\sigma^2 R/(L\Delta)$ and hence the last term in (3.9) would be dominated and negligible. Now we can
 198 observe the benefits of local iterations. Note that both (3.9) and (3.10) have the statistical lower bound
 199 $1/\sqrt{MKR}$. Hence when the statistical term dominates, both algorithms have similar worst-case rate.
 200 Once we leave the noise-dominated regime, then *Local* Adam converges faster than *Minibatch* Adam
 201 whenever $K \gtrsim \sigma^2 R/(L\Delta)$. And the gap will increase as K grows until $K \asymp L/\tau$.

202 Therefore, we conclude that in the large M and small τ regime, *Local* Adam would outperform
 203 *Minibatch* Adam. Since f is close to convex function when τ is small, this is consistent with
 204 Woodworth et al. [2020a]. Please see Appendix D.5 for more comparisons about Moreau envelop.

205 We defer further discussions on the choices of other important hyper-parameters including $\beta_1, \beta_2, \lambda$
 206 to Appendix D.5. The complete proof is in Appendix D.

207 4 Proof Sketch

208 In this section, we show high-level ideas in our proofs. We only demonstrate the *Local* Adam here
 209 since *Local* SGDM is a special case of *Local* Adam ($\beta_2 = 1$) and has similar patterns.

210 As a common practice in the study of weakly convex function [Davis and Drusvyatskiy, 2019, Mai
 211 and Johansson, 2020], the norm of the gradient of the Moreau envelope can serve as a proxy for
 212 near-stationarity. Here we use a generalized Moreau envelope for adaptive algorithms, proposed by
 213 Alacaoglu et al. [2020]. For any positive definite matrix H and $\gamma > 0$ such that $\gamma^{-1}H \succeq \tau I_d$, define
 214 the Moreau envelope of f as

$$f_\gamma^H(x) := \min_{y \in \mathbb{R}^d} f(y) + \frac{1}{2\gamma} \|x - y\|_H^2. \quad (4.1)$$

215 With some abuse of notation, we define $f_\gamma^\lambda(x) := f_\gamma^{\lambda I_d}(x) = f_{\gamma/\lambda}(x)$. The common convergence
 216 metric for weakly-convex function is correspondingly $\|\nabla f_\gamma^H(\cdot)\|_{H^{-1}}$, which can bound $\|\nabla f(\cdot)\|_{H^{-1}}$,
 217 as shown in the following lemma.

218 **Lemma 4** (Full version see Lemma D.4). *Let $z \in \Omega_0$ and $y := \arg \min_x f(x) + \frac{1}{2\gamma} \|x - z\|_H^2$ for
 219 some $H \succeq \lambda I_d$ and $L/\lambda \geq \gamma^{-1} \geq 2\tau/\lambda$. Then*

$$\nabla f_\gamma^H(z) = \nabla f(y) = H(z - y)/\gamma, \quad \|\nabla f(z)\|_{H^{-1}} \leq 2\gamma L \|\nabla f_\gamma^H(z)\|_{H^{-1}}/\lambda. \quad (4.2)$$

220 In the rest of this section, we provide the proof sketch for general Moreau envelop.

221 For any integer $0 \leq t \leq T-1$, we define $r(t), k(t) \in \mathbb{N}$ such that $t = r(t)K + k(t)$ and $k(t) \leq K-1$.
 222 We will omit the dependence on t and let $r = r(t), k = k(t)$ if not causing confusion. Further define

$$x_t^m := x_{r,k}^m, g_t^m := g_{r,k}^m, \widehat{g}_t^m := \widehat{g}_{r,k}^m, u_t^m = u_{r,k}^m, v_t^m = v_{r,k}^m, H_t^m := \mathbf{diag}(\sqrt{v_t^m + \lambda^2}) \quad (4.3)$$

223 Then Algorithm 1 is equivalent to the following update rule:

$$x_{t+1}^m = \begin{cases} x_t^m - \eta(H_t^m)^{-1}u_t^m & \text{if } t \bmod K \neq -1, \\ \bar{x}_t - \eta \mathbb{E}_m[(H_t^m)^{-1}u_t^m] & \text{otherwise.} \end{cases} \quad (4.4)$$

224 Define an auxiliary sequence $\{z_t^m\}$ as:

$$z_{t+1}^m = \begin{cases} (x_{t+1}^m - \beta_1 x_t^m)/(1 - \beta_1) & \text{if } t \bmod K \neq -1, \\ (x_{t+1}^m - \beta_1 \bar{x}_t)/(1 - \beta_1) & \text{otherwise.} \end{cases} \quad (4.5)$$

225 Let $y_t := \arg \min_y f(y) + \frac{1}{2\gamma} \|y - \bar{z}_t\|_{H_{r(t)}}^2$. Define filtration $\mathcal{F}_{-1} = \emptyset, \mathcal{F}_t := \sigma(\{g_{r,k}^m\}_m \cup \mathcal{F}_{t-1})$
 226 and conditional expectation $\mathbb{E}_t[\cdot] = \mathbb{E}[\cdot | \mathcal{F}_t]$.

227 As standard practice in distributed optimization, our proof mainly contains two parts: **contraction**
 228 and **descent**. Here contraction involves showing that the iterates of local training at different workers
 229 will not diverge to different points. And decent involves showing that the objective value decreases at
 230 each iteration.

231 Our strategy is to inductively prove that some probabilistic event $E_t \in \mathcal{F}_{t-1}$ holds with high
 232 probability, which are designed to ensure contraction and descent. And event E_T can directly imply
 233 the upper bound in Theorem 3. In fact, event E_t has the form of

$$E_t = \{\mathcal{A}_{j,i} \text{ holds for all } j \leq t-1, i \in \{1, 2, 3, 4\}\}, \quad (4.6)$$

234 where $\mathcal{A}_{j,i} \in \mathcal{F}_j$ (defined later) is also some probabilistic event. As the components of E_t , each $\mathcal{A}_{j,i}$
 235 is designed to ensure either contraction or descent. We will prove the high probability bound of these
 236 components in sequence.

237 4.1 Bounding the trajectory with high probability

238 Similar to [Sadiev et al. \[2023\]](#), we only make assumptions on f and noise in certain subset $\Omega \subset \mathbb{R}^d$.
 239 This is because we are able to show that all the iterates will not leave Ω with high probability.
 240 Specifically, if it holds for all iterates before time t , using standard techniques for weakly convex
 241 optimization, we can upper bound the function value and Moreau envelope at \bar{z}_{t+1} by

$$\begin{aligned} f_\gamma^{H_{r(t+1)}}(\bar{z}_{t+1}) &\leq f_\gamma^\lambda(x_0) - \Omega(\eta) \sum_{j=0}^t \|\nabla f_\gamma^{H_{r(j)}}(\bar{z}_j)\|_{H_{r(j)}^{-1}}^2 + \underbrace{\mathcal{O}(\eta^2) \sum_{j=0}^t \|\mathbb{E}_m[\nabla f(x_j^m) - \widehat{g}_j^m]\|^2}_{\text{stochastic noise}} \\ &\quad + \underbrace{\mathcal{O}(\eta) \sum_{j=0}^t \|\nabla f(\bar{z}_j) - \mathbb{E}_m[\nabla f(x_j^m)]\|^2}_{\text{discrepancy}} \\ &\quad + \underbrace{\mathcal{O}(\eta) \sum_{j=0}^t \left\langle \bar{z}_j - \eta H_{r(j)}^{-1} \nabla f(\bar{z}_j) - y_j, \mathbb{E}_m[\mathbb{E}_j[\widehat{g}_j^m] - \widehat{g}_j^m] \right\rangle}_{\text{martingale}} \\ &\quad + \text{higher order terms.} \end{aligned} \quad (4.7)$$

242 To see that the last term is a martingale, note that $H_{r(j)}$ is independent of \widehat{g}_j^m since the stochastic
 243 gradient \widehat{g}_j^m is drawn during round r . Further note that $\mathbb{E}_j[\widehat{g}_j^m] - \widehat{g}_j^m$ is almost surely bounded
 244 thanks to clipping. Now (4.7) allows us to inductively bound $f_\gamma^{H_{r(j)}}(\bar{z}_j)$ and thus bound $\|\bar{z}_j -$
 245 $\eta H_{r(j)}^{-1} \nabla f(\bar{z}_j) - y_j\|$. After these preliminaries, we are able to apply Berstein's inequality [[Bennett,](#)
 246 [1962](#), [Freedman, 1975](#)] to control this martingale. Hence the Moreau envelope at \bar{z}_{t+1} can be bounded
 247 by a constant with high probability. Combining this with contraction results below, we can show that
 248 all the iterates stay in Ω with high probability.

249 4.2 Contraction

250 Next, we aim to show contraction, *i.e.*, $\|x_t^m - x_t^n\|$ will not diverge during local iterations with
 251 high probability. This property is crucial for showing the benefits of local updates in distributed
 252 optimization. However, different from [Woodworth et al. \[2020a\]](#), [Khaled et al. \[2020\]](#), the update of
 253 x_t^m in Algorithm 1 is in the direction of $(H_t^m)^{-1} u_t^m$, which distorts the gradient by both exponential
 254 moving average (EMA) and coordinate-wise product. Thus, the weak monotonicity (3.6) can not be
 255 directly applied as in standard analysis of gradient descent. This will further impede contraction.

256 Our solution has two steps. Firstly, we try to diminish the negative effects of different denominators
 257 used in local iterations. Then we turn to deal with the EMA of past gradient in first order momentum.

258 **Lemma 5** (Informal). *Define probabilistic events*

$$\mathcal{A}_{t,1} := \left\{ \beta_2^{K/2} \preceq H_{r(t)}^{-1} H_t^m \preceq 1 + (1 - \beta_2)B \text{ and for all } m \in [M] \right\}, \quad (4.8)$$

259

$$\mathcal{A}_{t,2} := \left\{ \|H_{r(t)}((H_t^m)^{-1} - (H_t^n)^{-1})\| \leq (1 - \beta_2)B_1 \text{ for all } m, n \in [M] \right\}, \quad (4.9)$$

260 where B, B_1 are some constants. Define $E_{t,1} := E_t \cap \mathcal{A}_{t,1}, E_{t,2} := E_{t,1} \cap \mathcal{A}_{t,2}$. For $B =$
261 $\tilde{\mathcal{O}}(K), B_1 = \tilde{\mathcal{O}}(K)$, it holds that $\mathbb{P}(E_{t,1}) \geq \mathbb{P}(E_t) - \delta/(4T), \quad \mathbb{P}(E_{t,2}) \geq \mathbb{P}(E_{t,1}) - \delta/(4T)$.

262 Event $\mathcal{A}_{t,1}$ implies the denominator of each worker during local iterations tends to be stagnant
263 and close to the averaged one after communication. Event $\mathcal{A}_{t,2}$ suggests the denominator at each
264 worker is close to each other. Note that when there is no noise, all the workers will be exactly
265 the same and then event $\mathcal{A}_{t,2}$ will always hold. Therefore, although $\mathcal{A}_{t,2}$ seems to be implied by
266 $\mathcal{A}_{t,1}$, we will be able to take $B_1 \ll B$ as long as $\sigma \ll 1$ by handling them separately. The key
267 idea to prove Lemma 5 is to control the magnitude of the EMA of squared stochastic gradients, i.e.,

$$268 \quad v_t^m = (1 - \beta_2) \sum_{j=r(t)K}^t \beta_2^{t-j} \widehat{g}_j^m{}^2 + \beta_2^{k(t)+1} v_{r(t)}. \text{ Since all the iterates stay in } \mathbf{conv}(\mathbf{B}_{R_0}(\Omega_0)), \text{ the}$$

269 squared true gradient $\nabla f(x_j^m)^2$ can be bounded. Besides, we can again apply Bernstein's inequality
270 to handle the martingale induced by $\widehat{g}_j^m{}^2 - \mathbb{E}_j[\widehat{g}_j^m{}^2]$. The remaining term $\mathbb{E}_j[\widehat{g}_j^m{}^2] - \nabla f(x_j^m)^2$ is
271 controlled by the property of clipping operator.

272 Now that the denominator is relatively stagnant, the update of x_t^m is approximately preconditioned
273 by $H_{r(t)}$ for all m . Hence we can turn to handle the first order momentum. A vanilla idea is to do the
274 following expansion:

$$\|x_{t+1}^m - x_{t+1}^n\|_{H_r}^2 \approx \|x_t^m - x_t^n\|_{H_r}^2 - 2\eta \langle x_t^m - x_t^n, u_t^m - u_t^n \rangle + \mathcal{O}(\eta^2). \quad (4.10)$$

275 By the definition of u_t^m , however, it would be influenced by noises from past stochastic gradients. In
276 this way, $u_t^m - u_t^n$ is not independent of $x_t^m - x_t^n$ and thus it is difficult to construct a martingale
277 and apply Bernstein's inequality. This is the reason why we introduce the auxiliary sequence $\{z_t^m\}$
278 defined in (4.5). Fortunately, noticing that $x_t^m - x_t^n \in \mathbf{conv}(\{z_j^m - z_j^n\}_{j \leq t})$, it suffices to show that
279 $\|z_t^m - z_t^n\|$ will not get too large with high probability.

280 **Lemma 6** (Informal). *Define probabilistic event*

$$\mathcal{A}_{t,3} := \left\{ \|z_{t+1}^m - z_{t+1}^n\|_{H_r}^2 \leq \frac{\eta^2 \sigma^2}{\lambda} KA, \sum_{j=rK}^t \|\widehat{g}_j^m\|^2 \leq \frac{(1 - \beta_1)^2 \sigma^2 A}{2^{12} (1 - \beta_2)^2 B_1^2} \text{ for all } m, n \in [M] \right\}, \quad (4.11)$$

281 where A is some constant. Define $E_{t,3} := E_{t,2} \cap \mathcal{A}_{t,3}$. For $A = \tilde{\mathcal{O}}(1)$ and $\eta =$
282 $\tilde{\mathcal{O}}(\min\{1/(K\tau), 1/L\})$, it holds that $\mathbb{P}(E_{t,3}) \geq \mathbb{P}(E_{t,2}) - \delta/(4T)$.

283 Event $\mathcal{A}_{t,3}$ is the desired contraction property and can further imply that $\|x_{t+1}^m - x_{t+1}^n\|_{H_r}^2 \leq$
284 $\frac{\eta^2 \sigma^2}{\lambda} KA$ when combined with event E_t . In fact, for $\{z_t^m\}$, we can do the following expansion:

$$\|z_{t+1}^m - z_{t+1}^n\|_{H_r}^2 \approx \|z_t^m - z_t^n\|_{H_r}^2 - 2\eta \langle z_t^m - z_t^n, \widehat{g}_t^m - \widehat{g}_t^n \rangle + \mathcal{O}(\eta^2). \quad (4.12)$$

285 Informally speaking, $\mathbb{E}_t[\widehat{g}_t^m - \widehat{g}_t^n]$ is roughly $\nabla f(x_t^m) - \nabla f(x_t^n)$, which is close to $\nabla f(z_t^m) - \nabla f(z_t^n)$
286 since $\|z_t^m - x_t^m\|^2 = \mathcal{O}(\|x_t^m - x_{t-1}^m\|^2) = \mathcal{O}(\eta^2)$. In this way, the middle term $\mathcal{O}(\eta)$ of RHS above
287 can be turned to $-2\eta \langle z_t^m - z_t^n, \nabla f(z_t^m) - \nabla f(z_t^n) \rangle$, where the weak convexity can be applied. The
288 remaining part is to control the martingale induced by $\langle z_t^m - z_t^n, \widehat{g}_t^m - \widehat{g}_t^n - \mathbb{E}_t[\widehat{g}_t^m - \widehat{g}_t^n] \rangle$ through
289 Bernstein's inequality.

290 4.3 Descent

291 Finally, we are ready to prove the descent lemma, which is the last component of E_{t+1} . Define

$$\mathcal{A}_{t,4} := \left\{ f_{\gamma}^{H_{r(t+1)}}(\bar{z}_{t+1}) - f_* + \frac{\eta}{12} \sum_{j=0}^t \|\nabla f_{\gamma}^{H_{r(j)}}(\bar{z}_j)\|_{H_{r(j)}^{-1}}^2 \leq 2\Delta \right\}. \quad (4.13)$$

292 We proceed with (4.7) and control the stochastic noise term by subtracting its expectation to construct
 293 a martingale and apply Bernstein’s inequality. Its expectation can be bounded by properties of
 294 clipping operator and variance bound. As for the discrepancy overhead, we apply the upper bound of
 295 $\|x_j^m - x_j^n\|^2$, which is induced by event E_t and utilize the $\mathcal{O}(\eta^2)$ bound on $\|\bar{z}_j - \bar{x}_j\|^2$. Therefore,
 296 thanks to all the foundations beforehand, we are able to bound each of these terms.

297 **Lemma 7** (Informal). *For sufficiently small η , it holds that $\mathbb{P}(E_{t+1}) \geq \mathbb{P}(E_{t,3}) - \delta/(4T)$.*

298 Therefore, we prove that $\mathbb{P}(E_{t+1}) \geq \mathbb{P}(E_t) - \delta/T$. And by induction rule, $\mathbb{P}(E_T) \geq 1 - \delta$. After
 299 carefully choosing the learning rate η , we complete the proof of Theorem 3.

300 5 Conclusion

301 In this paper, we prove the benefits of local updates within distributed adaptive methods to reduce
 302 communication complexity compared to their minibatch counterparts. We study *Local* SGDM
 303 and *Local* Adam under convex and weakly convex setting, respectively. We consider generalized
 304 smoothness assumption and gradient clipping, and develop a novel technique to show contraction
 305 during local updates. Future works may include improved analysis of *Local* Adam, benefits of local
 306 adaptive algorithms in non-convex setting, advantages over non-adaptive methods, etc.

307 References

- 308 Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and
 309 Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC*
 310 *conference on computer and communications security*, pages 308–318, 2016.
- 311 Ahmet Alacaoglu, Yura Malitsky, and Volkan Cevher. Convergence of adaptive algorithms for weakly
 312 convex constrained optimization. *arXiv preprint arXiv:2006.06650*, 2020.
- 313 Yossi Arjevani, Yair Carmon, John C Duchi, Dylan J Foster, Nathan Srebro, and Blake Woodworth.
 314 Lower bounds for non-convex stochastic optimization. *Mathematical Programming*, 199(1-2):
 315 165–214, 2023.
- 316 George Bennett. Probability inequalities for the sum of independent random variables. *Journal of the*
 317 *American Statistical Association*, 57(297):33–45, 1962.
- 318 Jeremy Bernstein, Yu-Xiang Wang, Kamyar Azizzadenesheli, and Animashree Anandkumar. signsgd:
 319 Compressed optimisation for non-convex problems. In *International Conference on Machine*
 320 *Learning*, pages 560–569. PMLR, 2018.
- 321 Xiangyi Chen, Xiaoyun Li, and Ping Li. Toward communication efficient adaptive gradient method.
 322 In *Proceedings of the 2020 ACM-IMS on Foundations of Data Science Conference*, pages 119–128,
 323 2020.
- 324 Ziheng Cheng, Xinmeng Huang, Pengfei Wu, and Kun Yuan. Momentum benefits non-iid federated
 325 learning simply and provably. *arXiv preprint arXiv:2306.16504*, 2023.
- 326 Michael Crawshaw, Mingrui Liu, Francesco Orabona, Wei Zhang, and Zhenxun Zhuang. Robustness
 327 to unbounded smoothness of generalized signsgd. *Advances in Neural Information Processing*
 328 *Systems*, 35:9955–9968, 2022.
- 329 Ashok Cutkosky and Harsh Mehta. High-probability bounds for non-convex stochastic optimization
 330 with heavy tails. *Advances in Neural Information Processing Systems*, 34:4883–4895, 2021.
- 331 Damek Davis and Dmitriy Drusvyatskiy. Stochastic model-based minimization of weakly convex
 332 functions. *SIAM Journal on Optimization*, 29(1):207–239, 2019.
- 333 Alexandre Défossez, Léon Bottou, Francis Bach, and Nicolas Usunier. A simple convergence proof
 334 of adam and adagrad. *arXiv preprint arXiv:2003.02395*, 2020.
- 335 Qi Deng and Wenzhi Gao. Minibatch and momentum model-based methods for stochastic weakly
 336 convex optimization. *Advances in Neural Information Processing Systems*, 34:23115–23127, 2021.

- 337 Arthur Douillard, Qixuan Feng, Andrei A Rusu, Rachita Chhaparia, Yani Donchev, Adhiguna
338 Kuncoro, Marc’Aurelio Ranzato, Arthur Szlam, and Jiajun Shen. Diloco: Distributed low-
339 communication training of language models. *arXiv preprint arXiv:2311.08105*, 2023.
- 340 Matthew Faw, Litu Rout, Constantine Caramanis, and Sanjay Shakkottai. Beyond uniform smooth-
341 ness: A stopped analysis of adaptive sgd. In *The Thirty Sixth Annual Conference on Learning*
342 *Theory*, pages 89–160. PMLR, 2023.
- 343 David A Freedman. On tail probabilities for martingales. *the Annals of Probability*, pages 100–118,
344 1975.
- 345 Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, and Yann N Dauphin. Convolutional
346 sequence to sequence learning. In *International conference on machine learning*, pages 1243–1252.
347 PMLR, 2017.
- 348 Margalit R Glasgow, Honglin Yuan, and Tengyu Ma. Sharp bounds for federated averaging (local sgd)
349 and continuous perspective. In *International Conference on Artificial Intelligence and Statistics*,
350 pages 9050–9090. PMLR, 2022.
- 351 Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016.
- 352 Eduard Gorbunov, Marina Danilova, and Alexander Gasnikov. Stochastic optimization with heavy-
353 tailed noise via accelerated gradient clipping. *Advances in Neural Information Processing Systems*,
354 33:15042–15053, 2020.
- 355 Eduard Gorbunov, Abdurakhmon Sadiev, Marina Danilova, Samuel Horváth, Gauthier Gidel, Pavel
356 Dvurechensky, Alexander Gasnikov, and Peter Richtárik. High-probability convergence for
357 composite and distributed stochastic minimization and variational inequalities with heavy-tailed
358 noise. *arXiv preprint arXiv:2310.01860*, 2023.
- 359 Zhishuai Guo, Yi Xu, Wotao Yin, Rong Jin, and Tianbao Yang. A novel convergence analysis for
360 algorithms of the adam family. *arXiv preprint arXiv:2112.03459*, 2021.
- 361 Nicholas JA Harvey, Christopher Liaw, and Sikander Randhawa. Simple and optimal high-probability
362 bounds for strongly-convex stochastic gradient descent. *arXiv preprint arXiv:1909.00843*, 2019.
- 363 Sai Praneeth Karimireddy, Martin Jaggi, Satyen Kale, Mehryar Mohri, Sashank J Reddi, Sebastian U
364 Stich, and Ananda Theertha Suresh. Mime: Mimicking centralized stochastic algorithms in
365 federated learning. *arXiv preprint arXiv:2008.03606*, 2020a.
- 366 Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and
367 Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In
368 *International conference on machine learning*, pages 5132–5143. PMLR, 2020b.
- 369 Sai Praneeth Karimireddy, Lie He, and Martin Jaggi. Learning from history for byzantine robust
370 optimization. In *International Conference on Machine Learning*, pages 5311–5319. PMLR, 2021.
- 371 Ahmed Khaled, Konstantin Mishchenko, and Peter Richtárik. Tighter theory for local sgd on identical
372 and heterogeneous data. In *International Conference on Artificial Intelligence and Statistics*, pages
373 4519–4529. PMLR, 2020.
- 374 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
375 *arXiv:1412.6980*, 2014.
- 376 Frederik Kunstner, Jacques Chen, Jonathan Wilder Lavington, and Mark Schmidt. Noise is not the
377 main factor behind the gap between sgd and adam on transformers, but sign descent might be.
378 *arXiv preprint arXiv:2304.13960*, 2023.
- 379 Haochuan Li, Jian Qian, Yi Tian, Alexander Rakhlin, and Ali Jadbabaie. Convex and non-convex
380 optimization under generalized smoothness. *Advances in Neural Information Processing Systems*,
381 36, 2024a.
- 382 Haochuan Li, Alexander Rakhlin, and Ali Jadbabaie. Convergence of adam under relaxed assumptions.
383 *Advances in Neural Information Processing Systems*, 36, 2024b.

- 384 Xiaoyu Li and Francesco Orabona. A high probability analysis of adaptive sgd with momentum.
385 *arXiv preprint arXiv:2007.14294*, 2020.
- 386 Hong Liu, Zhiyuan Li, David Hall, Percy Liang, and Tengyu Ma. Sophia: A scalable stochastic
387 second-order optimizer for language model pre-training. *arXiv preprint arXiv:2305.14342*, 2023.
- 388 Mingrui Liu, Zhenxun Zhuang, Yunwen Lei, and Chunyang Liao. A communication-efficient
389 distributed gradient clipping algorithm for training deep neural networks. *Advances in Neural
390 Information Processing Systems*, 35:26204–26217, 2022.
- 391 Wei Liu, Li Chen, Yunfei Chen, and Wenyi Zhang. Accelerating federated learning via momentum
392 gradient descent. *IEEE Transactions on Parallel and Distributed Systems*, 31(8):1754–1766,
393 2020a.
- 394 Yanli Liu, Yuan Gao, and Wotao Yin. An improved analysis of stochastic gradient descent with
395 momentum. *Advances in Neural Information Processing Systems*, 33:18261–18271, 2020b.
- 396 Liangchen Luo, Yuanhao Xiong, Yan Liu, and Xu Sun. Adaptive gradient methods with dynamic
397 bound of learning rate. *arXiv preprint arXiv:1902.09843*, 2019.
- 398 Vien Mai and Mikael Johansson. Convergence of a stochastic gradient method with momentum for
399 non-smooth non-convex optimization. In *International conference on machine learning*, pages
400 6630–6639. PMLR, 2020.
- 401 Vien V Mai and Mikael Johansson. Stability and convergence of stochastic gradient clipping: Beyond
402 lipschitz continuity and smoothness. In *International Conference on Machine Learning*, pages
403 7325–7335. PMLR, 2021.
- 404 Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas.
405 Communication-efficient learning of deep networks from decentralized data. In *Artificial intelli-
406 gence and statistics*, pages 1273–1282. PMLR, 2017.
- 407 Stephen Merity, Nitish Shirish Keskar, and Richard Socher. Regularizing and optimizing lstm
408 language models. *arXiv preprint arXiv:1708.02182*, 2017.
- 409 Yan Pan and Yuanzhi Li. Toward understanding why adam converges faster than sgd for transformers.
410 *arXiv preprint arXiv:2306.00204*, 2023.
- 411 Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent neural
412 networks. In *International conference on machine learning*, pages 1310–1318. Pmlr, 2013.
- 413 Kumar Kshitij Patel, Lingxiao Wang, Blake E Woodworth, Brian Bullins, and Nati Srebro. Towards
414 optimal communication complexity in distributed non-convex optimization. *Advances in Neural
415 Information Processing Systems*, 35:13316–13328, 2022.
- 416 V. V. Petrov. Moments of sums of independent random variables. *Journal of Soviet Mathematics*, 61
417 (1):1905–1906, Aug 1992. ISSN 1573-8795. doi: 10.1007/BF01362802. URL [https://doi.
418 org/10.1007/BF01362802](https://doi.org/10.1007/BF01362802).
- 419 Sashank Reddi, Zachary Charles, Manzil Zaheer, Zachary Garrett, Keith Rush, Jakub Konečný,
420 Sanjiv Kumar, and H Brendan McMahan. Adaptive federated optimization. *arXiv preprint
421 arXiv:2003.00295*, 2020.
- 422 Sashank J Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. *arXiv
423 preprint arXiv:1904.09237*, 2019.
- 424 Abdurakhmon Sadiev, Marina Danilova, Eduard Gorbunov, Samuel Horváth, Gauthier Gidel, Pavel
425 Dvurechensky, Alexander Gasnikov, and Peter Richtárik. High-probability bounds for stochastic
426 optimization and variational inequalities: the case of unbounded variance. In *Proceedings of the
427 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine
428 Learning Research*, pages 29563–29648. PMLR, 23–29 Jul 2023. URL [https://proceedings.
429 mlr.press/v202/sadiev23a.html](https://proceedings.mlr.press/v202/sadiev23a.html).

- 430 Othmane Sebbouh, Robert M Gower, and Aaron Defazio. Almost sure convergence rates for stochastic
431 gradient descent and stochastic heavy ball. In *Conference on Learning Theory*, pages 3935–3971.
432 PMLR, 2021.
- 433 Naichen Shi, Dawei Li, Mingyi Hong, and Ruoyu Sun. Rmsprop converges with proper hyper-
434 parameter. In *International Conference on Learning Representations*, 2020.
- 435 Bohan Wang, Yushun Zhang, Huishuai Zhang, Qi Meng, Zhi-Ming Ma, Tie-Yan Liu, and Wei Chen.
436 Provable adaptivity in adam. *arXiv preprint arXiv:2208.09900*, 2022.
- 437 Bohan Wang, Jingwen Fu, Huishuai Zhang, Nanning Zheng, and Wei Chen. Closing the gap between
438 the upper bound and lower bound of adam’s iteration complexity. *Advances in Neural Information
439 Processing Systems*, 36, 2024.
- 440 Jianyu Wang, Vinayak Tantia, Nicolas Ballas, and Michael Rabbat. Slowmo: Improving
441 communication-efficient distributed sgd with slow momentum. *arXiv preprint arXiv:1910.00643*,
442 2019.
- 443 Jianyu Wang, Zheng Xu, Zachary Garrett, Zachary Charles, Luyang Liu, and Gauri Joshi. Local
444 adaptivity in federated learning: Convergence and consistency. *arXiv preprint arXiv:2106.02305*,
445 2021.
- 446 Jue Wang, Yucheng Lu, Binhang Yuan, Beidi Chen, Percy Liang, Christopher De Sa, Christopher
447 Re, and Ce Zhang. Cocktailsgd: fine-tuning foundation models over 500mbps networks. In
448 *International Conference on Machine Learning*, pages 36058–36076. PMLR, 2023.
- 449 Jianqiao Wangni, Jialei Wang, Ji Liu, and Tong Zhang. Gradient sparsification for communication-
450 efficient distributed optimization. *Advances in Neural Information Processing Systems*, 31, 2018.
- 451 Blake Woodworth, Kumar Kshitij Patel, Sebastian Stich, Zhen Dai, Brian Bullins, Brendan McMahan,
452 Ohad Shamir, and Nathan Srebro. Is local sgd better than minibatch sgd? In *International
453 Conference on Machine Learning*, pages 10334–10343. PMLR, 2020a.
- 454 Blake E Woodworth, Kumar Kshitij Patel, and Nati Srebro. Minibatch vs local sgd for heterogeneous
455 distributed learning. *Advances in Neural Information Processing Systems*, 33:6281–6292, 2020b.
- 456 Jing Xu, Sen Wang, Liwei Wang, and Andrew Chi-Chih Yao. Fedcm: Federated learning with
457 client-level momentum. *arXiv preprint arXiv:2106.10874*, 2021.
- 458 Haibo Yang, Minghong Fang, and Jia Liu. Achieving linear speedup with partial worker participation
459 in non-iid federated learning. *arXiv preprint arXiv:2101.11203*, 2021.
- 460 Hao Yu, Rong Jin, and Sen Yang. On the linear speedup analysis of communication efficient
461 momentum sgd for distributed non-convex optimization. In *International Conference on Machine
462 Learning*, pages 7184–7193. PMLR, 2019.
- 463 Honglin Yuan and Tengyu Ma. Federated accelerated stochastic gradient descent. *Advances in Neural
464 Information Processing Systems*, 33:5332–5344, 2020.
- 465 Jingzhao Zhang, Tianxing He, Suvrit Sra, and Ali Jadbabaie. Why gradient clipping accelerates
466 training: A theoretical justification for adaptivity. *arXiv preprint arXiv:1905.11881*, 2019.
- 467 Jingzhao Zhang, Sai Praneeth Karimireddy, Andreas Veit, Seungyeon Kim, Sashank Reddi, Sanjiv
468 Kumar, and Suvrit Sra. Why are adaptive methods good for attention models? *Advances in Neural
469 Information Processing Systems*, 33:15383–15393, 2020.
- 470 Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher
471 Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language
472 models. *arXiv preprint arXiv:2205.01068*, 2022a.
- 473 Yushun Zhang, Congliang Chen, Naichen Shi, Ruoyu Sun, and Zhi-Quan Luo. Adam can converge
474 without any modification on update rules. *Advances in Neural Information Processing Systems*, 35:
475 28386–28399, 2022b.

- 476 Weijie Zhao, Xuewu Jiao, Mingqing Hu, Xiaoyun Li, Xiangyu Zhang, and Ping Li. Communication-
477 efficient terabyte-scale model training framework for online advertising. *arXiv preprint*
478 *arXiv:2201.05500*, 2022.
- 479 Fangyu Zou, Li Shen, Zequn Jie, Weizhong Zhang, and Wei Liu. A sufficient condition for conver-
480 gences of adam and rmsprop. In *Proceedings of the IEEE/CVF Conference on computer vision*
481 *and pattern recognition*, pages 11127–11135, 2019.

482 A Related Work

483 **Theoretical benefits of local updates in distributed optimization.** Algorithms with local updates
484 have been used among practitioners for a long time to reduce communication complexity [McMahan
485 et al., 2017]. In the homogeneous and convex setting, *Local* SGD and its variants have been shown to
486 outperform the minibatch baseline, for a fixed amount of gradient computations and communication
487 rounds. Woodworth et al. [2020a] is the first to show that *Local* SGD can provably outperform
488 *Minibatch* SGD. Yuan and Ma [2020] develops FedAC to further accelerate *Local* SGD. In the
489 heterogeneous case, Woodworth et al. [2020b] demonstrates the advantages of *Local* SGD when
490 heterogeneity is very low. Algorithms with local updates have also been studied in the non-convex
491 setting [Karimireddy et al., 2020b, Yang et al., 2021, Glasgow et al., 2022], including momentum-
492 based and adaptive methods [Reddi et al., 2020, Karimireddy et al., 2020a], though no advantage of
493 local iterations over minibatch has been shown, without non-standard assumptions such as 3rd-order
494 smoothness. Notably, Liu et al. [2022] is one closely related work to ours, which considers *Local*
495 SGD with gradient clipping in homogeneous and non-convex setting and claims that the convergence
496 guarantee is better than naive parallel of centralized clipped-SGD. However, it still cannot outperform
497 minibatch baseline (with batch size K for each worker in each round) and thus fails to demonstrate
498 the benefits of local iterations.

499 **Convergence of centralized Adam.** Adam was first proposed by Kingma and Ba [2014] with
500 convergence guarantee in online convex optimization. However, Reddi et al. [2019] found a gap
501 in the original analysis of Adam and constructed a counter example to show its divergence. Since
502 then, many works have developed convergence analyses of Adam with various assumptions and
503 hyper-parameter settings. Guo et al. [2021] assumed the denominator is bounded from below and
504 above by two constants, which typically requires a bounded gradient assumption or the AdaBound
505 variant [Luo et al., 2019]. Défossez et al. [2020] assumed a bounded gradient and their convergence
506 guarantee depends on $\text{poly}(d)$. Zhang et al. [2022b], Wang et al. [2022] considered a finite sum
507 setting and showed that Adam converges to the neighborhood of stationary points. One closely related
508 work to ours is Li et al. [2024b], which established a high probability bound without a bounded
509 gradient assumption. However they assumed that noise is bounded almost surely. Another recent
510 work [Wang et al., 2024] provided a guarantee of $\mathcal{O}\left(\frac{1}{\varepsilon^4}\right)$ with dependence on $\text{poly}(d)$. Beyond the
511 guarantees on gradient norm given by non-convex analyses, no stronger bounds (e.g., on function
512 error) are known for Adam in the convex case.

513 **Convergence of distributed adaptive algorithms.** In the federated learning literature, Reddi et al.
514 [2020] introduced a framework, FedOPT, to leverage both worker optimizer and server optimizer.
515 Many works explored adaptive server optimizer while fixing worker side as vanilla SGD. The
516 theoretical results of local adaptive algorithms are much fewer. Some works have studied *Local*
517 Adam and *Local* AMSGrad with fixed momentum state during local iterations [Karimireddy et al.,
518 2020a, Chen et al., 2020, Zhao et al., 2022]. They also needed stringent assumptions such as a
519 huge batch size depending on the inverse of target error, bounded stochastic gradients, vanishing
520 difference between denominator, etc., which are not standard. Wang et al. [2021] explored adaptive
521 worker optimizer based on centralized algorithm, where the state of worker optimizer changes in local
522 updates. However, their analysis relied on an explicit assumptions [Wang et al., 2021, Assumption 1]
523 on the contraction property of worker optimizer. To the best of our knowledge, there is no end-to-end
524 convergence guarantee for distributed adaptive algorithms with local iterations.

525 **Gradient clipping.** Pascanu et al. [2013] first proposed gradient clipping technique to address the
526 issue of exploding gradient problem of deep neural networks. Since then, it has become standard
527 practice in the training of language models [Gehring et al., 2017, Merity et al., 2017, Zhang et al.,
528 2022a, Liu et al., 2023]. Furthermore, from theoretical perspective, gradient clipping is also used
529 for multiple purposes, including differential privacy [Abadi et al., 2016], distributed optimization
530 [Karimireddy et al., 2021, Liu et al., 2022], heavy-tailed noise [Zhang et al., 2020].

531 **Generalized smoothness.** The generalized smoothness condition was initially proposed by [Zhang
532 et al., 2019] to justify gradient clipping, and was called (L_0, L_1) -smoothness. The empirical evidence
533 therein illustrated that the norm of Hessian matrix of language models depends linearly on the
534 magnitude of gradient, contradicting the standard L -smoothness. A recent work [Li et al., 2024a]

535 further generalized this condition to ℓ -smoothness and proved convergence of classical SGD in this
536 setting. Apart from bounding the Hessian through gradient, [Sadiev et al. \[2023\]](#) proposed to assume
537 that the norm of Hessian is uniformly bounded in certain subset of whole space, in order to get
538 high probability bounds for (accelerated) clipped-SGD. [Gorbunov et al. \[2023\]](#) further extended this
539 setting to composite and distributed optimization without local updates. Here we follow the setting of
540 [\[Sadiev et al., 2023\]](#) since (L_0, L_1) -smoothness would reduce to it in most cases. See Section 2.1 for
541 details.

542 B Technical Lemmas

543 **Lemma B.1** ([\[Bennett, 1962, Freedman, 1975\]](#)). *Let the sequence of random variables $\{X_i\}_{i \geq 1}$
544 form a martingale difference sequence, i.e. $\mathbb{E}[X_i | X_{i-1}, \dots, X_1] = 0$ for all $i \geq 1$. Assume that
545 conditional variances $\sigma_i^2 \stackrel{\text{def}}{=} \mathbb{E}[X_i^2 | X_{i-1}, \dots, X_1]$ exist and are bounded and assume also that there
546 exists deterministic constant $c > 0$ such that $|X_i| \leq c$ almost surely for all $i \geq 1$. Then for all
547 $b > 0, V > 0$ and $n \geq 1$,*

$$\mathbb{P} \left\{ \left| \sum_{i=1}^n X_i \right| > b \text{ and } \sum_{i=1}^n \sigma_i^2 \leq V \right\} \leq 2 \exp \left(-\frac{b^2}{2V + 2cb/3} \right). \quad (\text{B.1})$$

548 **Lemma B.2.** *Let X be a random variable in \mathbb{R} and $\tilde{X} := \text{clip}(X, \rho)$, Then $\|\tilde{X} - \mathbb{E}\tilde{X}\| \leq 2\rho$.
549 Moreover, if for some $\sigma > 0$ and $\alpha \geq 2$,*

$$\mathbb{E}[X] = x \in \mathbb{R}, \quad \mathbb{E}|X - x|^\alpha \leq \sigma^\alpha, \quad (\text{B.2})$$

550 and $|x| \leq \frac{\rho}{2}, \rho \geq 3\sigma$, then

$$|\mathbb{E}[\tilde{X}] - x| \leq \frac{(2\sigma)^\alpha}{\rho^{\alpha-1}}, \quad \mathbb{E}|\tilde{X} - x|^\alpha \leq \sigma^\alpha, \quad \mathbb{E}|\tilde{X} - \mathbb{E}[\tilde{X}]|^\alpha \leq (2\sigma)^\alpha. \quad (\text{B.3})$$

551 *Proof.* The first claim is from [\[Sadiev et al., 2023\]](#) and we show the proof here for completeness. To
552 start the proof, we introduce two indicator random variables. Let

$$\chi = \mathbb{I}_{\{|X| > \rho\}} = \begin{cases} 1, & \text{if } |X| > \rho, \\ 0, & \text{otherwise} \end{cases}, \quad \eta = \mathbb{I}_{\{|X - x| > \frac{\rho}{2}\}} = \begin{cases} 1, & \text{if } |X - x| > \frac{\rho}{2}, \\ 0, & \text{otherwise} \end{cases}. \quad (\text{B.4})$$

553 Moreover, since $|X| \leq |x| + |X - x| \leq \frac{\rho}{2} + |X - x|$, we have $\chi \leq \eta$. Using that

$$\tilde{X} = \min \left\{ 1, \frac{\rho}{|X|} \right\} X = \chi \frac{\rho}{|X|} X + (1 - \chi)X, \quad (\text{B.5})$$

554 we obtain

$$\begin{aligned} |\mathbb{E}[\tilde{X}] - x| &= \left| \mathbb{E} \left[X + \chi \left(\frac{\rho}{|X|} - 1 \right) X \right] - x \right| \\ &= \left| \mathbb{E} \left[\chi \left(\frac{\rho}{|X|} - 1 \right) X \right] \right| \\ &= \mathbb{E} \left[\chi \left(1 - \frac{\rho}{|X|} \right) |X| \right]. \end{aligned} \quad (\text{B.6})$$

555 Since $1 - \frac{\rho}{|X|} \in (0, 1)$ when $\chi \neq 0$, we derive

$$\begin{aligned} |\mathbb{E}[\tilde{X}] - x| &\leq \mathbb{E}[\chi |X|] \\ &\leq \mathbb{E}[\eta |X|] \\ &\leq \mathbb{E}[\eta |X - x| + \eta |x|] \\ &\leq (\mathbb{E}[|X - x|^\alpha])^{\frac{1}{\alpha}} (\mathbb{E}[\eta^{\frac{\alpha-1}{\alpha}}])^{\frac{\alpha-1}{\alpha}} + |x| \mathbb{E}[\eta] \\ &\stackrel{\eta \in \{0,1\}}{\leq} \sigma (\mathbb{E}[\eta])^{\frac{\alpha-1}{\alpha}} + \frac{\rho}{2} \mathbb{E}[\eta], \end{aligned} \quad (\text{B.7})$$

556 By Markov's inequality,

$$\begin{aligned}\mathbb{E}[\eta] &= \mathbb{P}\left\{|X - x|^\alpha > \frac{\rho^\alpha}{2^\alpha}\right\} \\ &\leq \frac{2^\alpha}{\rho^\alpha} \mathbb{E}[|X - x|^\alpha] \\ &\leq \left(\frac{2\sigma}{\rho}\right)^\alpha.\end{aligned}\tag{B.8}$$

557 Thus, in combination with the previous chain of inequalities, we finally have

$$|\mathbb{E}[\tilde{X}] - x| \leq \sigma \left(\frac{2\sigma}{\rho}\right)^{\alpha-1} + \frac{\rho}{2} \left(\frac{2\sigma}{\rho}\right)^\alpha = \frac{2^\alpha \sigma^\alpha}{\rho^{\alpha-1}}.\tag{B.9}$$

558 For the second part, since

$$|\tilde{X} - x| = |\mathbf{clip}(X, \rho) - \mathbf{clip}(x, \rho)| \leq |X - x|,\tag{B.10}$$

559 hence $\mathbb{E}|\tilde{X} - x|^\alpha \leq \mathbb{E}|X - x|^\alpha \leq \sigma^\alpha$. By Jensen's inequality, we have for any $q \in (0, 1)$,

$$\begin{aligned}\mathbb{E}|\tilde{X} - \mathbb{E}[\tilde{X}]|^\alpha &\leq q^{1-\alpha} \mathbb{E}|\tilde{X} - x|^\alpha + (1-q)^{1-\alpha} \mathbb{E}|\tilde{X}] - x|^\alpha \\ &\leq q^{1-\alpha} \sigma^\alpha + (1-q)^{1-\alpha} \left(\frac{(2\sigma)^\alpha}{\rho^{\alpha-1}}\right)^\alpha.\end{aligned}\tag{B.11}$$

560 Choose the optimal $q = \frac{\sigma}{\sigma + \frac{(2\sigma)^\alpha}{\rho^{\alpha-1}}}$ and we can conclude that

$$\mathbb{E}|\tilde{X} - \mathbb{E}[\tilde{X}]|^\alpha \leq \left(\sigma + \frac{(2\sigma)^\alpha}{\rho^{\alpha-1}}\right)^\alpha \leq (2\sigma)^\alpha.\tag{B.12}$$

561 This completes the proof. \square

562 **Lemma B.3.** For M independent random vectors $X_1, \dots, X_M \in \mathbb{R}^d$ such that $\mathbb{E}[X_m] = 0$,
563 $\mathbb{E}[\|X_m\|^4] \leq \sigma^4$, the following holds

$$\mathbb{E}[\|\mathbb{E}_m X_m\|^2]^2 \leq \frac{4\sigma^4}{M^2}.\tag{B.13}$$

564 *Proof.* We prove by direct calculation as follows:

$$\begin{aligned}\mathbb{E}[\|\mathbb{E}_m X_m\|^2]^2 &\leq \mathbb{E}\left[\frac{1}{M^2} \sum_m \|X_m\|^2 + \frac{2}{M^2} \sum_{m < n} \langle X_m, X_n \rangle\right]^2 \\ &= \mathbb{E}\left[\frac{1}{M^2} \sum_m \|X_m\|^2\right]^2 + \mathbb{E}\left[\frac{2}{M^2} \sum_{m < n} \langle X_m, X_n \rangle\right]^2 \\ &\leq \frac{\sigma^4}{M^2} + \frac{4}{M^4} \mathbb{E} \sum_{m < n} \langle X_m, X_n \rangle^2 \\ &\leq \frac{4\sigma^4}{M^2}.\end{aligned}\tag{B.14}$$

565 \square

566 **Lemma B.4.** For any set $\Omega \in \mathbb{R}^d$ and $r > 0$, define $\mathbf{B}_r(\Omega) :=$
567 $\{x \in \mathbb{R}^d : \exists y \in \Omega, \text{ s.t., } \|x - y\| \leq r\}$. Then

$$\mathbf{B}_r(\mathbf{conv}(\Omega)) = \mathbf{conv}(\mathbf{B}_r(\Omega)).\tag{B.15}$$

568 *Proof.* For any $x \in \mathbf{B}_r(\mathbf{conv}(\Omega))$, there exist $y_1, \dots, y_N \in \Omega$ and $(\lambda_1, \dots, \lambda_N) \in \Delta^N$ for some N ,
569 such that

$$\|x - y\| \leq r, \quad y := \sum_{n=1}^N \lambda_n y_n. \quad (\text{B.16})$$

570 Then $x = y + (x - y) = \sum_{n=1}^N \lambda_n (y_n + x - y) = \sum_{n=1}^N \lambda_n x_n$, where

$$x_n = y_n + x - y \in B_r(\Omega). \quad (\text{B.17})$$

571 Hence $x \in \mathbf{conv}(\mathbf{B}_r(\Omega))$.

572 On the other hand, for any $x \in \mathbf{conv}(\mathbf{B}_r(\Omega))$, there exist $x_1, \dots, x_N \in \mathbf{B}_r(\Omega)$, $y_1, \dots, y_N \in \Omega$ and
573 $(\lambda_1, \dots, \lambda_N) \in \Delta^N$, such that

$$x = \sum_{n=1}^N \lambda_n x_n, \quad \|x_n - y_n\| \leq r. \quad (\text{B.18})$$

574 Let $y := \sum_{n=1}^N \lambda_n y_n \in \mathbf{conv}(\Omega)$. Then $\|x - y\| \leq \sum_{n=1}^N \lambda_n \|x_n - y_n\| \leq r$ and thus $x \in \mathbf{B}_r(\mathbf{conv}(\Omega))$.
575 □

576 C Proof of Local SGDM

577 We restate the *Local SGDM* algorithm here.

Algorithm 2 *Local SGDM*

Require: initial model x_0 , learning rate η , momentum $\beta_1 \in [0, 1]$

Set $x_{0,0}^m = x_0$, $u_{0,-1}^m = 0$ for each worker $m \in [M]$

for $r = 0, \dots, R - 1$ **do**

for each worker $m \in [M]$ **in parallel do**

for $k = 0, \dots, K - 1$ **do**

$g_{r,k}^m = \nabla F(x_{r,k}^m; \xi_{r,k}^m)$, $\widehat{g}_{r,k}^m = \mathbf{clip}(g_{r,k}^m, \rho)$ ▷ Compute clipped stochastic gradient

$u_{r,k}^m = \beta_1 u_{r,k-1}^m + (1 - \beta_1) \widehat{g}_{r,k}^m$ ▷ Update momentum

$x_{r,k+1}^m = x_{r,k}^m - \eta u_{r,k}^m$ ▷ Update model

end for

end for

$x_{r+1,0}^m = \mathbb{E}_m[x_{r,K}^m]$, $u_{r+1,-1}^m = \mathbb{E}_m[u_{r,K-1}^m]$ ▷ Communicate and average

end for

578 C.1 Overview and Main Theorem

579 For any integer $0 \leq t \leq T - 1$, we define $r(t), k(t) \in \mathbb{N}$ such that $t = r(t)K + k(t)$ and $k(t) \leq K - 1$.

580 We omit the dependence on t and let $r = r(t), k = k(t)$ through out the proof if not causing confusion.

581 Define $x_t^m := x_{r,k}^m$, $g_t^m := g_{r,k}^m$, $\widehat{g}_t^m := \widehat{g}_{r,k}^m$, $u_t^m = u_{r,k}^m$. Then Algorithm 2 is equivalent to the
582 following update rule:

$$u_t^m = \begin{cases} \beta_1 u_{t-1}^m + (1 - \beta_1) \widehat{g}_t^m & \text{if } t \bmod K \neq 0, \\ \beta_1 \bar{u}_{t-1} + (1 - \beta_1) \widehat{g}_t^m & \text{otherwise,} \end{cases} \quad (\text{C.1})$$

583

$$x_{t+1}^m = \begin{cases} x_t^m - \eta u_t^m & \text{if } t \bmod K \neq -1, \\ \bar{x}_t - \eta \bar{u}_t & \text{otherwise.} \end{cases} \quad (\text{C.2})$$

584 Define an auxiliary sequence $\{z_t^m\}$ as:

$$z_{t+1}^m = \begin{cases} \frac{1}{1 - \beta_1} x_{t+1}^m - \frac{\beta_1}{1 - \beta_1} x_t^m & \text{if } t \bmod K \neq -1, \\ \frac{1}{1 - \beta_1} x_{t+1}^m - \frac{\beta_1}{1 - \beta_1} \bar{x}_t & \text{otherwise.} \end{cases} \quad (\text{C.3})$$

585 Define probabilistic events (see (C.12) for definition of some parameters)

$$\mathcal{A}_{t,1} := \{\|z_{t+1}^m - z_{t+1}^n\|^2 \leq \eta^2 \sigma^2 K A \text{ for all } m, n \in [M]\}, \quad (\text{C.4})$$

586

$$\mathcal{A}_{t,2} := \left\{ \sum_{j=0}^t \frac{\eta}{2} (f(\bar{z}_j) - f_*) (1 - \frac{\eta\mu}{2})^{t-j} + \|\bar{z}_{t+1} - x_*\|^2 \leq 2(1 - \frac{\eta\mu}{2})^{t+1} D_0^2 \right\}. \quad (\text{C.5})$$

587 Besides, let

$$E_t := \{\mathcal{A}_{j,i} \text{ holds for all } j \leq t-1, i \in \{1, 2\}\}, \quad E_{t,1} := E_t \cap \mathcal{A}_{t,1}. \quad (\text{C.6})$$

588 Now we present two of our major lemmas, the first of which is to show contraction and the second is
589 a descent lemma.

590 **Lemma C.1.** *Let $A := \max \left\{ \frac{2^{10} \rho^2 d}{K \sigma^2} \log^2 \frac{MT}{\delta}, 2^9 \log \frac{MT}{\delta}, 2^{12} \frac{K \|2\sigma\|_{2\alpha}^{2\alpha}}{\sigma^2 \rho^{2(\alpha-1)}} \right\}$. If $\eta \leq$
591 $\min \left\{ \frac{(1-\beta_1)^2}{2L}, \frac{D_0}{4\sigma\sqrt{KA}} \right\}$ and $\rho \geq \max\{3\sigma_\infty, 2G_\infty\}$, then the following holds:*

$$\mathbb{P}(E_{t,1}) \geq \mathbb{P}(E_t) - \frac{\delta}{2T}. \quad (\text{C.7})$$

592 **Lemma C.2.** *For any $\varepsilon > 0$, let*

$$\rho \geq \begin{cases} \max \left\{ \left(\frac{2^8 \|2\sigma\|_{2\alpha}^{2\alpha}}{\mu\varepsilon} \right)^{\frac{1}{2(\alpha-1)}}, 3\sigma_\infty, 2G_\infty \right\}, & \text{if } \mu > 0, \\ \max \left\{ \left(\frac{2^8 D_0 \|2\sigma\|_{2\alpha}^\alpha}{\varepsilon} \right)^{\frac{1}{\alpha-1}}, 3\sigma_\infty, 2G_\infty \right\}, & \text{otherwise.} \end{cases} \quad (\text{C.8})$$

$$\eta := \begin{cases} \frac{2}{\mu T} \log \frac{4\mu D_0^2}{\varepsilon}, & \text{if } \mu > 0, \\ \frac{4D_0^2}{T\varepsilon}, & \text{otherwise.} \end{cases}$$

593 *If*

$$\eta \lesssim \begin{cases} \min \left\{ \frac{(1-\beta_1)^2}{L}, \frac{M\varepsilon}{\sigma^2 \log \frac{T}{\delta}}, \left(\frac{L\sigma^2 KA}{\varepsilon} \right)^{-1/2}, \frac{\sqrt{\varepsilon/\mu}}{\rho\sqrt{d} \log \frac{T}{\delta}} \right\}, & \text{if } \mu > 0, \\ \min \left\{ \frac{(1-\beta_1)^2}{L}, \frac{M\varepsilon}{\sigma^2 \log \frac{T}{\delta}}, \left(\frac{L\sigma^2 KA}{\varepsilon} \right)^{-1/2}, \frac{D_0}{\rho\sqrt{d} \log \frac{T}{\delta}} \right\}, & \text{otherwise,} \end{cases} \quad (\text{C.9})$$

594 *where A is defined in Lemma C.1, then the following holds*

$$\mathbb{P}(E_{t+1}) \geq \mathbb{P}(E_{t,1}) - \frac{\delta}{2T}. \quad (\text{C.10})$$

595 The following is our main result, from which we will parse the implications in Theorems 1 and 2.

596 **Theorem C.3.** *Let Assumption 1, 2, 3, 4 hold for $\Omega := \{\|x - x_*\| \leq \sqrt{3}D_0\}$. Further assume that
597 for any $x \in \Omega$, $\|\nabla f(x)\|_\infty \leq G_\infty$. Then with probability $\geq 1 - \delta$, Local SGDM yields $f(\hat{x}) - f_* \leq \varepsilon$
598 if*

$$T \gtrsim \begin{cases} \log \frac{\mu D_0^2}{\varepsilon} \left[\frac{L}{(1-\beta_1)^2 \mu} + \frac{\sigma^2}{\mu M \varepsilon} \log \frac{T}{\delta} + \sqrt{\frac{L\sigma^2 KA}{\mu^2 \varepsilon}} + \frac{\rho\sqrt{d}}{\sqrt{\mu\varepsilon}} \log \frac{T}{\delta} \right], & \text{if } \mu > 0, \\ \frac{D_0^2}{\varepsilon} \left[\frac{L}{(1-\beta_1)^2} + \frac{\sigma^2}{M\varepsilon} \log \frac{T}{\delta} + \sqrt{\frac{L\sigma^2 KA}{\varepsilon}} + \frac{\rho\sqrt{d}}{D_0} \log \frac{T}{\delta} \right], & \text{otherwise.} \end{cases} \quad (\text{C.11})$$

599 Here

$$\begin{aligned}
\rho &\geq \begin{cases} \max \left\{ \left(\frac{2^8 \|\mathbf{2}\boldsymbol{\sigma}\|_{2\alpha}^{2\alpha}}{\mu\varepsilon} \right)^{\frac{1}{2(\alpha-1)}}, 3\sigma_\infty, 2G_\infty \right\}, & \text{if } \mu > 0, \\ \max \left\{ \left(\frac{2^8 D_0 \|\mathbf{2}\boldsymbol{\sigma}\|_{2\alpha}^\alpha}{\varepsilon} \right)^{\frac{1}{\alpha-1}}, 3\sigma_\infty, 2G_\infty \right\}, & \text{otherwise,} \end{cases} \\
A &:= \max \left\{ \frac{2^{10} \rho^2 d}{K\sigma^2} \log^2 \frac{MT}{\delta}, 2^9 \log \frac{MT}{\delta}, 2^{12} \frac{K \|\mathbf{2}\boldsymbol{\sigma}\|_{2\alpha}^{2\alpha}}{\sigma^2 \rho^{2(\alpha-1)}} \right\}, \\
\eta &:= \begin{cases} \frac{2}{\mu T} \log \frac{4\mu D_0^2}{\varepsilon}, & \text{if } \mu > 0, \\ \frac{4D_0^2}{T\varepsilon}, & \text{otherwise.} \end{cases}
\end{aligned} \tag{C.12}$$

600 *Proof.* We prove by induction that $\mathbb{P}(E_t) \geq 1 - \frac{t\delta}{T}$ for $t = 0, \dots, T$.

601 When $t = 0$, this is trivial. Assume that the statement is true for some $t \leq T - 1$. We aim to prove
602 that $\mathbb{P}(E_{t+1}) \geq 1 - \frac{(t+1)\delta}{T}$. It is easy to verify the conditions in Lemma C.1, C.2 once (C.11) and
603 (C.12) hold. Hence we have

$$\mathbb{P}(E_{t+1}) \geq \mathbb{P}(E_t) - 2 \cdot \frac{\delta}{2T} \geq 1 - \frac{(t+1)\delta}{T}. \tag{C.13}$$

604 Therefore by induction rule, $\mathbb{P}(E_T) \geq 1 - \delta$ and this implies by event $\mathcal{A}_{T,2}$ that

$$\sum_{j=0}^{T-1} \frac{\eta}{2} (f(\bar{z}_j) - f_*) \left(1 - \frac{\eta\mu}{2}\right)^{T-j} \leq 2 \left(1 - \frac{\eta\mu}{2}\right)^T D_0^2. \tag{C.14}$$

605 Let $\hat{x} := \frac{\eta\mu \sum_{j=0}^{T-1} \left(1 - \frac{\eta\mu}{2}\right)^{T-j} \bar{z}_j}{2(1 - (1 - \frac{\eta\mu}{2})^T)}$. By convexity, we have

$$f(\hat{x}) - f_* \leq \frac{2(1 - \frac{\eta\mu}{2})^T \mu D_0^2}{1 - (1 - \frac{\eta\mu}{2})^T}. \tag{C.15}$$

606 (1) **Case** $\mu > 0$.

$$f(\hat{x}) - f_* \leq \frac{2(1 - \frac{\eta\mu}{2})^T \mu D_0^2}{1 - (1 - \frac{\eta\mu}{2})^T} \leq 4(1 - \frac{\eta\mu}{2})^T \mu D_0^2 \leq 4e^{-\eta\mu T/2} \mu D_0^2 = \varepsilon. \tag{C.16}$$

607 (2) **Case** $\mu = 0$.

$$f(\hat{x}) - f_* \leq \frac{2(1 - \frac{\eta\mu}{2})^T \mu D_0^2}{1 - (1 - \frac{\eta\mu}{2})^T} = \frac{4D_0^2}{\eta T} = \varepsilon. \tag{C.17}$$

608 \square

609 We now state and prove the implications of Theorem C.3 which yield the results stated in the main
610 body of our paper.

611 **Theorem C.4** (Complete version of Theorem 1). *Under the conditions of Theorem C.3 and $\mu > 0$,*

612 *assume $1 - \beta_1 = \Omega(1)$, $\left(\frac{\|\boldsymbol{\sigma}\|_{2\alpha}^{2\alpha}}{\mu\varepsilon}\right)^{\frac{1}{2(\alpha-1)}} \gtrsim G_\infty \vee \sigma_\infty$, and $K \gtrsim \log \frac{MT}{\delta} \left(\frac{\|\boldsymbol{\sigma}\|_{2\alpha} d^{\frac{1}{2} - \frac{1}{2\alpha}}}{\sigma}\right)^{\frac{2\alpha}{\alpha-2}}$.*

613 *Then with probability no less than $1 - \delta$, Local SGDM with optimal η, ρ yields $f(\hat{x}) - f_* \leq \varepsilon$, if*

$$T \gtrsim \log \frac{\mu D_0^2}{\varepsilon} \left[\frac{L}{\mu} + \frac{\sigma^2}{\mu M \varepsilon} \log \frac{T}{\delta} + \sqrt{\frac{L\sigma^2 K \log \frac{MT}{\delta}}{\mu^2 \varepsilon}} + \sqrt{\frac{Ld}{\mu^2 \varepsilon}} \log \frac{MT}{\delta} \left(\frac{\|\boldsymbol{\sigma}\|_{2\alpha}^{2\alpha}}{\mu\varepsilon}\right)^{\frac{1}{2(\alpha-1)}} \right]. \tag{C.18}$$

614 And equivalently, let $\kappa := L/\mu$,

$$f(\hat{x}) - f_* \lesssim \exp\left(-\Theta\left(\frac{\mu KR}{L}\right)\right) + \frac{\sigma^2 \log(MKR)}{\mu MKR} \log \frac{KR}{\delta} \\ + \frac{L\sigma^2 \log^2(KR)}{\mu^2 KR^2} \log \frac{MKR}{\delta} + \frac{\|\sigma\|_{2\alpha}^2 (\kappa d)^{\frac{\alpha-1}{\alpha}}}{\mu} \left(\frac{\log \frac{MKR}{\delta}}{KR}\right)^{\frac{2(\alpha-1)}{\alpha}}. \quad (\text{C.19})$$

615 *Proof.* Plug the definition of A in (C.11),

$$T \gtrsim \log \frac{\mu D_0^2}{\varepsilon} \left[\frac{L}{\mu} + \frac{\sigma^2}{\mu M \varepsilon} \log \frac{T}{\delta} + \sqrt{\frac{L\sigma^2 K \log \frac{MT}{\delta}}{\mu^2 \varepsilon}} + \frac{\rho \sqrt{d}}{\sqrt{\mu \varepsilon}} \log \frac{T}{\delta} \right] \\ + \log \frac{\mu D_0^2}{\varepsilon} \sqrt{\frac{LK}{\mu^2 \varepsilon}} \sqrt{\frac{\rho^2 d}{K} \log^2 \frac{MT}{\delta} + \frac{K \|2\sigma\|_{2\alpha}^{2\alpha}}{\rho^{2(\alpha-1)}}} \\ \asymp \log \frac{\mu D_0^2}{\varepsilon} \left[\frac{L}{\mu} + \frac{\sigma^2}{\mu M \varepsilon} \log \frac{T}{\delta} + \sqrt{\frac{L\sigma^2 K \log \frac{MT}{\delta}}{\mu^2 \varepsilon}} \right] \\ + \log \frac{\mu D_0^2}{\varepsilon} \sqrt{\frac{LK}{\mu^2 \varepsilon}} \sqrt{\frac{\rho^2 d}{K} \log^2 \frac{MT}{\delta} + \frac{K \|2\sigma\|_{2\alpha}^{2\alpha}}{\rho^{2(\alpha-1)}}}. \quad (\text{C.20})$$

616 Hence the optimal ρ is given by

$$\rho \asymp \max \left\{ \|\sigma\|_{2\alpha} \left(\frac{K}{\sqrt{d} \log \frac{MT}{\delta}} \right)^{1/\alpha}, \left(\frac{\|\sigma\|_{2\alpha}^{2\alpha}}{\mu \varepsilon} \right)^{\frac{1}{2(\alpha-1)}}, \sigma_\infty, G_\infty \right\}. \quad (\text{C.21})$$

617 Note that $\left(\frac{\|\sigma\|_{2\alpha}^{2\alpha}}{\mu \varepsilon} \right)^{\frac{1}{2(\alpha-1)}} \gtrsim G_\infty \vee \sigma_\infty$ and this implies

$$T \gtrsim \log \frac{\mu D_0^2}{\varepsilon} \left[\frac{L}{\mu} + \frac{\sigma^2}{\mu M \varepsilon} \log \frac{T}{\delta} + \sqrt{\frac{L\sigma^2 K \log \frac{MT}{\delta}}{\mu^2 \varepsilon}} \right] \\ + \log \frac{\mu D_0^2}{\varepsilon} \sqrt{\frac{L}{\mu^2 \varepsilon} \cdot \left[\|\sigma\|_{2\alpha}^2 K^{\frac{2}{\alpha}} \left(d \log^2 \frac{MT}{\delta} \right)^{1-\frac{1}{\alpha}} + \left(\frac{\|\sigma\|_{2\alpha}^{2\alpha}}{\mu \varepsilon} \right)^{\frac{1}{\alpha-1}} d \log^2 \frac{MT}{\delta} \right]} \\ \asymp \log \frac{\mu D_0^2}{\varepsilon} \left[\frac{L}{\mu} + \frac{\sigma^2}{\mu M \varepsilon} \log \frac{T}{\delta} + \sqrt{\frac{L\sigma^2 K \log \frac{MT}{\delta}}{\mu^2 \varepsilon}} + \sqrt{\frac{Ld}{\mu^2 \varepsilon}} \log \frac{MT}{\delta} \left(\frac{\|\sigma\|_{2\alpha}^{2\alpha}}{\mu \varepsilon} \right)^{\frac{1}{2(\alpha-1)}} \right]. \quad (\text{C.22})$$

618 In the last equation we use $K \gtrsim \log \frac{MT}{\delta} \left(\frac{\|\sigma\|_{2\alpha} d^{\frac{1}{2}-\frac{1}{2\alpha}}}{\sigma} \right)^{\frac{2\alpha}{\alpha-2}}$. This completes the proof. \square

619 **Theorem C.5** (Complete version of Theorem 2). Under the conditions of Theorem C.3 and $\mu = 0$,

620 assume $1 - \beta_1 = \Omega(1)$, $\left(\frac{D_0 \|\sigma\|_{2\alpha}^\alpha}{\varepsilon} \right)^{\frac{1}{\alpha-1}} \gtrsim G_\infty \vee \sigma_\infty$, and $K \gtrsim \log \frac{MT}{\delta} \left(\frac{\|\sigma\|_{2\alpha} d^{\frac{1}{2}-\frac{1}{2\alpha}}}{\sigma} \right)^{\frac{2\alpha}{\alpha-2}}$.

621 Then with probability no less than $1 - \delta$, Local SGDM with optimal η, ρ yields $f(\hat{x}) - f_* \leq \varepsilon$ if

$$T \gtrsim \frac{D_0^2}{\varepsilon} \left[L + \frac{\sigma^2}{M \varepsilon} \log \frac{T}{\delta} + \sqrt{\frac{L\sigma^2 K \log \frac{MT}{\delta}}{\varepsilon}} + \sqrt{\frac{dL}{\varepsilon}} \left(\frac{D_0 \|\sigma\|_{2\alpha}^\alpha}{\varepsilon} \right)^{\frac{1}{\alpha-1}} \log \frac{MT}{\delta} \right]. \quad (\text{C.23})$$

622 And equivalently,

$$\begin{aligned}
f(\hat{x}) - f_* &\lesssim \frac{LD_0^2}{KR} + \frac{\sigma D_0}{\sqrt{MKR}} \log^{\frac{1}{2}} \frac{KR}{\delta} \\
&\quad + \frac{L^{\frac{1}{3}} \sigma^{\frac{2}{3}} D_0^{\frac{4}{3}}}{K^{\frac{1}{3}} R^{\frac{2}{3}}} \log^{\frac{1}{3}} \frac{MKR}{\delta} + \left(\|\sigma\|_{2\alpha}^{\frac{2\alpha}{\alpha-1}} dLD_0 \right)^{\frac{\alpha-1}{3\alpha-1}} D_0 \left(\frac{\log \frac{MKR}{\delta}}{KR} \right)^{\frac{2(\alpha-1)}{3\alpha-1}}.
\end{aligned} \tag{C.24}$$

623 *Proof.* Plug the definition of A in (C.11),

$$\begin{aligned}
T &\gtrsim \frac{D_0^2}{\varepsilon} \left[L + \frac{\sigma^2}{M\varepsilon} \log \frac{T}{\delta} + \sqrt{\frac{L\sigma^2 K \log \frac{MT}{\delta}}{\varepsilon}} + \frac{\rho\sqrt{d}}{D_0} \log \frac{T}{\delta} \right] \\
&\quad + \frac{D_0^2}{\varepsilon} \sqrt{\frac{LK}{\varepsilon}} \sqrt{\frac{\rho^2 d}{K} \log^2 \frac{MT}{\delta} + \frac{K\|\sigma\|_{2\alpha}^{2\alpha}}{\rho^{2(\alpha-1)}}} \\
&\asymp \frac{D_0^2}{\varepsilon} \left[L + \frac{\sigma^2}{M\varepsilon} \log \frac{T}{\delta} + \sqrt{\frac{L\sigma^2 K \log \frac{MT}{\delta}}{\varepsilon}} + \sqrt{\frac{LK}{\varepsilon}} \sqrt{\frac{\rho^2 d}{K} \log^2 \frac{MT}{\delta} + \frac{K\|\sigma\|_{2\alpha}^{2\alpha}}{\rho^{2(\alpha-1)}}} \right].
\end{aligned} \tag{C.25}$$

624 Hence the optimal ρ is given by

$$\rho \asymp \max \left\{ \|\sigma\|_{2\alpha} \left(\frac{K}{\sqrt{d} \log \frac{MT}{\delta}} \right)^{1/\alpha}, \left(\frac{D_0 \|\sigma\|_{2\alpha}^\alpha}{\varepsilon} \right)^{\frac{1}{\alpha-1}}, \sigma_\infty, G_\infty \right\}. \tag{C.26}$$

625 Note that $\left(\frac{D_0 \|\sigma\|_{2\alpha}^\alpha}{\varepsilon} \right)^{\frac{1}{\alpha-1}} \gtrsim G_\infty \vee \sigma_\infty$ and this implies

$$\begin{aligned}
T &\gtrsim \frac{D_0^2}{\varepsilon} \left[L + \frac{\sigma^2}{M\varepsilon} \log \frac{T}{\delta} + \sqrt{\frac{L\sigma^2 K \log \frac{MT}{\delta}}{\varepsilon}} \right] \\
&\quad + \frac{D_0^2}{\varepsilon} \sqrt{\frac{L}{\varepsilon} \cdot \left[\|\sigma\|_{2\alpha}^2 K^{\frac{2}{\alpha}} \left(d \log^2 \frac{MT}{\delta} \right)^{1-\frac{1}{\alpha}} + \left(\frac{D_0 \|\sigma\|_{2\alpha}^\alpha}{\varepsilon} \right)^{\frac{2}{\alpha-1}} d \log^2 \frac{MT}{\delta} \right]} \\
&\asymp \frac{D_0^2}{\varepsilon} \left[L + \frac{\sigma^2}{M\varepsilon} \log \frac{T}{\delta} + \sqrt{\frac{L\sigma^2 K \log \frac{MT}{\delta}}{\varepsilon}} + \sqrt{\frac{dL}{\varepsilon}} \left(\frac{D_0 \|\sigma\|_{2\alpha}^\alpha}{\varepsilon} \right)^{\frac{1}{\alpha-1}} \log \frac{MT}{\delta} \right].
\end{aligned} \tag{C.27}$$

626 In the last equation we use $K \gtrsim \log \frac{MT}{\delta} \left(\frac{\|\sigma\|_{2\alpha} d^{\frac{1}{2} - \frac{1}{2\alpha}}}{\sigma} \right)^{\frac{2\alpha}{\alpha-2}}$. Solve ε and we get the upper
627 bound of $f(\hat{x}) - f_*$. This completes the proof. \square

628 C.2 Preliminaries

629 In this subsection, we show that event E_t implies all the iterates remain in certain area, so that we
630 can apply all kinds of properties of f afterwards.

631 **Lemma C.6.** *If $\eta\sigma\sqrt{KA} \leq (\sqrt{3} - \sqrt{2})D_0$, Event E_t implies that for all $j \leq t, m \in [M]$, we have
632 $x_j^m, \bar{x}_j, z_j^m, \bar{z}_j \in \Omega$. And $\|x_j^m - x_j^n\| \leq \eta\sigma\sqrt{KA}$ for all m, n .*

633 *Proof.* Event E_t implies that for all $j \leq t$,

$$\|\bar{z}_j - x_*\| \leq \sqrt{2}D_0, \quad \|z_j^m - z_j^n\| \leq \eta\sigma\sqrt{KA} \leq (\sqrt{3} - \sqrt{2})D_0. \tag{C.28}$$

634 Hence $\bar{z}_j \in \Omega$, $\|z_j^m - x_*\| \leq \sqrt{3}D_0$ and $z_j^m \in \Omega$. Also, notice that $\bar{x}_j \in \mathbf{conv}\{\bar{z}_i\}_{i \leq j}$ and
 635 $x_j^m - x_j^n \in \mathbf{conv}\{z_i^m - z_i^n\}_{i \leq j}$. We have

$$\|\bar{x}_j - x_*\| \leq \sqrt{2}D_0, \|x_j^m - x_j^n\| \leq \eta\sigma\sqrt{KA}, \|x_j^m - \bar{x}_j\| \leq \eta\sigma\sqrt{KA} \leq (\sqrt{3} - \sqrt{2})D_0. \quad (\text{C.29})$$

636 Therefore $x_j^m, \bar{x}_j \in \Omega$. This completes the proof. \square

637 C.3 Proof of Contraction Lemma C.1

638 In this subsection, we aim to show contraction, *i.e.*, $\|x_t^m - x_t^n\|$ won't be too large during local
 639 iterations with high probability. This property is crucial for showing the benefits of local updates in
 640 distributed optimization. However, different from [Woodworth et al., 2020a, Khaled et al., 2020],
 641 the update of x_t^m is in the direction of momentum u_t^m , which incorporates information from all past
 642 gradient. Therefore, we cannot directly apply $\langle x_t^m - x_t^n, \mathbb{E}_t[u_t^m - u_t^n] \rangle \geq 0$. Fortunately, noticing
 643 that $x_t^m - x_t^n \in \mathbf{conv}\{z_j^m - z_j^n\}_{j \leq t}$, it suffices to show that $\|z_t^m - z_t^n\|$ won't get too large with
 644 high probability. Besides, the update rule of z_t^m is much easier to handle.

645 *Proof.* First note that by the upper bound of η , Lemma C.6 holds. Since $z_{t+1}^m = z_t^m - \eta\widehat{g}_t^m$,

$$\begin{aligned} \|z_{t+1}^m - z_{t+1}^n\|^2 &= \|z_t^m - z_t^n\|^2 - 2\eta \langle z_t^m - z_t^n, \widehat{g}_t^m - \widehat{g}_t^n \rangle + \eta^2 \|\widehat{g}_t^m - \widehat{g}_t^n\|^2 \\ &\leq \|z_t^m - z_t^n\|^2 - 2\eta \langle z_t^m - z_t^n, \nabla f(x_t^m) - \nabla f(x_t^n) \rangle + 2\eta^2 \|\nabla f(x_t^m) - \nabla f(x_t^n)\|^2 \\ &\quad + 2\eta \langle z_t^m - z_t^n, \nabla f(x_t^m) - \nabla f(x_t^n) - \widehat{g}_t^m + \widehat{g}_t^n \rangle + 2\eta^2 \|\nabla f(x_t^m) - \nabla f(x_t^n) - \widehat{g}_t^m + \widehat{g}_t^n\|^2. \end{aligned} \quad (\text{C.30})$$

646 Event E_t implies $z_t^m, x_t^m \in \Omega$ and thus by $\forall x, y \in \Omega, \langle x - y, \nabla f(x) - \nabla f(y) \rangle \geq \frac{1}{L} \|\nabla f(x) -$
 647 $\nabla f(y)\|^2$,

$$\begin{aligned} \langle z_t^m - z_t^n, \nabla f(x_t^m) - \nabla f(x_t^n) \rangle &= \langle x_t^m - x_t^n, \nabla f(x_t^m) - \nabla f(x_t^n) \rangle + \langle z_t^m - z_t^n - (x_t^m - x_t^n), \nabla f(x_t^m) - \nabla f(x_t^n) \rangle \\ &\geq \langle x_t^m - x_t^n, \nabla f(x_t^m) - \nabla f(x_t^n) \rangle \\ &\quad - \left[L \|z_t^m - z_t^n - (x_t^m - x_t^n)\|^2 + \frac{1}{4L} \|\nabla f(x_t^m) - \nabla f(x_t^n)\|^2 \right] \\ &\geq \frac{3}{4L} \|\nabla f(x_t^m) - \nabla f(x_t^n)\|^2 - L \|z_t^m - z_t^n - (x_t^m - x_t^n)\|^2. \end{aligned} \quad (\text{C.31})$$

648 Therefore, for the second and third term in the RHS of (C.30),

$$\begin{aligned} -2\eta \langle z_t^m - z_t^n, \nabla f(x_t^m) - \nabla f(x_t^n) \rangle + 2\eta^2 \|\nabla f(x_t^m) - \nabla f(x_t^n)\|^2 \\ \leq -\frac{\eta}{L} \|\nabla f(x_t^m) - \nabla f(x_t^n)\|^2 + 2\eta L \|z_t^m - z_t^n - (x_t^m - x_t^n)\|^2. \end{aligned} \quad (\text{C.32})$$

649 By the update rule,

$$\begin{aligned} \|z_t^m - z_t^n - (x_t^m - x_t^n)\|^2 &= \left(\frac{\eta\beta_1}{1 - \beta_1} \right)^2 \|u_{t-1}^m - u_{t-1}^n\|^2 \\ &\leq \left(\frac{\eta\beta_1}{1 - \beta_1} \right)^2 \left\| (1 - \beta_1) \sum_{j=rK}^{t-1} \beta_1^{t-j-1} [\widehat{g}_k^m - \widehat{g}_k^n] \right\|^2 \\ &\leq \frac{2(\eta\beta_1)^2}{1 - \beta_1} \sum_{j=rK}^{t-1} \beta_1^{t-j-1} \left[\|\nabla f(x_j^m) - \nabla f(x_j^n)\|^2 + \|\widehat{g}_j^m - \widehat{g}_j^n - \nabla f(x_j^m) + \nabla f(x_j^n)\|^2 \right]. \end{aligned} \quad (\text{C.33})$$

650 Let $S_t := \sum_{j=rK}^t \beta_1^{t-j} \|\nabla f(x_j^m) - \nabla f(x_j^n)\|^2$. We further get

$$\begin{aligned} \text{LHS of (C.32)} &\leq -\frac{\eta}{L}(S_t - \beta_1 S_{t-1}) + \frac{4\eta L(\eta\beta_1)^2}{1-\beta_1} \left[S_{t-1} + \sum_{j=rK}^{t-1} \beta_1^{t-j-1} [\|\widehat{g}_j^m - \widehat{g}_j^n - \nabla f(x_j^m) + \nabla f(x_j^n)\|^2] \right] \\ &= -\frac{\eta}{L}(S_t - S_{t-1}) + \frac{4\eta L(\eta\beta_1)^2}{1-\beta_1} \left[\sum_{j=rK}^{t-1} \beta_1^{t-j-1} [\|\widehat{g}_j^m - \widehat{g}_j^n - \nabla f(x_j^m) + \nabla f(x_j^n)\|^2] \right] \end{aligned} \quad (\text{C.34})$$

651 Then plug in (C.30),

$$\begin{aligned} \|z_{t+1}^m - z_{t+1}^n\|^2 &\leq \|z_t^m - z_t^n\|^2 - \frac{\eta}{L}(S_t - S_{t-1}) \\ &\quad + \frac{4\eta L(\eta\beta_1)^2}{1-\beta_1} \left[\sum_{j=rK}^{t-1} \beta_1^{t-j-1} [\|\widehat{g}_j^m - \widehat{g}_j^n - \nabla f(x_j^m) + \nabla f(x_j^n)\|^2] \right] \\ &\quad + 2\eta \left\langle z_t^m - z_t^n, \nabla f(x_t^m) - \nabla f(x_t^n) - \widehat{g}_t^m + \widehat{g}_t^n \right\rangle + 2\eta^2 \|\widehat{g}_t^m - \widehat{g}_t^n - \nabla f(x_t^m) + \nabla f(x_t^n)\|^2. \end{aligned} \quad (\text{C.35})$$

652 Notice that this recursive bound holds for any $rK \leq i \leq t$. Unroll it and recalculate the coefficients
653 using $\eta L \leq (1-\beta_1)^2/2$,

$$\begin{aligned} \|z_{t+1}^m - z_{t+1}^n\|^2 + \frac{\eta}{L} S_t &\leq \sum_{j=rK}^t 2\eta \left\langle z_j^m - z_j^n, \nabla f(x_j^m) - \nabla f(x_j^n) - \widehat{g}_j^m + \widehat{g}_j^n \right\rangle \\ &\quad + \sum_{j=rK}^t 4\eta^2 \|\nabla f(x_j^m) - \nabla f(x_j^n) - \widehat{g}_j^m + \widehat{g}_j^n\|^2 \\ &\leq \underbrace{\sum_{j=rK}^t 2\eta \left\langle z_j^m - z_j^n, \mathbb{E}_j[\widehat{g}_j^m - \widehat{g}_j^n] - [\widehat{g}_j^m - \widehat{g}_j^n] \right\rangle}_{\text{①: martingale}} \\ &\quad + \underbrace{\sum_{j=rK}^t 2\eta \left\langle z_j^m - z_j^n, \nabla f(x_j^m) - \nabla f(x_j^n) - \mathbb{E}_j[\widehat{g}_j^m - \widehat{g}_j^n] \right\rangle}_{\text{②: clipping bias}} \\ &\quad + \underbrace{\sum_{j=rK}^t 4\eta^2 \left[\|\nabla f(x_j^m) - \nabla f(x_j^n) - \widehat{g}_j^m + \widehat{g}_j^n\|^2 - \mathbb{E}_j[\|\nabla f(x_j^m) - \nabla f(x_j^n) - [\widehat{g}_j^m - \widehat{g}_j^n]\|^2] \right]}_{\text{③: martingale}} \\ &\quad + 4\eta^2 K \cdot 2\sigma^2. \end{aligned} \quad (\text{C.36})$$

654 For ①, define

$$\zeta_j^{m,n} = \begin{cases} 2\eta \left\langle z_j^m - z_j^n, \mathbb{E}_j[\widehat{g}_j^m - \widehat{g}_j^n] - [\widehat{g}_j^m - \widehat{g}_j^n] \right\rangle, & \text{if event } E_j \text{ holds,} \\ 0, & \text{otherwise.} \end{cases} \quad (\text{C.37})$$

655 Then since event E_j implies $\|z_j^m - z_j^n\| \leq \eta\sigma\sqrt{KA}$,

$$|\zeta_j^{m,n}| \leq 2\eta \cdot \eta\sigma\sqrt{KA} \cdot 2\rho\sqrt{d} = 4\eta^2\sigma\rho\sqrt{dKA} \stackrel{\text{def}}{=} c, \quad (\text{C.38})$$

656

$$\text{Var}_j(\zeta_j^{m,n}) \leq 4\eta^2 \cdot \eta^2\sigma^2 KA \cdot 2\sigma^2 = 8\eta^4\sigma^4 KA. \quad (\text{C.39})$$

657 Let $b = \frac{1}{4}\eta^2\sigma^2KA$, $V = 8\eta^4\sigma^4K^2A$. By Lemma B.1, $|\sum_{j=0}^t \zeta_j^{m,n}| \leq b$ with probability no less than

$$1 - 2\exp\left(\frac{b^2}{2V + 2cb/3}\right) \geq 1 - \frac{\delta}{4M^2T}. \quad (\text{C.40})$$

658 For ②,

$$|\textcircled{2}| \leq 2\eta K \cdot \eta\sigma\sqrt{KA} \cdot 2\frac{\|2\sigma\|_{2\alpha}^\alpha}{\rho^{(\alpha-1)}} \leq \frac{1}{4}\eta^2\sigma^2KA. \quad (\text{C.41})$$

659 For ③, define

$$\theta_j^{m,n} = \begin{cases} 4\eta^2 \left[\|\nabla f(x_j^m) - \nabla f(x_j^n) - \widehat{g}_j^m + \widehat{g}_j^n\|^2 - \mathbb{E}_j[\|\nabla f(x_j^m) - \nabla f(x_j^n) - [\widehat{g}_j^m - \widehat{g}_j^n]\|^2] \right], & \text{if event } E_j \text{ holds,} \\ 0, & \text{otherwise.} \end{cases} \quad (\text{C.42})$$

660 Then,

$$|\theta_j^{m,n}| \leq 4\eta^2 \cdot 4\rho^2d = 16\eta^2\rho^2d \stackrel{\text{def}}{=} c, \quad (\text{C.43})$$

661

$$\text{Var}_j(\theta_j^{m,n}) \leq 16\eta^4 \cdot \mathbb{E}_j[\|\nabla f(x_j^m) - \nabla f(x_j^n) - [\widehat{g}_j^m - \widehat{g}_j^n]\|^2]^2 \leq 64\eta^4\sigma^4. \quad (\text{C.44})$$

662 Let $b = \frac{1}{4}\eta^2\sigma^2KA$, $V = 64K\eta^4\sigma^4$. By Lemma B.1, $|\sum_{j=0}^t \theta_j^{m,n}| \leq b$ with probability no less than

$$1 - 2\exp\left(\frac{b^2}{2V + 2cb/3}\right) \geq 1 - \frac{\delta}{4M^2T}. \quad (\text{C.45})$$

663 Combine ①, ②, ③ and thus we can conclude that with probability no less than $\mathbb{P}(E_t) - 2 \cdot \frac{\delta}{4T}$, event
664 E_t holds and $\|z_{t+1}^m - z_{t+1}^n\|^2 \leq \eta^2\sigma^2KA$ for all m, n . This completes the proof. \square

665 C.4 Proof of Descent Lemma C.2

666 Now we are ready to state the main descent lemma of *Local* SGDM.

667 *Proof.* Again, note that by the upper bound of η , Lemma C.6 holds. Under event E_t ,

$$\begin{aligned} \|\bar{z}_{t+1} - x_*\|^2 &= \|\bar{z}_t - x_*\|^2 - 2\eta \langle \bar{z}_t - x_*, \mathbb{E}_m[\widehat{g}_t^m] \rangle + \eta^2 \|\mathbb{E}_m[\widehat{g}_t^m]\|^2 \\ &\leq \|\bar{z}_t - x_*\|^2 - 2\eta \langle \bar{z}_t - x_*, \mathbb{E}_m[\nabla f(x_t^m)] \rangle - 2\eta \langle \bar{z}_t - x_*, \mathbb{E}_m[\widehat{g}_t^m - \nabla f(x_t^m)] \rangle \\ &\quad + 2\eta^2 \|\mathbb{E}_m[\widehat{g}_t^m - \nabla f(x_t^m)]\|^2 + 2\eta^2 \|\mathbb{E}_m[\nabla f(x_t^m)]\|^2. \end{aligned} \quad (\text{C.46})$$

668 Since $x_t^m, \bar{x}_t, \bar{z}_t \in \Omega$, for the second term,

$$\begin{aligned} \langle \bar{z}_t - x_*, \mathbb{E}_m[\nabla f(x_t^m)] \rangle &= \langle \bar{x}_t - x_*, \mathbb{E}_m[\nabla f(x_t^m)] \rangle + \langle \bar{z}_t - \bar{x}_t, \mathbb{E}_m[\nabla f(x_t^m)] \rangle \\ &= \mathbb{E}_m[\langle \bar{x}_t - x_t^m, \nabla f(x_t^m) \rangle + \langle x_t^m - x_*, \nabla f(x_t^m) \rangle] \\ &\quad + \langle \bar{z}_t - \bar{x}_t, \nabla f(\bar{x}_t) \rangle + \langle \bar{z}_t - \bar{x}_t, \mathbb{E}_m[\nabla f(x_t^m) - \nabla f(\bar{x}_t)] \rangle. \end{aligned} \quad (\text{C.47})$$

669 By smoothness,

$$\mathbb{E}_m[\langle \bar{x}_t - x_t^m, \nabla f(x_t^m) \rangle] \geq -L\mathbb{E}_m[\|x_t^m - \bar{x}_t\|^2], \quad (\text{C.48})$$

670

$$f(\bar{z}_t) \leq f(\bar{x}_t) + \langle \bar{z}_t - \bar{x}_t, \nabla f(\bar{x}_t) \rangle + \frac{L}{2}\|\bar{x}_t - \bar{z}_t\|^2. \quad (\text{C.49})$$

671 By μ -strong convexity,

$$\begin{aligned} \mathbb{E}_m[\langle x_t^m - x_*, \nabla f(x_t^m) \rangle] &\geq \mathbb{E}_m[f(x_t^m) - f_* + \frac{\mu}{2}\|x_t^m - x_*\|^2] \\ &\geq f(\bar{x}_t) - f_* + \frac{\mu}{2}\|\bar{x}_t - x_*\|^2. \end{aligned} \quad (\text{C.50})$$

672 Therefore,

$$\begin{aligned}
\langle \bar{z}_t - x_*, \mathbb{E}_m[\nabla f(x_t^m)] \rangle &= \langle \bar{x}_t - x_*, \mathbb{E}_m[\nabla f(x_t^m)] \rangle + \langle \bar{z}_t - \bar{x}_t, \mathbb{E}_m[\nabla f(x_t^m)] \rangle \\
&\stackrel{(C.48), (C.50)}{\geq} f(\bar{x}_t) - f_* + \frac{\mu}{2} \|\bar{x}_t - x_*\|^2 - L \mathbb{E}_m[\|x_t^m - \bar{x}_t\|^2] \\
&\quad + \langle \bar{z}_t - \bar{x}_t, \nabla f(\bar{x}_t) \rangle + \langle \bar{z}_t - \bar{x}_t, \mathbb{E}_m[\nabla f(x_t^m) - \nabla f(\bar{x}_t)] \rangle \\
&\stackrel{(C.49), \text{AM-GM}}{\geq} f(\bar{z}_t) - f_* + \frac{\mu}{2} \|\bar{x}_t - x_*\|^2 - \frac{L}{2} \|\bar{z}_t - \bar{x}_t\|^2 - L \mathbb{E}_m[\|x_t^m - \bar{x}_t\|^2] \\
&\quad - \frac{L}{2} (\|\bar{z}_t - \bar{x}_t\|^2 + \mathbb{E}_m[\|x_t^m - \bar{x}_t\|^2]) \\
&\stackrel{\text{AM-GM}}{\geq} f(\bar{z}_t) - f_* + \frac{\mu}{4} \|\bar{z}_t - x_*\|^2 - \frac{3L}{2} (\|\bar{z}_t - \bar{x}_t\|^2 + \mathbb{E}_m[\|x_t^m - \bar{x}_t\|^2]).
\end{aligned} \tag{C.51}$$

673 For the last term in (C.46),

$$\begin{aligned}
2\eta^2 \|\mathbb{E}_m[\nabla f(x_t^m)]\|^2 &\leq 6\eta^2 [L^2 \|x_t^m - \bar{x}_t\|^2 + L^2 \|\bar{x}_t - \bar{z}_t\|^2 + \|\nabla f(\bar{z}_t)\|^2] \\
&\leq 6\eta^2 \left[L^2 \|x_t^m - \bar{x}_t\|^2 + L^2 \|\bar{x}_t - \bar{z}_t\|^2 + \frac{1}{2L} (f(\bar{z}_t) - f_*) \right]
\end{aligned} \tag{C.52}$$

674 Combine all these inequalities plugging in (C.46) and notice that $\eta \leq \frac{1}{6L}$,

$$\begin{aligned}
\|\bar{z}_{t+1} - x_*\|^2 &\leq (1 - \frac{\eta\mu}{2}) \|\bar{z}_t - x_*\|^2 - \eta(f(\bar{z}_t) - f_*) + 4\eta L [\|\bar{z}_t - \bar{x}_t\|^2 + \mathbb{E}_m[\|x_t^m - \bar{x}_t\|^2]] \\
&\quad - 2\eta \langle \bar{z}_t - x_*, \mathbb{E}_m[\widehat{g}_t^m - \nabla f(x_t^m)] \rangle + 2\eta^2 \|\mathbb{E}_m[\widehat{g}_t^m - \nabla f(x_t^m)]\|^2.
\end{aligned} \tag{C.53}$$

675 Define $\Lambda_t := \sum_{j=0}^{t-1} a_{t,j} \|\bar{x}_j - \bar{x}_{j+1}\|^2$, where $a_{t,j} := \beta_1^{t-j-1} (t-j + \frac{\beta_1}{1-\beta_1})$. By Lemma C.7,

676 we plug (C.85) in the above inequality and compute (C.53) + $\frac{2^8(\eta L)^3 \beta_1^2}{(1-\beta_1)^4} \times$ (C.84). Now let

677 $\Phi_t := \|\bar{z}_t - x_*\|^2 + \frac{2^8(\eta L)^3 \beta_1^2}{(1-\beta_1)^4} \Lambda_{t-1}$. Hence we obtain

$$\begin{aligned}
\Phi_{t+1} &\leq (1 - \frac{\eta\mu}{2}) \Phi_t - \eta(f(\bar{z}_t) - f_*) + 4\eta L \left[\mathbb{E}_m[\|x_t^m - \bar{x}_t\|^2] + 64 \left(\frac{\eta\beta_1}{1-\beta_1} \right)^2 \|\nabla f(\bar{z}_t)\|^2 \right] \\
&\quad + 32\eta L \left(\frac{\eta\beta_1}{1-\beta_1} \right)^2 \left[(1-\beta_1) \sum_{j=0}^{t-1} \beta_1^{t-j-1} \left[2L^2 \mathbb{E}_m[\|x_j^m - \bar{x}_j\|^2] + \|\mathbb{E}_m[\widehat{g}_j^m - \nabla f(x_j^m)]\|^2 \right] \right] \\
&\quad - 2\eta \langle \bar{z}_t - x_*, \mathbb{E}_m[\widehat{g}_t^m - \nabla f(x_t^m)] \rangle + 2\eta^2 \|\mathbb{E}_m[\widehat{g}_t^m - \nabla f(x_t^m)]\|^2 \\
&\leq (1 - \frac{\eta\mu}{2}) \Phi_t - \frac{\eta}{2} (f(\bar{z}_t) - f_*) + 4\eta L \mathbb{E}_m[\|x_t^m - \bar{x}_t\|^2] \\
&\quad + 32\eta L \left(\frac{\eta\beta_1}{1-\beta_1} \right)^2 \left[(1-\beta_1) \sum_{j=0}^{t-1} \beta_1^{t-j-1} \left[2L^2 \mathbb{E}_m[\|x_j^m - \bar{x}_j\|^2] + \|\mathbb{E}_m[\widehat{g}_j^m - \nabla f(x_j^m)]\|^2 \right] \right] \\
&\quad - 2\eta \langle \bar{z}_t - x_*, \mathbb{E}_m[\widehat{g}_t^m - \nabla f(x_t^m)] \rangle + 2\eta^2 \|\mathbb{E}_m[\widehat{g}_t^m - \nabla f(x_t^m)]\|^2 \\
&\leq (1 - \frac{\eta\mu}{2}) \Phi_t - \frac{\eta}{2} (f(\bar{z}_t) - f_*) + 16\eta L \cdot \eta^2 \sigma^2 K A \\
&\quad + 32\eta L \left(\frac{\eta\beta_1}{1-\beta_1} \right)^2 \left[(1-\beta_1) \sum_{j=0}^{t-1} \beta_1^{t-j-1} \|\mathbb{E}_m[\widehat{g}_j^m - \nabla f(x_j^m)]\|^2 \right] \\
&\quad - 2\eta \langle \bar{z}_t - x_*, \mathbb{E}_m[\widehat{g}_t^m - \nabla f(x_t^m)] \rangle + 2\eta^2 \|\mathbb{E}_m[\widehat{g}_t^m - \nabla f(x_t^m)]\|^2.
\end{aligned} \tag{C.54}$$

678 Here in the second inequality we use $\|\nabla f(\bar{z}_t)\|^2 \leq 2L(f(\bar{z}_t) - f_*)$. In the last inequality, we apply
 679 contraction results implied by event $E_{t,1}$.

680 Unroll this recursive bound and re-calculate the coefficients,

$$\begin{aligned} \sum_{j=0}^t \frac{\eta}{2} (f(\bar{z}_j) - f_*) (1 - \frac{\eta\mu}{2})^{t-j} + \Phi_{t+1} &\leq (1 - \frac{\eta\mu}{2})^{t+1} \Phi_0 + \frac{32\eta^2 L\sigma^2 K A}{\mu} \\ &\quad - 2\eta \sum_{j=0}^t (1 - \frac{\eta\mu}{2})^{t-j} \left\langle \bar{z}_j - x_*, \mathbb{E}_m[\widehat{g}_j^m - \nabla f(x_j^m)] \right\rangle \\ &\quad + 4\eta^2 \sum_{j=0}^t (1 - \frac{\eta\mu}{2})^{t-j} \|\mathbb{E}_m[\widehat{g}_j^m - \nabla f(x_j^m)]\|^2 \end{aligned} \tag{C.55}$$

681 Simplify Φ_{t+1} term,

$$\begin{aligned} \sum_{j=0}^t \frac{\eta}{2} (f(\bar{z}_j) - f_*) (1 - \frac{\eta\mu}{2})^{t-j} + \|\bar{z}_{t+1} - x_*\|^2 &\leq (1 - \frac{\eta\mu}{2})^{t+1} \|x_0 - x_*\|^2 + \frac{32\eta^2 L\sigma^2 K A}{\mu} \\ &\quad - 2\eta \underbrace{\sum_{j=0}^t (1 - \frac{\eta\mu}{2})^{t-j} \left\langle \bar{z}_j - x_*, \mathbb{E}_m[\widehat{g}_j^m - \mathbb{E}_j[\widehat{g}_j^m]] \right\rangle}_{\textcircled{1}: \text{martingale}} \\ &\quad - 2\eta \underbrace{\sum_{j=0}^t (1 - \frac{\eta\mu}{2})^{t-j} \left\langle \bar{z}_j - x_*, \mathbb{E}_m[\mathbb{E}_j[\widehat{g}_j^m] - \nabla f(x_j^m)] \right\rangle}_{\textcircled{2}: \text{clipping bias}} \\ &\quad + 4\eta^2 \sum_{j=0}^t (1 - \frac{\eta\mu}{2})^{t-j} \|\mathbb{E}_m[\widehat{g}_j^m - \nabla f(x_j^m)]\|^2. \end{aligned} \tag{C.56}$$

682 For the last term,

$$\begin{aligned} 4\eta^2 \sum_{j=0}^t (1 - \frac{\eta\mu}{2})^{t-j} \|\mathbb{E}_m[\widehat{g}_j^m - \nabla f(x_j^m)]\|^2 &\leq 8\eta^2 \underbrace{\sum_{j=0}^t (1 - \frac{\eta\mu}{2})^{t-j} \left[\|\mathbb{E}_m[\widehat{g}_j^m - \mathbb{E}_j[\widehat{g}_j^m]]\|^2 - \mathbb{E}_j[\|\mathbb{E}_m[\widehat{g}_j^m] - \mathbb{E}_j[\widehat{g}_j^m]\|^2] \right]}_{\textcircled{3}: \text{martingale}} \\ &\quad + 8\eta^2 \underbrace{\sum_{j=0}^t (1 - \frac{\eta\mu}{2})^{t-j} \mathbb{E}_j[\|\mathbb{E}_m[\widehat{g}_j^m] - \mathbb{E}_j[\widehat{g}_j^m]\|^2]}_{\text{Lemma B.2}} \\ &\quad + 8\eta^2 \underbrace{\sum_{j=0}^t (1 - \frac{\eta\mu}{2})^{t-j} \|\mathbb{E}_m[\mathbb{E}_j[\widehat{g}_j^m] - \nabla f(x_j^m)]\|^2}_{\textcircled{4}: \text{clipping bias}}, \end{aligned} \tag{C.57}$$

683 we finally get

$$\begin{aligned} \sum_{j=0}^t \frac{\eta}{2} (f(\bar{z}_j) - f_*) (1 - \frac{\eta\mu}{2})^{t-j} + \|\bar{z}_{t+1} - x_*\|^2 &\leq (1 - \frac{\eta\mu}{2})^{t+1} D_0^2 + 32 \left[\eta L K A + \frac{1}{M} \right] \frac{\eta\sigma^2}{\mu} \\ &\quad + \textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4}. \end{aligned} \tag{C.58}$$

684 (1) **Case** $\mu > 0$.

685 For ①, define

$$\zeta_j = \begin{cases} -2\eta(1 - \frac{\eta\mu}{2})^{t-j} \left\langle \bar{z}_j - x_*, \mathbb{E}_m[\widehat{g}_j^m] - \mathbb{E}_j[\widehat{g}_j^m] \right\rangle, & \text{if event } E_j \text{ holds,} \\ 0, & \text{otherwise.} \end{cases} \quad (\text{C.59})$$

686 Then since event E_j implies $\|\bar{z}_j - x_*\| \leq \sqrt{2}(1 - \frac{\eta\mu}{2})^{j/2}D_0$,

$$|\zeta_j| \leq 2\eta \cdot \sqrt{2}(1 - \frac{\eta\mu}{2})^{t/2}D_0 \cdot 2\rho\sqrt{d} = 4(1 - \frac{\eta\mu}{2})^{t/2}\eta\rho\sqrt{2d}D_0 \stackrel{\text{def}}{=} c, \quad (\text{C.60})$$

687

$$\text{Var}_j(\zeta_j) \leq 4\eta^2(1 - \frac{\eta\mu}{2})^{2(t-j)} \cdot 2(1 - \frac{\eta\mu}{2})^j D_0^2 \cdot \frac{\sigma^2}{M} = 8(1 - \frac{\eta\mu}{2})^{2t-j} \frac{\eta^2 D_0^2 \sigma^2}{M}. \quad (\text{C.61})$$

688 Let $b = \frac{(1 - \frac{\eta\mu}{2})^{t+1}D_0^2}{5}$, $V = 16(1 - \frac{\eta\mu}{2})^t \frac{\eta D_0^2 \sigma^2}{\mu M}$. By Lemma B.1, $|\sum_{j=0}^t \zeta_j| \leq b$ with probability

689 no less than

$$1 - 2 \exp\left(\frac{b^2}{2V + 2cb/3}\right) \geq 1 - \frac{\delta}{4T}. \quad (\text{C.62})$$

690 For ②, since by Lemma B.2,

$$\|\mathbb{E}_j[\widehat{g}_j^m] - \nabla f(x_j^m)\|^2 \leq \frac{\|2\sigma\|_{2\alpha}^{2\alpha}}{\rho^{2(\alpha-1)}}, \quad (\text{C.63})$$

691 event E_t implies that

$$\begin{aligned} |\textcircled{2}| &\leq 2\eta \sum_{j=0}^t (1 - \frac{\eta\mu}{2})^{t-j} \cdot \sqrt{2}(1 - \frac{\eta\mu}{2})^{j/2} D_0 \cdot \frac{\|2\sigma\|_{2\alpha}^\alpha}{\rho^{\alpha-1}} \\ &\leq 4\sqrt{2}(1 - \frac{\eta\mu}{2})^{t/2} \frac{D_0 \|2\sigma\|_{2\alpha}^\alpha}{\mu\rho^{\alpha-1}} \\ &\leq \frac{(1 - \frac{\eta\mu}{2})^{t+1}D_0^2}{5}. \end{aligned} \quad (\text{C.64})$$

692 Here we use the definition of η and conditions of ρ in (C.12).

693 For ③, define

$$\theta_j = \begin{cases} 8\eta^2(1 - \frac{\eta\mu}{2})^{t-j} \left[\|\mathbb{E}_m[\widehat{g}_j^m] - \mathbb{E}_j[\widehat{g}_j^m]\|^2 - \mathbb{E}_j[\|\mathbb{E}_m[\widehat{g}_j^m] - \mathbb{E}_j[\widehat{g}_j^m]\|^2] \right], & \text{if event } E_j \text{ holds,} \\ 0, & \text{otherwise.} \end{cases} \quad (\text{C.65})$$

694 Then

$$|\theta_j| \leq 8\eta^2 \cdot 4\rho^2 d = 32\eta^2 \rho^2 d \stackrel{\text{def}}{=} c, \quad (\text{C.66})$$

695

$$\text{Var}_j(\theta_j) \leq 64\eta^4(1 - \frac{\eta\mu}{2})^{2(t-j)} \cdot \mathbb{E}_j[\|\mathbb{E}_m[\widehat{g}_j^m] - \mathbb{E}_j[\widehat{g}_j^m]\|^2]^2 \stackrel{\text{Lemma B.3}}{\leq} 64\eta^4(1 - \frac{\eta\mu}{2})^{2(t-j)} \cdot \frac{4(2\sigma)^4}{M^2}. \quad (\text{C.67})$$

696 Let $b = \frac{(1 - \frac{\eta\mu}{2})^{t+1}D_0^2}{5}$, $V = \frac{2^{13}\eta^3\sigma^4}{\mu M^2}$. By Lemma B.1, $|\sum_{j=0}^t \theta_j| \leq b$ with probability no less than

$$1 - 2 \exp\left(\frac{b^2}{2V + 2cb/3}\right) \geq 1 - \frac{\delta}{4T}. \quad (\text{C.68})$$

697 For ④, by Lemma B.2,

$$|\textcircled{4}| \leq \frac{16\eta}{\mu} \cdot \frac{\|2\sigma\|_{2\alpha}^{2\alpha}}{\rho^{2(\alpha-1)}} \leq \frac{(1 - \frac{\eta\mu}{2})^{t+1}D_0^2}{5}. \quad (\text{C.69})$$

698 Combine the above claims, with probability no less than $\mathbb{P}(E_{t,1}) - 2 \cdot \frac{\delta}{4T}$, we have $|\textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4}| \leq$
699 $\frac{4}{5}(1 - \frac{\eta\mu}{2})^{t+1}D_0^2$. By (C.58), these implies

$$\begin{aligned} \sum_{j=0}^t \frac{\eta}{2} (f(\bar{z}_j) - f_*) (1 - \frac{\eta\mu}{2})^{t-j} + \|\bar{z}_{t+1} - x_*\|^2 &\leq (1 - \frac{\eta\mu}{2})^{t+1} D_0^2 + 32 \left[\eta LKA + \frac{1}{M} \right] \frac{\eta\sigma^2}{\mu} \\ &\quad + \frac{4}{5} (1 - \frac{\eta\mu}{2})^{t+1} D_0^2 \\ &\leq 2(1 - \frac{\eta\mu}{2})^{t+1} D_0^2. \end{aligned} \tag{C.70}$$

700 Therefore, we conclude that $\mathbb{P}(E_{t+1}) \geq \mathbb{P}(E_{t,1}) - \frac{\delta}{2T}$.

701 (2) **Case** $\mu = 0$.

702 In this case, (C.58) reduces to

$$\frac{\eta}{2} \sum_{j=0}^t (f(\bar{z}_j) - f_*) + \|\bar{z}_{t+1} - x_*\|^2 \leq D_0^2 + 16 \left[\eta LKA + \frac{1}{M} \right] \eta^2 \sigma^2 (t+1) + \textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4}. \tag{C.71}$$

703 For $\textcircled{1}$, define

$$\zeta_j = \begin{cases} -2\eta \langle \bar{z}_j - x_*, \mathbb{E}_m[\widehat{g}_j^m] - \mathbb{E}_j[\widehat{g}_j^m] \rangle, & \text{if event } E_j \text{ holds,} \\ 0, & \text{otherwise.} \end{cases} \tag{C.72}$$

704 Then since event E_j implies $\|\bar{z}_j - x_*\| \leq \sqrt{2}D_0$,

$$|\zeta_j| \leq 2\eta \cdot \sqrt{2}D_0 \cdot 2\rho\sqrt{d} = 4\eta\rho\sqrt{2d}D_0 \stackrel{\text{def}}{=} c, \tag{C.73}$$

705

$$\text{Var}_j(\zeta_j) \leq 4\eta^2 \cdot 2D_0^2 \cdot \frac{\sigma^2}{M} = \frac{8\eta^2 D_0^2 \sigma^2}{M}. \tag{C.74}$$

706 Let $b = \frac{D_0^2}{5}$, $V = \frac{8\eta^2 D_0^2 \sigma^2 T}{M}$. By Lemma B.1, $|\sum_{j=0}^t \zeta_j| \leq b$ with probability no less than

$$1 - 2 \exp\left(-\frac{b^2}{2V + 2cb/3}\right) \geq 1 - \frac{\delta}{4T}. \tag{C.75}$$

707 For $\textcircled{2}$, since by Lemma B.2,

$$\|\mathbb{E}_j[\widehat{g}_j^m] - \nabla f(x_j^m)\|^2 \leq \frac{\|2\sigma\|_{2\alpha}^{2\alpha}}{\rho^{2(\alpha-1)}}, \tag{C.76}$$

708 event E_t implies that

$$|\textcircled{2}| \leq 2\eta(t+1) \cdot \sqrt{2}D_0 \cdot \frac{\|2\sigma\|_{2\alpha}^\alpha}{\rho^{(\alpha-1)}} \leq \frac{D_0^2}{5}. \tag{C.77}$$

709 Here we again use definitions and conditions in (C.12).

710 For $\textcircled{3}$, define

$$\theta_j = \begin{cases} 8\eta^2 \left[\|\mathbb{E}_m[\widehat{g}_j^m] - \mathbb{E}_j[\widehat{g}_j^m]\|^2 - \mathbb{E}_j[\|\mathbb{E}_m[\widehat{g}_j^m] - \mathbb{E}_j[\widehat{g}_j^m]\|^2] \right], & \text{if event } E_j \text{ holds,} \\ 0, & \text{otherwise.} \end{cases} \tag{C.78}$$

711 Then

$$|\theta_j| \leq 8\eta^2 \cdot 4\rho^2 d = 32\eta^2 \rho^2 d \stackrel{\text{def}}{=} c, \tag{C.79}$$

712

$$\text{Var}_j(\theta_j) \leq 64\eta^4 \cdot \mathbb{E}_j[\|\mathbb{E}_m[\widehat{g}_j^m] - \mathbb{E}_j[\widehat{g}_j^m]\|^2]^2 \stackrel{\text{Lemma B.3}}{\leq} 64\eta^4 \cdot \frac{4(2\sigma)^4}{M^2}. \tag{C.80}$$

713 Let $b = \frac{D_0^2}{5}$, $V = \frac{2^{12}\eta^4\sigma^4}{M^2}$. By Lemma B.1, $|\sum_{j=0}^t \theta_j| \leq b$ with probability no less than

$$1 - 2 \exp\left(\frac{b^2}{2V + 2cb/3}\right) \geq 1 - \frac{\delta}{4T}. \quad (\text{C.81})$$

714 For ④, by Lemma B.2,

$$|\textcircled{4}| \leq 8\eta^2(t+1) \cdot \frac{\|2\sigma\|_{2\alpha}^{2\alpha}}{\rho^{2(\alpha-1)}} \leq \frac{D_0^2}{5}. \quad (\text{C.82})$$

715 Combine the above claims, with probability no less than $\mathbb{P}(E_{t,1}) - 2 \cdot \frac{\delta}{4T}$, we have $|\textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4}| \leq$
 716 $\frac{4}{5}D_0^2$. By (C.58), these implies

$$\begin{aligned} \frac{\eta}{2} \sum_{j=0}^t (f(\bar{z}_j) - f_*) + \|\bar{z}_{t+1} - x_*\|^2 &\leq D_0^2 + 16 \left[\eta LKA + \frac{1}{M} \right] \eta^2 \sigma^2 (t+1) + \frac{4}{5} D_0^2 \\ &\leq 2D_0^2. \end{aligned} \quad (\text{C.83})$$

717 Therefore, we conclude that $\mathbb{P}(E_{t+1}) \geq \mathbb{P}(E_{t,1}) - \frac{\delta}{2T}$.

718 □

719 **Lemma C.7.** Let $\Lambda_t := \sum_{j=0}^{t-1} a_{t,j} \|\bar{x}_j - \bar{x}_{j+1}\|^2$, where $a_{t,j} := \beta_1^{t-j-1} (t-j + \frac{\beta_1}{1-\beta_1})$. Under the
 720 conditions in Lemma C.2, then the following holds:

$$\begin{aligned} \Lambda_t &\leq \left(1 - \frac{(1-\beta_1)^2}{2}\right) \Lambda_{t-1} + \frac{32\eta^2}{1-\beta_1} \|\nabla f(\bar{z}_t)\|^2 \\ &\quad + 4\eta^2 \sum_{j=0}^{t-1} \beta_1^{t-j-1} \left[2L^2 \mathbb{E}_m[\|x_j^m - \bar{x}_j\|^2] + \|\mathbb{E}_m[\widehat{g}_j^m - \nabla f(x_j^m)]\|^2 \right]. \end{aligned} \quad (\text{C.84})$$

721

$$\begin{aligned} \|\bar{z}_t - \bar{x}_t\|^2 &\leq \left(\frac{\eta\beta_1}{1-\beta_1}\right)^2 [16L^2 \Lambda_{t-1} + 32\|\nabla f(\bar{z}_t)\|^2] \\ &\quad + \frac{4(\eta\beta_1)^2}{1-\beta_1} \sum_{j=0}^{t-1} \beta_1^{t-j-1} \left[2L^2 \mathbb{E}_m[\|x_j^m - \bar{x}_j\|^2] + \|\mathbb{E}_m[\widehat{g}_j^m - \nabla f(x_j^m)]\|^2 \right]. \end{aligned} \quad (\text{C.85})$$

722 *Proof.* By definition, $\|\bar{z}_t - \bar{x}_t\|^2 = \left(\frac{\beta_1}{1-\beta_1}\right)^2 \|\bar{x}_t - \bar{x}_{t-1}\|^2$ and

$$\begin{aligned} \|\bar{x}_t - \bar{x}_{t-1}\|^2 &= \eta^2 \|\bar{u}_{t-1}\|^2 \\ &= \eta^2 \left\| (1-\beta_1) \sum_{j=0}^{t-1} \beta_1^{t-j-1} \mathbb{E}_m[\widehat{g}_j^m] \right\|^2 \\ &\leq 2\eta^2 \left[\left\| (1-\beta_1) \sum_{j=0}^{t-1} \beta_1^{t-j-1} \mathbb{E}_m[\nabla f(x_j^m)] \right\|^2 + \left\| (1-\beta_1) \sum_{j=0}^{t-1} \beta_1^{t-j-1} \mathbb{E}_m[\widehat{g}_j^m - \nabla f(x_j^m)] \right\|^2 \right] \\ &\leq 4\eta^2 \left\| (1-\beta_1) \sum_{j=0}^{t-1} \beta_1^{t-j-1} \nabla f(\bar{x}_j) \right\|^2 \\ &\quad + 2\eta^2 (1-\beta_1) \sum_{j=0}^{t-1} \beta_1^{t-j-1} \left[2L^2 \mathbb{E}_m[\|x_j^m - \bar{x}_j\|^2] + \|\mathbb{E}_m[\widehat{g}_j^m - \nabla f(x_j^m)]\|^2 \right]. \end{aligned} \quad (\text{C.86})$$

723 Note that

$$\begin{aligned}
\left\| (1 - \beta_1) \sum_{j=0}^{t-1} \beta_1^{t-j-1} \nabla f(\bar{x}_j) \right\|^2 &\leq 2 \left\| (1 - \beta_1) \sum_{j=0}^{t-1} \beta_1^{t-j-1} [\nabla f(\bar{x}_j) - \nabla f(\bar{x}_t)] \right\|^2 + 2 \|\nabla f(\bar{x}_t)\|^2 \\
&\leq 2(1 - \beta_1) \sum_{j=0}^{t-1} \beta_1^{t-j-1} L^2 \|\bar{x}_j - \bar{x}_t\|^2 + 2 \|\nabla f(\bar{x}_t)\|^2 \\
&\leq 2(1 - \beta_1) \sum_{j=0}^{t-1} \beta_1^{t-j-1} L^2 \cdot (t - j) \sum_{i=j}^{t-1} [\|\bar{x}_i - \bar{x}_{i+1}\|^2] + 2 \|\nabla f(\bar{x}_t)\|^2 \\
&\leq 2L^2 \sum_{j=0}^{t-1} a_{t,j} \|\bar{x}_j - \bar{x}_{j+1}\|^2 + 4 \|\nabla f(\bar{z}_t)\|^2 + 4L^2 \|\bar{x}_t - \bar{z}_t\|^2 \\
&\leq 2L^2 \sum_{j=0}^{t-2} a_{t-1,j} \|\bar{x}_j - \bar{x}_{j+1}\|^2 + 4 \|\nabla f(\bar{z}_t)\|^2 + \frac{4L^2}{(1 - \beta_1)^2} \|\bar{x}_t - \bar{x}_{t-1}\|^2
\end{aligned} \tag{C.87}$$

724 Here $a_{t,j} = \beta_1^{t-j-1} (t - j + \frac{\beta_1}{1 - \beta_1})$. For $j \leq t - 2$, we have $a_{t,j} \leq \beta_1(2 - \beta_1)a_{t-1,j}$. Since

725 $\Lambda_t = \sum_{j=0}^{t-1} a_{t,j} \|\bar{x}_j - \bar{x}_{j+1}\|^2$, we can conclude that

$$\begin{aligned}
\|\bar{x}_t - \bar{x}_{t-1}\|^2 &\leq 16\eta^2 L^2 \Lambda_{t-1} + 32\eta^2 \|\nabla f(\bar{z}_t)\|^2 \\
&\quad + 4\eta^2 (1 - \beta_1) \sum_{j=0}^{t-1} \beta_1^{t-j-1} \left[2L^2 \mathbb{E}_m [\|x_j^m - \bar{x}_j\|^2] + \|\mathbb{E}_m[\widehat{g}_j^m - \nabla f(x_j^m)]\|^2 \right],
\end{aligned} \tag{C.88}$$

726 which implies (C.85). We complete the proof by plugging the above inequality in

$$\Lambda_t \leq \beta_1(2 - \beta_1)\Lambda_{t-1} + \frac{1}{1 - \beta_1} \|\bar{x}_t - \bar{x}_{t-1}\|^2. \tag{C.89}$$

727 □

728 C.5 Further Discussion

729 **Coordinate-wise clipping and global clipping.** Lemma B.2 can be easily extended to \mathbb{R}^d , similar
730 to Sadiev et al. [2023, Lemma 5.1]. Therefore, our results can be easily generalized to global
731 clipping operator $\mathbf{clip}_g(X, \rho_g) := \min \left\{ 1, \frac{\rho_g}{\|X\|} \right\} X$ with threshold $\rho_g := \rho\sqrt{d}$. We omit the
732 details in this paper. Readers may also wonder why our Theorem C.4 and Theorem C.5 depend on
733 $\mathbf{poly}(d)$. However, if we assume $\|\sigma\|_{2\alpha} d^{\frac{1}{2} - \frac{1}{2\alpha}} = \mathcal{O}(\sigma)$, both of which are of order $\mathcal{O}(d^{\frac{1}{2}})$, then
734 our convergence guarantee will not depend on $\mathbf{poly}(d)$ explicitly. Zhang et al. [2020, Corollary
735 7] claims that coordinate-wise clipping has better dependence on dimension d . But they simply
736 upper bound $\mathbb{E}_{\xi \sim \mathcal{D}} \|\nabla F(x, \xi)\|^\alpha$ by $d^{\alpha/2} \mathbb{E}_{\xi \sim \mathcal{D}} \|\nabla F(x, \xi)\|_\alpha^\alpha$, which is too pessimistic. In fact, if we
737 assume $\mathbb{E}_{\xi \sim \mathcal{D}} \|\nabla F(x, \xi)\|^\alpha = \mathcal{O}(d^{\alpha/2-1} \mathbb{E}_{\xi \sim \mathcal{D}} \|\nabla F(x, \xi)\|_\alpha^\alpha)$, both of which are of order $\mathcal{O}(d^{\frac{\alpha}{2}})$,
738 then there is still no difference between coordinate-wise clipping and global clipping in their setting.

739 **Prior works on distributed SGDM with local updates.** There are many works on *Local* SGDM in
740 distributed setting. Liu et al. [2020a] studies *Local* SGDM in convex setting and rely on some strong
741 assumptions to show convergence. Xu et al. [2021] analyze *Local* SGDM with bounded gradient
742 assumption and the use a global momentum parameter during local iterations. Yu et al. [2019]
743 considers non-convex *Local* SGDM but is only able to prove linear speedup. Wang et al. [2019],
744 Cheng et al. [2023] also study non-convex problem and use momentum to handle heterogeneity in
745 federated learning. All these works fail to show the benefits of local iterations compared to minibatch
746 baseline.

747 **D Proof of Local Adam**

748 **D.1 Overview and Main Theorem**

749 For any integer $0 \leq t \leq T-1$, we define $r(t), k(t) \in \mathbb{N}$ such that $t = r(t)K + k(t)$ and $k(t) \leq K-1$.
 750 We omit the dependence on t and let $r = r(t), k = k(t)$ through out the proof if not causing confusion.

751 Define $x_t^m := x_{r,k}^m, g_t^m := g_{r,k}^m, \widehat{g}_t^m := \widehat{g}_{r,k}^m, u_t^m = u_{r,k}^m$. Then Algorithm 2 is equivalent to the
 752 following update rule:

$$u_t^m = \begin{cases} \beta_1 u_{t-1}^m + (1 - \beta_1) \widehat{g}_t^m & \text{if } t \bmod K \neq 0, \\ \beta_1 \bar{u}_{t-1} + (1 - \beta_1) \widehat{g}_t^m & \text{otherwise,} \end{cases} \quad (\text{D.1})$$

753

$$v_t^m = \begin{cases} \beta_2 v_{t-1}^m + (1 - \beta_2) \widehat{g}_t^m{}^2 & \text{if } t \bmod K \neq 0, \\ \beta_2 \bar{v}_{t-1} + (1 - \beta_2) \widehat{g}_t^m{}^2 & \text{otherwise,} \end{cases} \quad (\text{D.2})$$

754

$$x_{t+1}^m = \begin{cases} x_t^m - \eta(H_t^m)^{-1} u_t^m & \text{if } t \bmod K \neq -1, \\ \bar{x}_t - \eta \mathbb{E}_m[(H_t^m)^{-1} u_t^m] & \text{otherwise.} \end{cases} \quad (\text{D.3})$$

755 Define an auxiliary sequence $\{z_t^m\}$ as:

$$z_{t+1}^m = \begin{cases} \frac{1}{1 - \beta_1} x_{t+1}^m - \frac{\beta_1}{1 - \beta_1} x_t^m & \text{if } t \bmod K \neq -1, \\ \frac{1}{1 - \beta_1} x_{t+1}^m - \frac{\beta_1}{1 - \beta_1} \bar{x}_t & \text{otherwise.} \end{cases} \quad (\text{D.4})$$

756 Let

$$e_t^m := \frac{\beta_1}{1 - \beta_1} (I_d - H_t^m (H_{t-1}^m)^{-1}) u_{t-1}^m. \quad (\text{D.5})$$

757 Then the definition of $\{z_t^m\}$ implies

$$\begin{aligned} z_{t+1}^m - z_t^m &= -\frac{\eta(H_t^m)^{-1} u_t^m}{1 - \beta_1} + \frac{\eta\beta_1(H_{t-1}^m)^{-1} u_{t-1}^m}{1 - \beta_1} \\ &= -\frac{\eta\beta_1}{1 - \beta_1} [(H_t^m)^{-1} - (H_{t-1}^m)^{-1}] u_{t-1}^m - \eta(H_t^m)^{-1} \widehat{g}_t^m \\ &=: -\eta(H_t^m)^{-1} (\widehat{g}_t^m + e_t^m). \end{aligned} \quad (\text{D.6})$$

758 Finally, let $y_t := \arg \min_y f(y) + \frac{1}{2\gamma} \|y - \bar{z}_t\|_{H_{r(t)}}^2$.

759 Define probabilistic events (see (D.15) for definition of some parameters)

$$\mathcal{A}_{t,1} := \left\{ \beta_2^{K/2} \preceq H_{r(t)}^{-1} H_t^m \preceq 1 + (1 - \beta_2)B \text{ and for all } m \in [M] \right\}, \quad (\text{D.7})$$

760

$$\mathcal{A}_{t,2} := \left\{ \|H_{r(t)}((H_t^m)^{-1} - (H_t^n)^{-1})\| \leq (1 - \beta_2)B_1 \text{ for all } m, n \in [M] \right\}, \quad (\text{D.8})$$

761

$$\mathcal{A}_{t,3} := \left\{ \|z_{t+1}^m - z_{t+1}^n\|_{H_r}^2 \leq \frac{\eta^2 \sigma^2}{\lambda} K A, \sum_{j=rK}^t \|\widehat{g}_j^m\|^2 \leq \frac{(1 - \beta_1)^2 \sigma^2 A}{2^{12} (1 - \beta_2)^2 B_1^2} \text{ for all } m, n \in [M] \right\}, \quad (\text{D.9})$$

762

$$\mathcal{A}_{t,4} := \left\{ f_\gamma^{H_{r(t+1)}}(\bar{z}_{t+1}) - \min f_\gamma^\lambda + \frac{\eta}{12} \sum_{j=0}^t \|\nabla f_\gamma^{H_{r(j)}}(\bar{z}_j)\|_{H_{r(j)}^{-1}}^2 \leq 2\Delta \right\}. \quad (\text{D.10})$$

763 Here $\Delta := f_\gamma^\lambda(x_0) - \min f_\gamma^\lambda$. Besides, let

$$E_t := \{\mathcal{A}_{j,i} \text{ holds for all } j \leq t-1, i \in \{1, 2, 3, 4\}\}, \quad (\text{D.11})$$

764

$$E_{t,1} := E_t \cap \mathcal{A}_{t,1}, E_{t,2} := E_{t,1} \cap \mathcal{A}_{t,2}, E_{t,3} := E_{t,2} \cap \mathcal{A}_{t,3}. \quad (\text{D.12})$$

765 **Theorem D.1.** For $L/\lambda \geq \gamma^{-1} \geq 2\tau/\lambda$, let Assumption 1, 2, 3, 5 hold for $\Omega = \text{conv}(\mathbf{B}_{R_0}(\Omega_0))$,
766 where $\Omega_0 := \{f_\gamma^\lambda(x) - \min f_\gamma^\lambda \leq 2\Delta\}$, $\Delta = f_\gamma^\lambda(x_0) - \min f_\gamma^\lambda$ and $R_0 = \sqrt{\frac{\Delta\gamma}{160\lambda}}$. Further assume
767 that for any $x \in \Omega$, $\|\nabla f(x)\| \leq G$, $\|\nabla f(x)\|_\infty \leq G_\infty$, and

$$1 - \beta_2 \lesssim \min \left\{ \frac{1 - \beta_1}{K^{1/2}B_1} \frac{(1 - \beta_1)\sigma\sqrt{A}}{K^{1/2}B_1G}, \frac{\eta}{\gamma B}, \frac{1 - \beta_1}{K^{1/2}B}, \frac{1}{K} \right\}. \quad (\text{D.13})$$

768 If $\eta = \frac{24\lambda\Delta}{\varepsilon T}$, then with probability no less than $1 - \delta$, Local Adam yields

$$769 \frac{\lambda}{KR} \sum_{r=0}^{R-1} \sum_{k=0}^{K-1} \|\nabla f_\gamma^{H_r}(\bar{z}_{r,k})\|_{H_r^{-1}}^2 \leq \varepsilon \text{ if}$$

$$T \gtrsim \frac{\lambda\Delta\sigma^2}{\gamma M\varepsilon^2} \log^{\frac{1}{2}} \frac{T}{\delta} + \frac{\Delta}{\varepsilon} \cdot \sqrt{\frac{L^2\sigma^2KA}{\min\{\varepsilon, \sigma_\infty^2/G_\infty\}}} + \frac{L\Delta}{(1 - \beta_1)^2\varepsilon} + \frac{K\tau\Delta}{\varepsilon} + \frac{\sqrt{L\Delta\rho^2d \log \frac{T}{\delta}}}{(\sqrt{\beta_2} - \beta_1)\varepsilon}. \quad (\text{D.14})$$

770 Here

$$\begin{aligned} \rho &\geq \max \left\{ \left(\frac{2^6 \|2\sigma\|_{2\alpha}^{2\alpha}}{\varepsilon} \right)^{\frac{1}{2(\alpha-1)}}, 3\sigma_\infty, 2G_\infty \right\}, \\ B &:= \max \left\{ \frac{6K(G_\infty^2 + \sigma_\infty^2)}{\lambda^2}, \frac{16\rho^2}{\lambda^2} \log \frac{dMT}{\delta}, 2^6 \frac{\sqrt{K}(G_\infty + \sigma_\infty)\sigma_\infty}{\lambda^2} \log^{1/2} \frac{dMT}{\delta} \right\}, \\ B_1 &:= \max \left\{ \frac{16K\sigma_\infty^2}{\lambda^2}, \frac{16\rho^2}{\lambda^2} \log \frac{dMT}{\delta}, 2^6 \frac{\sqrt{K}(G_\infty + \sigma_\infty)\sigma_\infty}{\lambda^2} \log^{1/2} \frac{dMT}{\delta} \right\}, \\ A &:= \max \left\{ \frac{2^{20}\rho^2d}{K\sigma^2} \log \frac{MT}{\delta}, 2^{20} \log^2 \frac{MT}{\delta}, \frac{2^8 K \|2\sigma\|_{2\alpha}^{2\alpha}}{\sigma^2\rho^{2(\alpha-1)}} \right\}. \end{aligned} \quad (\text{D.15})$$

771 *Proof.* We prove by induction that $\mathbb{P}(E_t) \geq 1 - \frac{t\delta}{T}$ for $t = 0, \dots, T$.

772 When $t = 0$, this is trivial. Assume that the statement is true for some $t \leq T - 1$. We aim to prove
773 that $\mathbb{P}(E_{t+1}) \geq 1 - \frac{(t+1)\delta}{T}$. By Lemma D.8, D.9, D.10, D.11, we have

$$\mathbb{P}(E_{t+1}) \geq \mathbb{P}(E_t) - 4 \cdot \frac{\delta}{4T} \geq 1 - \frac{(t+1)\delta}{T}. \quad (\text{D.16})$$

774 Therefore by induction rule, $\mathbb{P}(E_T) \geq 1 - \delta$ and this implies

$$\frac{\lambda}{T} \sum_{t=0}^{T-1} \|\nabla f_\gamma^{H_{r(t)}}(\bar{z}_t)\|_{H_{r(t)}^{-1}}^2 \leq \frac{24\lambda\Delta}{\eta T} = \varepsilon. \quad (\text{D.17})$$

775 Now we verify the conditions in all the lemmas. In Lemma D.7,

$$\frac{\eta}{\lambda} \lesssim \sqrt{\frac{\Delta\gamma}{\lambda\sigma^2KA}} \iff T \gtrsim \frac{\sigma}{\varepsilon} \sqrt{L\Delta KA}. \quad (\text{D.18})$$

776 In Lemma D.9,

$$\frac{\eta}{\lambda} \lesssim \frac{\sigma_\infty^2}{G_\infty L\sigma\sqrt{KA}} \iff T \gtrsim \frac{\Delta}{\varepsilon} \cdot \sqrt{\frac{L^2\sigma^2KA}{\sigma_\infty^2/G_\infty}}. \quad (\text{D.19})$$

777 In Lemma D.10,

$$\frac{\eta}{\lambda} \lesssim \min \left\{ \frac{1}{K\tau}, \frac{(1 - \beta_1)^2}{L} \right\} \iff T \gtrsim \frac{L\Delta}{(1 - \beta_1)^2\varepsilon} + \frac{K\tau\Delta}{\varepsilon}. \quad (\text{D.20})$$

778 In Lemma D.11, by noticing that $\frac{24\Delta\lambda}{\eta T} = \varepsilon$, (D.113) is equivalent to $\rho \gtrsim \left(\frac{\|2\sigma\|_{2\alpha}^{2\alpha}}{\varepsilon}\right)^{\frac{1}{2(\alpha-1)}}$ and

$$\frac{\eta}{\lambda} \lesssim \min \left\{ \frac{(1-\beta_1)^2}{L}, \frac{M\gamma\varepsilon}{\lambda\sigma^2 \log^{1/2} \frac{T}{\delta}}, \left(\frac{L^2\sigma^2 KA}{\varepsilon}\right)^{-1/2}, \frac{M\Delta}{\sigma^2 \log \frac{T}{\delta}}, \sqrt{\frac{\gamma\Delta}{\lambda\rho^2 d \log \frac{T}{\delta}}}, \frac{\sqrt{T\varepsilon}(\sqrt{\beta_2}-\beta_1)}{L\rho\sqrt{d} \log^{1/2} \frac{T}{\delta}} \right\}, \quad (\text{D.21})$$

779 which can be ensured as long as

$$T \gtrsim \max \left\{ \frac{L\Delta}{(1-\beta_1)^2\varepsilon}, \frac{\lambda\Delta\sigma^2}{\gamma M\varepsilon^2} \log^{\frac{1}{2}} \frac{T}{\delta}, \frac{\Delta}{\varepsilon} \cdot \sqrt{\frac{L^2\sigma^2 KA}{\varepsilon}}, \frac{\sqrt{L\Delta\rho^2 d \log \frac{T}{\delta}}}{(\sqrt{\beta_2}-\beta_1)\varepsilon} \right\}. \quad (\text{D.22})$$

780 Here we use the fact that $\gamma \geq \frac{\lambda}{L}$. Therefore we can conclude that all the lemmas hold if

$$T \gtrsim \frac{\lambda\Delta\sigma^2}{\gamma M\varepsilon^2} \log^{\frac{1}{2}} \frac{T}{\delta} + \frac{\Delta}{\varepsilon} \cdot \sqrt{\frac{L^2\sigma^2 KA}{\min\{\varepsilon, \sigma_\infty^2/G_\infty\}}} + \frac{L\Delta}{(1-\beta_1)^2\varepsilon} + \frac{K\tau\Delta}{\varepsilon} + \frac{\sqrt{L\Delta\rho^2 d \log \frac{T}{\delta}}}{\varepsilon}. \quad (\text{D.23})$$

781 Finally, we verify the upper bound of $1-\beta_2$ in Lemma D.9, D.10 and D.11 as:

$$1-\beta_2 \lesssim \min \left\{ \frac{1-\beta_1}{K^{1/2}B_1}, \frac{(1-\beta_1)\sigma\sqrt{A}}{K^{1/2}B_1G}, \frac{\eta}{\gamma B}, \frac{1-\beta_1}{K^{1/2}B}, \frac{1}{K} \right\}. \quad (\text{D.24})$$

782

□

783 **Theorem D.2.** Under the conditions of Theorem D.1, assume $1-\beta_1 = \Omega(1)$ and

$$1-\beta_2 = \tilde{\mathcal{O}} \left(\frac{1}{K^{3/2}R^{1/2}} \right), \quad \left(\frac{\|\sigma\|_{2\alpha}^{2\alpha}}{\varepsilon} \right)^{\frac{1}{2(\alpha-1)}} \gtrsim G_\infty \vee \sigma_\infty, \varepsilon \lesssim \frac{\sigma_\infty^2}{G_\infty}, \quad (\text{D.25})$$

$$K \gtrsim \log \frac{MT}{\delta} \left(\frac{\|\sigma\|_{2\alpha} d^{\frac{1}{2}-\frac{1}{2\alpha}}}{\sigma} \right)^{\frac{2\alpha}{\alpha-2}}.$$

784 Then with probability no less than $1-\delta$, Local Adam with optimal η, ρ yields

785 $\frac{\lambda}{KR} \sum_{r=0}^{R-1} \sum_{k=0}^{K-1} \|\nabla f_\gamma^{H_r}(\bar{z}_{r,k})\|_{H_r^{-1}}^2 \leq \varepsilon$ if

$$T \gtrsim \frac{\lambda\Delta\sigma^2}{\gamma M\varepsilon^2} \log^{\frac{1}{2}} \frac{T}{\delta} + \frac{L\Delta}{\varepsilon^{\frac{3}{2}}} \cdot \sqrt{\sigma^2 K \log \frac{MT}{\delta}} + \frac{(L+K\tau)\Delta}{\varepsilon} + \frac{L\Delta}{\varepsilon^{\frac{3}{2}}} \left(\frac{\|\sigma\|_{2\alpha}^{2\alpha}}{\varepsilon} \right)^{\frac{1}{2(\alpha-1)}} d^{\frac{1}{2}} \log \frac{MT}{\delta}. \quad (\text{D.26})$$

786 And equivalently,

$$\begin{aligned} \frac{\lambda}{KR} \sum_{r=0}^{R-1} \sum_{k=0}^{K-1} \|\nabla f_\gamma^{H_r}(\bar{z}_{r,k})\|_{H_r^{-1}}^2 &\lesssim \frac{\tau\Delta}{R} + \frac{L\Delta}{KR} + \sqrt{\frac{\lambda\Delta\sigma^2}{\gamma MKR}} \log^{\frac{1}{4}} \frac{KR}{\delta} \\ &+ \frac{(L\Delta\sigma)^{\frac{2}{3}}}{K^{\frac{1}{3}}R^{\frac{2}{3}}} \log^{\frac{1}{3}} \frac{MKR}{\delta} + \left(\|\sigma\|_{2\alpha} d^{\frac{1}{2}-\frac{1}{2\alpha}} \right)^{\frac{2\alpha}{3\alpha-2}} \left(\frac{L\Delta \log \frac{MKR}{\delta}}{KR} \right)^{\frac{2(\alpha-1)}{3\alpha-2}}. \end{aligned} \quad (\text{D.27})$$

787 *Proof.* Plug the definition of A in (D.14),

$$\begin{aligned}
T &\gtrsim \frac{\lambda\Delta\sigma^2}{\gamma M\varepsilon^2} \log^{\frac{1}{2}} \frac{T}{\delta} + \frac{\Delta}{\varepsilon} \cdot \sqrt{\frac{L^2\sigma^2 K \log \frac{MT}{\delta}}{\varepsilon} + \frac{(L+K\tau)\Delta}{\varepsilon}} + \frac{\sqrt{L\Delta\rho^2 d \log \frac{T}{\delta}}}{\varepsilon} \\
&\quad + \frac{\Delta}{\varepsilon} \cdot \sqrt{\frac{L^2 K}{\varepsilon}} \sqrt{\frac{d \log^2 \frac{MT}{\delta}}{K} \rho^2 + K \|\sigma\|_{2\alpha}^{2\alpha} \cdot \rho^{2(1-\alpha)}} \\
&\asymp \frac{\lambda\Delta\sigma^2}{\gamma M\varepsilon^2} \log^{\frac{1}{2}} \frac{T}{\delta} + \frac{\Delta}{\varepsilon} \cdot \sqrt{\frac{L^2\sigma^2 K \log \frac{MT}{\delta}}{\varepsilon} + \frac{(L+K\tau)\Delta}{\varepsilon}} \\
&\quad + \frac{\Delta}{\varepsilon} \cdot \sqrt{\frac{L^2 K}{\varepsilon}} \sqrt{\frac{d \log^2 \frac{MT}{\delta}}{K} \rho^2 + K \|\sigma\|_{2\alpha}^{2\alpha} \cdot \rho^{2(1-\alpha)}}.
\end{aligned} \tag{D.28}$$

788 Hence the optimal ρ is given by

$$\rho \asymp \max \left\{ \|\sigma\|_{2\alpha} \left(\frac{K}{\sqrt{d} \log \frac{MT}{\delta}} \right)^{1/\alpha}, \left(\frac{\|\sigma\|_{2\alpha}^{2\alpha}}{\varepsilon} \right)^{\frac{1}{2(\alpha-1)}}, \sigma_\infty, G_\infty \right\}. \tag{D.29}$$

789 Note that $\left(\frac{\|\sigma\|_{2\alpha}^{2\alpha}}{\varepsilon} \right)^{\frac{1}{2(\alpha-1)}} \gtrsim G_\infty \vee \sigma_\infty$ and this implies

$$\begin{aligned}
T &\gtrsim \frac{\lambda\Delta\sigma^2}{\gamma M\varepsilon^2} \log^{\frac{1}{2}} \frac{T}{\delta} + \frac{\Delta}{\varepsilon} \cdot \sqrt{\frac{L^2\sigma^2 K \log \frac{MT}{\delta}}{\varepsilon} + \frac{(L+K\tau)\Delta}{\varepsilon}} \\
&\quad + \frac{L\Delta}{\varepsilon^{\frac{3}{2}}} \left[\|\sigma\|_{2\alpha} d^{\frac{1}{2} - \frac{1}{2\alpha}} K^{\frac{1}{\alpha}} \log^{1-\frac{1}{\alpha}} \frac{MT}{\delta} + \left(\frac{\|\sigma\|_{2\alpha}^{2\alpha}}{\varepsilon} \right)^{\frac{1}{2(\alpha-1)}} d^{\frac{1}{2}} \log \frac{MT}{\delta} \right] \\
&\asymp \frac{\lambda\Delta\sigma^2}{\gamma M\varepsilon^2} \log^{\frac{1}{2}} \frac{T}{\delta} + \frac{L\Delta}{\varepsilon^{\frac{3}{2}}} \cdot \sqrt{\sigma^2 K \log \frac{MT}{\delta}} + \frac{(L+K\tau)\Delta}{\varepsilon} + \frac{L\Delta}{\varepsilon^{\frac{3}{2}}} \left(\frac{\|\sigma\|_{2\alpha}^{2\alpha}}{\varepsilon} \right)^{\frac{1}{2(\alpha-1)}} d^{\frac{1}{2}} \log \frac{MT}{\delta}.
\end{aligned} \tag{D.30}$$

790 In the last equation we use $K \gtrsim \log \frac{MT}{\delta} \left(\frac{\|\sigma\|_{2\alpha} d^{\frac{1}{2} - \frac{1}{2\alpha}}}{\sigma} \right)^{\frac{2\alpha}{\alpha-2}}$. Solve ε and we get the upper

791 bound of $\frac{\lambda}{KR} \sum_{r=0}^{R-1} \sum_{k=0}^{K-1} \|\nabla f_\gamma^{H_r}(\bar{z}_{r,k})\|_{H_r}^2$.

792 Further note that $A = \tilde{O}(1)$, $B = \tilde{O}(K)$, $B_1 = \tilde{O}(K)$, $\eta = \tilde{O}(1/\sqrt{T})$ and we can get the upper
793 bound of $1 - \beta_2$ as:

$$1 - \beta_2 = \tilde{O} \left(\frac{1}{K^{3/2} R^{1/2}} \right). \tag{D.31}$$

794 This completes the proof. \square

795 **Theorem D.3** (Complete version of Theorem 3). *Under the conditions of Theorem D.2, let $\gamma = \frac{\lambda}{L}$
796 and thus $\Omega_0 \subset \{x : f(x) - f_* \leq 4(f(x_0) - f_*)\}$, $\Delta \asymp f(x_0) - f_*$. Then with probability no less
797 than $1 - \delta$, Local Adam with optimal η, ρ yields $\frac{\lambda}{KR} \sum_{r=0}^{R-1} \sum_{k=0}^{K-1} \|\nabla f(\bar{z}_{r,k})\|_{H_r}^2 \leq \varepsilon$ if*

$$T \gtrsim \frac{L\Delta\sigma^2}{M\varepsilon^2} \log^{\frac{1}{2}} \frac{T}{\delta} + \frac{L\Delta}{\varepsilon^{\frac{3}{2}}} \cdot \sqrt{\sigma^2 K \log \frac{MT}{\delta}} + \frac{(L+K\tau)\Delta}{\varepsilon} + \frac{L\Delta}{\varepsilon^{\frac{3}{2}}} \left(\frac{\|\sigma\|_{2\alpha}^{2\alpha}}{\varepsilon} \right)^{\frac{1}{2(\alpha-1)}} d^{\frac{1}{2}} \log \frac{MT}{\delta}. \tag{D.32}$$

798 And equivalently,

$$\begin{aligned} \frac{\lambda}{KR} \sum_{r=0}^{R-1} \sum_{k=0}^{K-1} \|\nabla f(\bar{z}_{r,k})\|_{H_r^{-1}}^2 &\lesssim \frac{\tau\Delta}{R} + \frac{L\Delta}{KR} + \sqrt{\frac{L\Delta\sigma^2}{MKR}} \log^{\frac{1}{4}} \frac{KR}{\delta} \\ &\quad + \frac{(L\Delta\sigma)^{\frac{2}{3}}}{K^{\frac{1}{3}}R^{\frac{2}{3}}} \log^{\frac{1}{3}} \frac{MKR}{\delta} + \left(\|\sigma\|_{2\alpha} d^{\frac{1}{2} - \frac{1}{2\alpha}}\right)^{\frac{2\alpha}{3\alpha-2}} \left(\frac{L\Delta \log \frac{MKR}{\delta}}{KR}\right)^{\frac{2(\alpha-1)}{3\alpha-2}}. \end{aligned} \quad (\text{D.33})$$

799 Further, if $1 - \beta_2 \lesssim \frac{G_\infty^2 + \sigma_\infty^2}{\rho^2 \log \frac{dR}{\delta}}$, where ρ is defined in (D.29), then with probability no less than
800 $1 - 2\delta$,

$$\begin{aligned} \frac{1}{KR} \sum_{r=0}^{R-1} \sum_{k=0}^{K-1} \|\nabla f(\bar{z}_{r,k})\|^2 &\lesssim \left(1 + \frac{G_\infty + \sigma_\infty}{\lambda}\right) \left[\frac{\tau\Delta}{R} + \frac{L\Delta}{KR} + \sqrt{\frac{L\Delta\sigma^2}{MKR}} \log^{\frac{1}{4}} \frac{KR}{\delta} + \frac{(L\Delta\sigma)^{\frac{2}{3}}}{K^{\frac{1}{3}}R^{\frac{2}{3}}} \log^{\frac{1}{3}} \frac{MKR}{\delta} \right. \\ &\quad \left. + \left(\|\sigma\|_{2\alpha} d^{\frac{1}{2} - \frac{1}{2\alpha}}\right)^{\frac{2\alpha}{3\alpha-2}} \left(\frac{L\Delta \log \frac{MKR}{\delta}}{KR}\right)^{\frac{2(\alpha-1)}{3\alpha-2}} \right]. \end{aligned} \quad (\text{D.34})$$

801 *Proof.* By Lemma D.6, we have $\Omega_0 \subset \{x : f(x) - f_* \leq 4(f(x_0) - f_*)\}$, $\Delta \asymp f(x_0) - f_*$. By Lemma
802 D.4, we have $\|\nabla f(\bar{z}_{r,k})\|_{H_r^{-1}} \leq 2\|\nabla f_\gamma^{H_r}(\bar{z}_{r,k})\|_{H_r^{-1}}$. Therefore, the bound for T in Theorem D.2

803 will reduce to (D.32). Solve ε and we get the upper bound of $\frac{\lambda}{KR} \sum_{r=0}^{R-1} \sum_{k=0}^{K-1} \|\nabla f(\bar{z}_{r,k})\|_{H_r^{-1}}^2$.

804 Now we turn to bound $\|H_r\|$. Note that $H_{r+1} = \mathbf{diag}(\sqrt{v_{r+1} + \lambda^2})$ and

$$\begin{aligned} [v_{r+1}]_i &= (1 - \beta_2) \sum_{j=0}^{rK-1} \beta_2^{rK-j-1} \mathbb{E}_m [\widehat{g}_j^m]_i^2 \\ &= (1 - \beta_2) \sum_{j=0}^{rK-1} \beta_2^{rK-j-1} \left(\mathbb{E}_m \left[[\widehat{g}_j^m]_i^2 - \mathbb{E}_j [\widehat{g}_j^m]_i^2 \right] + \mathbb{E}_m \mathbb{E}_j [\widehat{g}_j^m]_i^2 \right) \\ &\leq (1 - \beta_2) \sum_{j=0}^{rK-1} \beta_2^{rK-j-1} \mathbb{E}_m \left[[\widehat{g}_j^m]_i^2 - \mathbb{E}_j [\widehat{g}_j^m]_i^2 \right] + \sigma_\infty^2 + 3G_\infty^2, \end{aligned} \quad (\text{D.35})$$

805 where the last inequality is due to Lemma B.2. Define

$$[\theta_j]_i = \begin{cases} (1 - \beta_2) \beta_2^{rK-j-1} \mathbb{E}_m \left[[\widehat{g}_j^m]_i^2 - \mathbb{E}_j [\widehat{g}_j^m]_i^2 \right], & \text{if event } E_j \text{ holds,} \\ 0, & \text{otherwise.} \end{cases} \quad (\text{D.36})$$

806 Further note that

$$|[\theta_j]_i| \leq (1 - \beta_2) \rho^2 \stackrel{\text{def}}{=} c, \quad (\text{D.37})$$

807

$$\begin{aligned} \text{Var}_j([\theta_j]_i) &\leq \frac{(1 - \beta_2)^2 \beta_2^{2(rK-j-1)}}{M} \mathbb{E}_m \mathbb{E}_j \left[[\widehat{g}_j^m]_i^2 - \mathbb{E}_j [\widehat{g}_j^m]_i^2 \right]^2 \\ &\leq \frac{(1 - \beta_2)^2 \beta_2^{2(rK-j-1)}}{M} \mathbb{E}_m \mathbb{E}_j \left[[\widehat{g}_j^m]_i^2 - [\nabla f(x_j^m)]_i^2 \right]^2 \\ &\leq \frac{(1 - \beta_2)^2 \beta_2^{2(rK-j-1)}}{M} (2\sigma_\infty^4 + 8\sigma_\infty^2 G_\infty^2). \end{aligned} \quad (\text{D.38})$$

808 Let $b = G_\infty^2 + 3\sigma_\infty^2$, $V = \frac{2(1 - \beta_2)\sigma_\infty^2(\sigma_\infty^2 + 4G_\infty^2)}{M}$. If $1 - \beta_2 \lesssim \frac{G_\infty^2 + \sigma_\infty^2}{\rho^2 \log \frac{dR}{\delta}}$, then by Lemma

809 **B.1**, we have $|\sum_{j=0}^{rK-1} [\theta_j]_i| \leq b$ with probability no less than

$$1 - 2 \exp\left(-\frac{b^2}{2V + 2cb/3}\right) \geq 1 - \frac{\delta}{dR}, \quad (\text{D.39})$$

810 which implies $[H_r]_{i,i} \leq \lambda + 2G_\infty + 2\sigma_\infty$. Therefore, we have

$$\mathbb{P}\{E_T \text{ and } \|H_r\| \leq \lambda + 2G_\infty + 2\sigma_\infty \text{ for all } r \leq R\} \geq 1 - 2\delta. \quad (\text{D.40})$$

811 And thus

$$\begin{aligned} \frac{1}{KR} \sum_{r=0}^{R-1} \sum_{k=0}^{K-1} \|\nabla f(\bar{z}_{r,k})\|^2 &\lesssim \left(1 + \frac{G_\infty + \sigma_\infty}{\lambda}\right) \left[\frac{\tau\Delta}{R} + \frac{L\Delta}{KR} + \sqrt{\frac{L\Delta\sigma^2}{MKR}} \log^{\frac{1}{4}} \frac{T}{\delta} + \frac{(L\Delta\sigma)^{\frac{2}{3}}}{K^{\frac{1}{3}}R^{\frac{2}{3}}} \log^{\frac{1}{3}} \frac{MKR}{\delta}\right. \\ &\quad \left. + \left(\|\sigma\|_{2\alpha} d^{\frac{1}{2} - \frac{1}{2\alpha}}\right)^{\frac{2\alpha}{3\alpha-2}} \left(\frac{L\Delta \log \frac{MKR}{\delta}}{KR}\right)^{\frac{2(\alpha-1)}{3\alpha-2}}\right]. \end{aligned} \quad (\text{D.41})$$

812 \square

813 D.2 Preliminaries

814 We start with theoretical properties of weakly convex function and Moreau envelop, which are
815 repeatedly used in our proof.

816 **Lemma D.4.** Let $z \in \mathbb{R}^d$ and $y = y(z) := \arg \min_x f(x) + \frac{1}{2\gamma}\|x - z\|_H^2$ for some $H \succeq \lambda I_d$ and
817 $L/\lambda \geq \gamma^{-1} \geq 2\tau/\lambda$. Then

$$\nabla f_\gamma^H(z) = \nabla f(y) = \frac{H(z - y)}{\gamma}. \quad (\text{D.42})$$

818 If further assume $f_\gamma^H(z) - \min f_\gamma^\lambda \leq 2\Delta$, $0 \leq \eta \leq \frac{\lambda}{L}$, then $z, y \in \Omega_0$, and

$$\|\nabla f(z)\|_{H^{-1}} \leq \frac{2\gamma L}{\lambda} \|\nabla f_\gamma^H(z)\|_{H^{-1}}, \quad (\text{D.43})$$

819

$$\|H(z - y) - \eta \nabla f(z)\|_{H^{-1}} \leq \gamma \|\nabla f(y)\|_{H^{-1}}. \quad (\text{D.44})$$

820

$$\|\nabla f_\gamma^H(z)\|_{H^{-1}}^2 \leq \frac{2}{\gamma} (f_\gamma^H(z) - \min f_\gamma^\lambda). \quad (\text{D.45})$$

821 *Proof.* Since y is the minimizer,

$$0 = \nabla_y \left[f(y) + \frac{1}{2\gamma} \|y - z\|_H^2 \right] = \nabla f(y) + \frac{H(y - z)}{\gamma}, \quad (\text{D.46})$$

822 and note that

$$\nabla f_\gamma^H(z) = \nabla_z \left[f(y(z)) + \frac{1}{2\gamma} \|y(z) - z\|_H^2 \right] = \frac{H(z - y)}{\gamma}. \quad (\text{D.47})$$

823 If $f_\gamma^H(z) - \min f_\gamma^\lambda \leq 2\Delta$, then $f_\gamma^\lambda(z) \leq f_\gamma^H(z)$ and

$$f_\gamma^\lambda(y) \leq f_\gamma^H(y) \leq f(y) \leq f_\gamma^H(z) \leq f(z), \quad (\text{D.48})$$

824 which implies $y, z \in \Omega_0$.

825 By mean value theorem, there exists a symmetric matrix $-\tau I_d \leq H_g \preceq LI_d$, such that

$$\nabla f(z) - \nabla f(y) = H_g(z - y) = \gamma H_g H^{-1} \nabla f(y). \quad (\text{D.49})$$

826 Hence,

$$\|\nabla f(z) - \nabla f(y)\|_{H^{-1}} \leq \gamma \|H^{-1} \nabla f(y)\|_{H_g H^{-1} H_g} \leq \frac{\gamma L}{\lambda} \|\nabla f_\gamma^H(z)\|_{H^{-1}}. \quad (\text{D.50})$$

827

$$\|\nabla f(z)\|_{H^{-1}} \leq (1 + \frac{\gamma L}{\lambda}) \|\nabla f_\gamma^H(z)\|_{H^{-1}} \leq \frac{2\gamma L}{\lambda} \|\nabla f_\gamma^H(z)\|_{H^{-1}}. \quad (\text{D.51})$$

828 Also,

$$H(z - y) - \eta \nabla f(z) = (\gamma I_d - \eta(I_d + \gamma H_g H^{-1})) \nabla f(y) =: \gamma \Lambda \nabla f(y). \quad (\text{D.52})$$

829 By noticing that

$$-I_d \preceq H^{-1/2} \Lambda H^{1/2} = I_d - \eta \gamma^{-1} - \eta H^{-1/2} H_g H^{-1/2} \preceq I_d, \quad (\text{D.53})$$

830 we have $\|H(z - y) - \eta \nabla f(z)\|_{H^{-1}} \leq \gamma \|\nabla f(y)\|_{H^{-1}}$.

831 Last,

$$\min f_\gamma^\lambda \leq f_\gamma^\lambda(y) \leq f(y) = f_\gamma^H(z) - \frac{1}{2\gamma} \|y - z\|_H^2 = f_\gamma^H(z) - \frac{\gamma}{2} \|\nabla f_\gamma^H(z)\|_{H^{-1}}^2. \quad (\text{D.54})$$

832 This completes the proof. \square

833 **Lemma D.5.** *If $x, y \in \Omega$, then*

$$-\langle x - y, \nabla f(x) - \nabla f(y) \rangle + \frac{1}{L} \|\nabla f(x) - \nabla f(y)\|^2 \leq 2\tau \|x - y\|^2. \quad (\text{D.55})$$

834 *Proof.* By mean value theorem, there exists a symmetric matrix $-\tau I_d \preceq H \preceq L I_d$, such that

$$\nabla f(x) - \nabla f(y) = H(x - y). \quad (\text{D.56})$$

835 Therefore,

$$\begin{aligned} -\langle x - y, \nabla f(x) - \nabla f(y) \rangle + \frac{1}{L} \|\nabla f(x) - \nabla f(y)\|^2 &= (x - y)^T (-H + \frac{H^2}{L})(x - y) \\ &\leq (\tau + \frac{\tau^2}{L}) \|x - y\|^2 \\ &\leq 2\tau \|x - y\|^2. \end{aligned} \quad (\text{D.57})$$

836 \square

837 **Lemma D.6.** *If $\gamma = \frac{\lambda}{L}$, then for $z \in \Omega_0$, it holds that $\frac{f(z) - f_*}{2} \leq f_{1/L}(z) - f_* \leq f(z) - f_*$.*

838 *Proof.* By definition of Moreau envelop, the second inequality is trivial. Let $y = \arg \min_x f(x) +$

839 $\frac{L}{2} \|x - z\|^2$. Note that $x \rightarrow f(x) + \frac{L}{2} \|x - z\|^2$ is $2L$ -smooth. Then we have

$$f(z) \leq f(y) + \frac{L}{2} \|y - z\|^2 + L \|y - z\|^2 = f_{1/L}(z) + L \|y - z\|^2. \quad (\text{D.58})$$

840 Furthermore, by Lemma D.4

$$\frac{L}{2} \|y - z\|^2 = \frac{1}{2L} \|\nabla f(y)\|^2 \leq f(y) - f_*. \quad (\text{D.59})$$

841 Therefore, $f(z) - f_* \leq f_{1/L}(z) - f_* + L \|y - z\|^2 \leq 2(f_{1/L}(z) - f_*)$. \square

842 Next, we show that event E_t implies all the iterates remain in certain area.

843 **Lemma D.7.** *If $\frac{\eta\sigma}{\lambda} \sqrt{KA} \leq \sqrt{\frac{\Delta\gamma}{160\lambda}}$, then event E_t implies that for all $j \leq t, m \in [M]$, we have*

844 $\bar{z}_j \in \Omega_0, x_j^m, \bar{x}_j, z_j^m \in \Omega$. And $\|x_j^m - x_j^n\| \leq \frac{\eta\sigma}{\lambda} \sqrt{KA}$ for all m, n .

845 *Proof.* Event E_t implies that for all $j \leq t$,

$$f_\gamma^\lambda(\bar{z}_j) - \min f_\gamma^\lambda \leq 2\Delta, \|z_j^m - z_j^n\| \leq \frac{\eta\sigma}{\lambda}\sqrt{KA} \leq \sqrt{\frac{\Delta\gamma}{160\lambda}}. \quad (\text{D.60})$$

846 Hence $\bar{z}_j \in \Omega_0, \|z_j^m - \bar{z}_j\| \leq \frac{\eta\sigma}{\lambda}\sqrt{KA}$ and $z_j^m \in \mathbf{B}_{R_0}(\Omega_0) \subset \Omega$. Also, notice that $\bar{x}_j \in$
847 $\mathbf{conv}\{\bar{z}_i\}_{i \leq j} \subset \mathbf{conv}(\Omega_0) \subset \Omega$ and $x_j^m - x_j^n \in \mathbf{conv}\{z_i^m - z_i^n\}_{i \leq j}$. We have

$$\|x_j^m - x_j^n\| \leq \frac{\eta\sigma}{\lambda}\sqrt{KA}, \|x_j^m - \bar{x}_j\| \leq \frac{\eta\sigma}{\lambda}\sqrt{KA} \leq \sqrt{\frac{\Delta\gamma}{160\lambda}}. \quad (\text{D.61})$$

848 Therefore by Lemma B.4, $x_j^m \in \mathbf{B}_{R_0}(\mathbf{conv}(\Omega_0)) = \Omega$. \square

849 The following lemma shows that the second order momentum v_t^m does not change too much from
850 $v_{r(t)}$ during local training with high probability, which is also repeatedly used in our proof.

851 **Lemma D.8.** Let $B := \max \left\{ \frac{6K(G_\infty^2 + \sigma_\infty^2)}{\lambda^2}, \frac{16\rho^2}{\lambda^2} \log \frac{dMT}{\delta}, 2^6 \frac{\sqrt{K}(G_\infty + \sigma_\infty)\sigma_\infty}{\lambda^2} \log^{1/2} \frac{dMT}{\delta} \right\}$.

852 If $\rho \geq \max\{3\sigma_\infty, 2G_\infty\}$, then the following holds

$$\mathbb{P}(E_{t,1}) \geq \mathbb{P}(E_t) - \frac{\delta}{4T}. \quad (\text{D.62})$$

853 *Proof.* Let $t = rK + k$. By the update rule of local Adam, we have

$$v_t^m = \beta_2^{k+1}v_r + (1 - \beta_2) \sum_{j=rK}^t \beta_2^{t-j} \widehat{g}_j^m \odot \widehat{g}_j^m \succeq \beta_2^K v_r, \quad (\text{D.63})$$

854 and hence

$$H_t^m = \mathbf{diag}(\sqrt{v_t^m + \lambda^2}) \succeq \beta_2^{K/2} \mathbf{diag}(\sqrt{v_r + \lambda^2}) = \beta_2^{K/2} H_r. \quad (\text{D.64})$$

855 For the upper bound, for any index $i \in [d]$, by Lemma B.2,

$$\mathbb{E}_j[\widehat{g}_j^m]_i^2 \leq \sigma_i^2 + \mathbb{E}_j[g_j^m]_i^2 \leq \sigma_\infty^2 + 3G_\infty^2. \quad (\text{D.65})$$

856 Therefore,

$$[v_t^m]_i \leq [v_r]_i + (1 - \beta_2)K(\sigma_\infty^2 + 3G_\infty^2) + (1 - \beta_2) \sum_{j=rK}^t \left[[\widehat{g}_j^m]_i^2 - \mathbb{E}_j[\widehat{g}_j^m]_i^2 \right]. \quad (\text{D.66})$$

857 Define

$$[\theta_j^m]_i = \begin{cases} [\widehat{g}_j^m]_i^2 - \mathbb{E}_j[\widehat{g}_j^m]_i^2, & \text{if event } E_j \text{ holds,} \\ 0, & \text{otherwise.} \end{cases} \quad (\text{D.67})$$

858 Event E_t implies $[\theta_j^m]_i = [\widehat{g}_j^m]_i^2 - \mathbb{E}_j[\widehat{g}_j^m]_i^2$. Further note that $|[\theta_j^m]_i| \leq \rho^2 \stackrel{\text{def}}{=} c$,

$$\begin{aligned} \text{Var}_j([\theta_j^m]_i) &\leq \mathbb{E}_j \left[[\widehat{g}_j^m]_i^2 - [\nabla f(x_j^m)]_i^2 \right]^2 \\ &= \mathbb{E}_j \left[[\widehat{g}_j^m]_i - [\nabla f(x_j^m)]_i \right]^2 \left[[\widehat{g}_j^m]_i + [\nabla f(x_j^m)]_i \right]^2 \\ &\stackrel{\text{AM-GM}}{\leq} 2\mathbb{E}_j \left[[\widehat{g}_j^m]_i - [\nabla f(x_j^m)]_i \right]^4 + 8\mathbb{E}_j \left[[\widehat{g}_j^m]_i - [\nabla f(x_j^m)]_i \right]^2 [\nabla f(x_j^m)]_i^2 \\ &\stackrel{\text{Lemma B.2}}{\leq} 2\sigma_\infty^4 + 8\sigma_\infty^2 G_\infty^2. \end{aligned} \quad (\text{D.68})$$

859 Let $b = B\lambda^2/2, V = 2K\sigma_\infty^2(\sigma_\infty^2 + 4G_\infty^2)$. Applying Lemma B.1, we have $|\sum_{j=rK}^t [\theta_j^m]_i| \leq b$ with
860 probability no less than

$$1 - 2 \exp\left(-\frac{b^2}{2V + 2cb/3}\right) \geq 1 - \frac{\delta}{4dMT}, \quad (\text{D.69})$$

861 which implies with probability no less than $1 - \frac{\delta}{4T}$, for any $m \in [M]$,

$$v_t^m \preceq v_r + (1 - \beta_2)K(\sigma_\infty^2 + 3G_\infty^2) + (1 - \beta_2)B\lambda^2/2 \preceq v_r + (1 - \beta_2)B\lambda^2. \quad (\text{D.70})$$

862 and thus

$$H_t^m \preceq \sqrt{1 + (1 - \beta_2)BH_r}. \quad (\text{D.71})$$

863 \square

864 D.3 Proof of Contraction

865 In this subsection, we aim to show contraction, *i.e.*, $\|x_t^m - x_t^n\|$ will not get too large during local
866 iterations with high probability. However, since the update of x_t^m involves the coupling of both first
867 order momentum and second order momentum, it is much harder than showing the contraction of
868 *Local* SGDM. Our solution below is in two folds.

869 We begin with showing contraction of the second order momentum in some sense.

870 **Lemma D.9.** *Let $B_1 := \max \left\{ \frac{16K\sigma_\infty^2}{\lambda^2}, \frac{16\rho^2}{\lambda^2} \log \frac{dMT}{\delta}, 2^6 \frac{\sqrt{K}(G_\infty + \sigma_\infty)\sigma_\infty}{\lambda^2} \log^{1/2} \frac{dMT}{\delta} \right\}$*

871 *and $1 - \beta_2 \leq \frac{1}{4K}$. If $\rho \geq \max\{3\sigma_\infty, 2G_\infty\}$, $\frac{\eta L\sigma}{\lambda} \sqrt{KA}G_\infty \leq 2\sigma_\infty^2$, then the following holds:*

$$\mathbb{P}(E_{t,2}) \geq \mathbb{P}(E_{t,1}) - \frac{\delta}{4T} \quad (\text{D.72})$$

872 *Proof.* Event $E_{t,1}$ implies for all $j \leq t$, $x_j^m, x_j^n \in \Omega$ and for any index $i \in [d]$,

$$\begin{aligned} \left| [v_t^m - v_t^n]_i \right| &= \left| (1 - \beta_2) \sum_{j=rK}^t \beta_2^{t-j} \left[[\widehat{g}_j^m]_i^2 - [\widehat{g}_j^n]_i^2 \right] \right| \\ &\leq \left| (1 - \beta_2) \sum_{j=rK}^t \beta_2^{t-j} \left[[\widehat{g}_j^m]_i^2 - [\widehat{g}_j^n]_i^2 - \mathbb{E}_j \left[[\widehat{g}_j^m]_i^2 - [\widehat{g}_j^n]_i^2 \right] \right] \right| \\ &\quad + \left| (1 - \beta_2) \sum_{j=rK}^t \beta_2^{t-j} \left[\mathbb{E}_j \left[[\widehat{g}_j^m]_i^2 - [\widehat{g}_j^n]_i^2 \right] - [\nabla f(x_j^m)]_i^2 - [\nabla f(x_j^n)]_i^2 \right] \right| \\ &\quad + \left| (1 - \beta_2) \sum_{j=rK}^t \beta_2^{t-j} \left[[\nabla f(x_j^m)]_i^2 - [\nabla f(x_j^n)]_i^2 \right] \right| \\ &\leq \left| (1 - \beta_2) \sum_{j=rK}^t \beta_2^{t-j} \left[[\widehat{g}_j^m]_i^2 - [\widehat{g}_j^n]_i^2 - \mathbb{E}_j \left[[\widehat{g}_j^m]_i^2 - [\widehat{g}_j^n]_i^2 \right] \right] \right| \\ &\quad + (1 - \beta_2)K \cdot 4\sigma_\infty^2 + (1 - \beta_2)K \cdot 2G_\infty \frac{\eta L\sigma}{\lambda} \sqrt{KA} \\ &\leq \left| (1 - \beta_2) \sum_{j=rK}^t \beta_2^{t-j} \left[[\widehat{g}_j^m]_i^2 - [\widehat{g}_j^n]_i^2 - \mathbb{E}_j \left[[\widehat{g}_j^m]_i^2 - [\widehat{g}_j^n]_i^2 \right] \right] \right| + 8(1 - \beta_2)K \cdot \sigma_\infty^2. \end{aligned} \quad (\text{D.73})$$

873 Here in the second inequality we apply Lemma B.2 and contraction results implied by $E_{t,1}$.

874 Define

$$[\Xi_j^{m,n}]_i = \begin{cases} \beta_2^{t-j} \left[[\widehat{g}_j^m]_i^2 - [\widehat{g}_j^n]_i^2 - \mathbb{E}_j \left[[\widehat{g}_j^m]_i^2 - [\widehat{g}_j^n]_i^2 \right] \right], & \text{if event } E_j \text{ holds,} \\ 0, & \text{otherwise.} \end{cases} \quad (\text{D.74})$$

875 Then we have

$$\left| [\Xi_j^{m,n}]_i \right| \leq 2\rho^2 \stackrel{\text{def}}{=} c, \quad (\text{D.75})$$

876

$$\begin{aligned}
\text{Var}_j([\Xi_j^{m,n}]_i) &\leq 2\mathbb{E}_j \left[[\widehat{g}_j^m]_i^2 - \mathbb{E}_j [\widehat{g}_j^m]_i^2 \right]^2 \\
&\leq 2\mathbb{E}_j \left[[\widehat{g}_j^m]_i^2 - [\nabla f(x_j^m)]_i^2 \right]^2 \\
&\leq 4\mathbb{E}_j \left[[\widehat{g}_j^m]_i - [\nabla f(x_j^m)]_i \right]^2 \cdot \left[[\widehat{g}_j^m]_i - [\nabla f(x_j^m)]_i \right]^2 + 4[\nabla f(x_j^m)]_i^2 \\
&\stackrel{\text{Lemma B.2}}{\leq} 4\sigma_\infty^4 + 16\sigma_\infty^2 G_\infty^2.
\end{aligned} \tag{D.76}$$

877 Let $b = B_1\lambda^2/2$, $V = 4K\sigma_\infty^2(\sigma_\infty^2 + 4G_\infty^2)$ and by Lemma B.1, we have $|\sum_{j=rK}^t [\Xi_j^{m,n}]_i| \leq b$ with
878 probability no less than

$$1 - 2 \exp\left(\frac{b^2}{2V + 2cb/3}\right) \geq 1 - \frac{\delta}{4dM^2T}. \tag{D.77}$$

879 This implies with probability no less than $1 - \frac{\delta}{4M^2T}$,

$$\left|v_t^m - v_t^n\right| \preceq (1 - \beta_2)B_1\lambda^2/2 + 8(1 - \beta_2)K \cdot \sigma_\infty^2 \preceq (1 - \beta_2)B_1\lambda^2. \tag{D.78}$$

880 Combine this inequality and event $E_{t,1}$,

$$\begin{aligned}
\left|\frac{H_r}{H_t^m} - \frac{H_r}{H_t^n}\right| &= \frac{\sqrt{v_r + \lambda^2}|v_t^n - v_t^m|}{\sqrt{v_t^m + \lambda^2}\sqrt{v_t^n + \lambda^2}(\sqrt{v_t^m + \lambda^2} + \sqrt{v_t^n + \lambda^2})} \\
&\preceq (1 - \beta_2)B_1 \frac{\sqrt{v_r + \lambda^2}}{(\sqrt{v_t^m + \lambda^2} + \sqrt{v_t^n + \lambda^2})} \\
&\preceq (1 - \beta_2)B_1.
\end{aligned} \tag{D.79}$$

881 The last inequality is due to event $E_{t,1}$ and $1 - \beta_2 \leq \frac{1}{4K}$. We can conclude that under event $E_{t,1}$,

882 with probability no less than $1 - \frac{\delta}{4T}$, the inequality above holds for any $m, n \in [M]$, which implies

$$883 \mathbb{P}(E_{t,2}) \geq \mathbb{P}(E_{t,1}) - \frac{\delta}{4T}. \quad \square$$

884 Now we are ready to prove contraction of z_t^m .

885 **Lemma D.10.** Let $A := \max\left\{\frac{2^{20}\rho^2 d}{K\sigma^2} \log \frac{MT}{\delta}, 2^{20} \log \frac{MT}{\delta}, \frac{2^8 K \|2\sigma\|_{2\alpha}^{2\alpha}}{\sigma^2 \rho^{2(\alpha-1)}}\right\}$. If $\eta \leq$

886 $\min\left\{\frac{\lambda}{60K\tau}, \frac{(1 - \beta_1)^2 \lambda}{64L}\right\}$, $\rho \geq \max\{3\sigma_\infty, 2G_\infty\}$, and

$$(1 - \beta_2)K^{1/2} \leq \min\left\{\frac{(1 - \beta_1)}{4B_1}, \frac{(1 - \beta_1)\sigma}{2^{12}B_1G} \sqrt{A}, \frac{1 - \beta_1}{4B}\right\}, \tag{D.80}$$

887 then the following holds:

$$\mathbb{P}(E_{t,3}) \geq \mathbb{P}(E_{t,2}) - \frac{\delta}{4T}. \tag{D.81}$$

888 *Proof.* If $t \bmod K \equiv -1$, then $z_{t+1}^m = z_{t+1}^n$ for all m, n and the claim is trivial. Below we assume
889 that $t \bmod K \not\equiv -1$. The update rules implies

$$\begin{aligned}
\|z_{t+1}^m - z_{t+1}^n\|_{H_r}^2 &\stackrel{(D.6)}{=} \|z_t^m - z_t^n\|_{H_r}^2 - 2\eta \left\langle z_t^m - z_t^n, (H_t^m)^{-1}(\widehat{g}_t^m + e_t^m) - (H_t^n)^{-1}(\widehat{g}_t^n + e_t^n) \right\rangle_{H_r} \\
&\quad + \eta^2 \underbrace{\left\| (H_t^m)^{-1}(\widehat{g}_t^m + e_t^m) - (H_t^n)^{-1}(\widehat{g}_t^n + e_t^n) \right\|_{H_r}^2}_{\textcircled{1}}.
\end{aligned} \tag{D.82}$$

890 Note that the first order term is

$$\begin{aligned}
& \left\langle z_t^m - z_t^n, (H_t^m)^{-1}(\widehat{g}_t^m + e_t^m) - (H_t^n)^{-1}(\widehat{g}_t^n + e_t^n) \right\rangle_{H_r} \\
&= \langle z_t^m - z_t^n, \nabla f(x_t^m) - \nabla f(x_t^n) \rangle \\
&+ \left\langle z_t^m - z_t^n, \widehat{g}_t^m - \widehat{g}_t^n - \nabla f(x_t^m) + \nabla f(x_t^n) \right\rangle \\
&+ \underbrace{\left\langle z_t^m - z_t^n, (H_t^m)^{-1}e_t^m - (H_t^n)^{-1}e_t^n \right\rangle_{H_r}}_{\textcircled{2}} \\
&+ \underbrace{\left\langle z_t^m - z_t^n, (H_r(H_t^m)^{-1} - I_d)\widehat{g}_t^m - (H_r(H_t^n)^{-1} - I_d)\widehat{g}_t^n \right\rangle}_{\textcircled{3}}.
\end{aligned} \tag{D.83}$$

891 And for the first term above,

$$\begin{aligned}
\langle z_t^m - z_t^n, \nabla f(x_t^m) - \nabla f(x_t^n) \rangle &= \langle x_t^m - x_t^n, \nabla f(x_t^m) - \nabla f(x_t^n) \rangle \\
&+ \langle z_t^m - z_t^n - (x_t^m - x_t^n), \nabla f(x_t^m) - \nabla f(x_t^n) \rangle \\
&\geq \langle x_t^m - x_t^n, \nabla f(x_t^m) - \nabla f(x_t^n) \rangle \\
&\quad - \frac{L}{\lambda} \|(z_t^m - z_t^n) - (x_t^m - x_t^n)\|_{H_r}^2 - \frac{\lambda}{4L} \|\nabla f(x_t^m) - \nabla f(x_t^n)\|_{H_r^{-1}}^2
\end{aligned} \tag{D.84}$$

892 By definition of $\{z_t^m\}$ and event $E_{t,2}$,

$$\begin{aligned}
\|(z_t^m - z_t^n) - (x_t^m - x_t^n)\|_{H_r}^2 &= \left(\frac{\eta\beta_1}{1-\beta_1} \right)^2 \|(H_t^m)^{-1}u_t^m - (H_t^n)^{-1}u_t^n\|_{H_r}^2 \\
&\leq 2 \left(\frac{\eta\beta_1}{1-\beta_1} \right)^2 \left[\|(H_t^m)^{-1} - (H_t^n)^{-1}\|_{H_r}^2 \|u_t^m\|_{H_r}^2 + \|(H_t^n)^{-1}\|_{H_r}^2 \|u_t^m - u_t^n\|_{H_r}^2 \right] \\
&\stackrel{\mathcal{A}_{t,1}, \mathcal{A}_{t,2}}{\leq} 2 \left(\frac{\eta\beta_1}{1-\beta_1} \right)^2 \left[[(1-\beta_2)B_1]^2 \|u_t^m\|_{H_r^{-1}}^2 + 4\|u_t^m - u_t^n\|_{H_r^{-1}}^2 \right].
\end{aligned} \tag{D.85}$$

893 Besides,

$$\begin{aligned}
\textcircled{1} &\leq 4 \underbrace{\|(H_t^m)^{-1}e_t^m - (H_t^n)^{-1}e_t^n\|_{H_r}^2}_{(*)} + 4 \underbrace{\|(H_r(H_t^m)^{-1} - I_d)\widehat{g}_t^m - (H_r(H_t^n)^{-1} - I_d)\widehat{g}_t^n\|_{H_r^{-1}}^2}_{(**)} \\
&\quad + 4\|\widehat{g}_t^m - \widehat{g}_t^n - \nabla f(x_t^m) + \nabla f(x_t^n)\|_{H_r^{-1}}^2 + 4\|\nabla f(x_t^m) - \nabla f(x_t^n)\|_{H_r^{-1}}^2,
\end{aligned} \tag{D.86}$$

894

$$|\textcircled{2}| \leq \frac{1}{8\eta K} \|z_t^m - z_t^n\|_{H_r}^2 + 2\eta K \cdot (*). \tag{D.87}$$

895

$$|\textcircled{3}| \leq \frac{1}{8\eta K} \|z_t^m - z_t^n\|_{H_r}^2 + 2\eta K \cdot (**). \tag{D.88}$$

896

$$\begin{aligned}
(*) &\stackrel{(D.5)}{=} \left(\frac{\beta_1}{1-\beta_1} \right)^2 \left\| [(H_t^m)^{-1} - (H_{t-1}^m)^{-1}] u_t^m - [(H_t^n)^{-1} - (H_{t-1}^n)^{-1}] u_t^n \right\|_{H_r}^2 \\
&\leq 2 \left(\frac{\beta_1}{1-\beta_1} \right)^2 \left[\left\| [(H_t^m)^{-1} - (H_{t-1}^m)^{-1} - (H_t^n)^{-1} + (H_{t-1}^n)^{-1}] u_t^m \right\|_{H_r}^2 \right. \\
&\quad \left. + \left\| [(H_t^n)^{-1} - (H_{t-1}^n)^{-1}] (u_t^m - u_t^n) \right\|_{H_r}^2 \right] \\
&\stackrel{\mathcal{A}_{t,1}, \mathcal{A}_{t,2}}{\leq} 2 \left(\frac{\beta_1}{1-\beta_1} \right)^2 \left[4[(1-\beta_2)B_1]^2 \|u_t^m\|_{H_r^{-1}}^2 + 4[(1-\beta_2)B]^2 \|(u_t^m - u_t^n)\|_{H_r^{-1}}^2 \right] \\
&= 8 \left(\frac{\beta_1(1-\beta_2)}{1-\beta_1} \right)^2 \left[B_1^2 \|u_t^m\|_{H_r^{-1}}^2 + B^2 \|(u_t^m - u_t^n)\|_{H_r^{-1}}^2 \right]
\end{aligned} \tag{D.89}$$

$$\begin{aligned}
(**) &\leq 2 \left[\left\| H_r((H_t^m)^{-1} - (H_t^n)^{-1})\widehat{g}_t^m \right\|_{H_r^{-1}}^2 + \left\| (H_r(H_t^n)^{-1} - I_d)(\widehat{g}_t^m - \widehat{g}_t^n) \right\|_{H_r^{-1}}^2 \right] \\
&\stackrel{\mathcal{A}_{t,1}, \mathcal{A}_{t,2}}{\leq} 2 \left[[(1 - \beta_2)B_1]^2 \|\widehat{g}_t^m\|_{H_r^{-1}}^2 + [(1 - \beta_2)B]^2 \|\widehat{g}_t^m - \widehat{g}_t^n\|_{H_r^{-1}}^2 \right] \\
&\leq 2(1 - \beta_2)^2 \left[B_1^2 \|\widehat{g}_t^m\|_{H_r^{-1}}^2 + 2B^2 \left(\|\widehat{g}_t^m - \widehat{g}_t^n - \nabla f(x_t^m) + \nabla f(x_t^n)\|_{H_r^{-1}}^2 + \|\nabla f(x_t^m) - \nabla f(x_t^n)\|_{H_r^{-1}}^2 \right) \right] \\
&\hspace{15em} \text{(D.90)}
\end{aligned}$$

898 Here we repeatedly apply $\|H_r(H_t^n)^{-1} - I_d\| \leq (1 - \beta_2)B$ and $\|H_r((H_t^m)^{-1} - (H_t^n)^{-1})\| \leq$
899 $(1 - \beta_2)B_1$ by event $E_{t,2}$. Plug in (D.82),

$$\begin{aligned}
\|z_{t+1}^m - z_{t+1}^n\|_{H_r}^2 &\leq \|z_t^m - z_t^n\|_{H_r}^2 \underbrace{\stackrel{\text{(D.83)}}{-} 2\eta \langle z_t^m - z_t^n, \widehat{g}_t^m - \widehat{g}_t^n - \nabla f(x_t^m) + \nabla f(x_t^n) \rangle}_{(***)} \stackrel{\text{(D.84)}}{-} 2\eta \langle x_t^m - x_t^n, \nabla f(x_t^m) - \nabla f(x_t^n) \rangle \\
&\stackrel{\text{(D.84)}}{+} 2\eta \left[\frac{L}{\lambda} \|(z_t^m - z_t^n) - (x_t^m - x_t^n)\|_{H_r}^2 + \frac{\lambda}{4L} \|\nabla f(x_t^m) - \nabla f(x_t^n)\|_{H_r^{-1}}^2 \right] \\
&\stackrel{\text{(D.84)}}{-} 2\eta \cdot (\textcircled{2} + \textcircled{3}) + \eta^2 \cdot \textcircled{1} \\
&\leq \|z_t^m - z_t^n\|_{H_r}^2 + (***) + 2\eta \left[\frac{L}{\lambda} \|(z_t^m - z_t^n) - (x_t^m - x_t^n)\|_{H_r}^2 + \frac{\lambda}{4L} \|\nabla f(x_t^m) - \nabla f(x_t^n)\|_{H_r^{-1}}^2 \right] \\
&\quad + 2\eta \left[\frac{1}{4\eta K} \|z_t^m - z_t^n\|_{H_r}^2 \stackrel{\text{(D.87)}}{+} 2\eta K \cdot (*) + 2\eta K \cdot (**) \right] \\
&\stackrel{\text{(D.86)}}{+} 4\eta^2 \left[(*) + (**) + \|\widehat{g}_t^m - \widehat{g}_t^n - \nabla f(x_t^m) + \nabla f(x_t^n)\|_{H_r^{-1}}^2 + \|\nabla f(x_t^m) - \nabla f(x_t^n)\|_{H_r^{-1}}^2 \right] \\
&\leq \left(1 + \frac{1}{2K}\right) \|z_t^m - z_t^n\|_{H_r}^2 + (***) + \frac{2\eta L}{\lambda} \|(z_t^m - z_t^n) - (x_t^m - x_t^n)\|_{H_r}^2 \\
&\quad + \left(\frac{\eta}{2L} + \frac{4\eta^2}{\lambda}\right) \|\nabla f(x_t^m) - \nabla f(x_t^n)\|_{H_r^{-1}}^2 + 4\eta^2 \underbrace{\|\widehat{g}_t^m - \widehat{g}_t^n - \nabla f(x_t^m) + \nabla f(x_t^n)\|_{H_r^{-1}}^2}_{(\#)} \\
&\quad + 8\eta^2 K ((*) + (**)) \\
&\leq \left(1 + \frac{1}{2K}\right) \|z_t^m - z_t^n\|_{H_r}^2 - 2\eta \underbrace{\langle x_t^m - x_t^n, \nabla f(x_t^m) - \nabla f(x_t^n) \rangle}_{(\#\#)} + \frac{\eta}{L} \|\nabla f(x_t^m) - \nabla f(x_t^n)\|_{H_r^{-1}}^2 \\
&\quad - 2\eta \langle z_t^m - z_t^n, \widehat{g}_t^m - \widehat{g}_t^n - \nabla f(x_t^m) + \nabla f(x_t^n) \rangle + 8\eta^2 \cdot (\#) \\
&\stackrel{\text{(D.85)}}{+} \frac{4\eta L}{\lambda} \left(\frac{\eta\beta_1}{1 - \beta_1} \right)^2 \left[[(1 - \beta_2)B_1]^2 \|u_t^m\|_{H_r^{-1}}^2 + 4\|u_t^m - u_t^n\|_{H_r^{-1}}^2 \right] \\
&\stackrel{\text{(D.89)}}{+} 64\eta^2 K \left(\frac{\beta_1(1 - \beta_2)}{1 - \beta_1} \right)^2 \left[B_1^2 \|u_t^m\|_{H_r^{-1}}^2 + B^2 \|(u_t^m - u_t^n)\|_{H_r^{-1}}^2 \right] \stackrel{\text{(D.90)}}{+} 16\eta^2 K(1 - \beta_2)^2 B_1^2 \|\widehat{g}_t^m\|_{H_r^{-1}}^2 \\
&\leq \left(1 + \frac{1}{2K}\right) \|z_t^m - z_t^n\|_{H_r}^2 + (\#\#) + 8\eta^2 \cdot (\#) \\
&\quad - 2\eta \langle z_t^m - z_t^n, \widehat{g}_t^m - \widehat{g}_t^n - \mathbb{E}_t[\widehat{g}_t^m - \widehat{g}_t^n] \rangle - 2\eta \langle z_t^m - z_t^n, \mathbb{E}_t[\widehat{g}_t^m - \widehat{g}_t^n] - \nabla f(x_t^m) + \nabla f(x_t^n) \rangle \\
&\quad + 24\eta^2 \|u_t^m - u_t^n\|_{H_r^{-1}}^2 + 65\eta^2 K \left(\frac{\beta_1(1 - \beta_2)}{1 - \beta_1} \right)^2 \underbrace{B_1^2 \|u_t^m\|_{H_r^{-1}}^2 + 16\eta^2 K(1 - \beta_2)^2 B_1^2 \|\widehat{g}_t^m\|_{H_r^{-1}}^2}_{(\#\#\#)} \\
&\leq \left(1 + \frac{1}{K}\right) \|z_t^m - z_t^n\|_{H_r}^2 + (\#\#) + 8\eta^2 \cdot (\#) - 2\eta \langle z_t^m - z_t^n, \widehat{g}_t^m - \widehat{g}_t^n - \mathbb{E}_t[\widehat{g}_t^m - \widehat{g}_t^n] \rangle \\
&\quad + (\#\#\#) \stackrel{\text{Lemma B.2}}{+} \frac{8\eta^2 K}{\lambda} \cdot \frac{\|2\sigma\|_{2\alpha}^{2\alpha}}{\rho^{2(\alpha-1)}}.
\end{aligned} \tag{D.91}$$

900 In the second to last inequality we apply $8K(1 - \beta_2)^2 B^2 \leq (1 - \beta_1)^2$ and $\frac{\eta L}{\lambda} \leq (1 - \beta_1)^2$. Also
 901 notice that by definition of $\{u_t^m\}$,

$$u_t^m = (1 - \beta_1) \sum_{j=rK}^t \beta_1^{t-j} \widehat{g}_j^m + \beta_1^{t-rK+1} u_r, \quad (\text{D.92})$$

902 which implies

$$\|u_t^m\|_{H_r^{-1}}^2 \leq (1 - \beta_1) \sum_{j=rK}^t \beta_1^{t-j} \|\widehat{g}_j^m\|_{H_r^{-1}}^2 + \beta_1^{t-rK+1} \|u_r\|_{H_r^{-1}}^2. \quad (\text{D.93})$$

903

$$\begin{aligned} \|u_t^m - u_t^n\|_{H_r^{-1}}^2 &\leq (1 - \beta_1) \sum_{j=rK}^t \beta_1^{t-j} \|\widehat{g}_j^m - \widehat{g}_j^n\|_{H_r^{-1}}^2 \\ &\leq 2(1 - \beta_1) \sum_{j=rK}^t \beta_1^{t-j} \left[\|\nabla f(x_j^m) - \nabla f(x_j^n)\|_{H_r^{-1}}^2 + \|\widehat{g}_j^m - \widehat{g}_j^n - [\nabla f(x_j^m) - \nabla f(x_j^n)]\|_{H_r^{-1}}^2 \right]. \end{aligned} \quad (\text{D.94})$$

904 And thus

$$\sum_{j=rK}^t \|u_j^m - u_j^n\|_{H_r^{-1}}^2 \leq 2 \sum_{j=rK}^t \left[\|\nabla f(x_j^m) - \nabla f(x_j^n)\|_{H_r^{-1}}^2 + \|\widehat{g}_j^m - \widehat{g}_j^n - [\nabla f(x_j^m) - \nabla f(x_j^n)]\|_{H_r^{-1}}^2 \right]. \quad (\text{D.95})$$

905 Unroll the recursive bound (D.91) and note that $(1 + \frac{1}{K})^K \leq 3$,

$$\begin{aligned} \|z_{t+1}^m - z_{t+1}^n\|_{H_r}^2 &\leq \underbrace{- \sum_{j=rK}^t 2\eta(1 + \frac{1}{K})^{t-j} \langle z_j^m - z_j^n, \widehat{g}_j^m - \widehat{g}_j^n - \mathbb{E}_j[\widehat{g}_j^m - \widehat{g}_j^n] \rangle}_{\textcircled{1}: \text{martingale}} \\ &\quad + \sum_{j=rK}^t (1 + \frac{1}{K})^{t-j} \left[-2\eta \langle x_j^m - x_j^n, \nabla f(x_j^m) - \nabla f(x_j^n) \rangle + \frac{\eta}{L} \|\nabla f(x_j^m) - \nabla f(x_j^n)\|^2 \right] \\ &\quad + 24 \sum_{j=rK}^t \eta^2 \|\widehat{g}_j^m - \widehat{g}_j^n - \nabla f(x_j^m) + \nabla f(x_j^n)\|_{H_r^{-1}}^2 + 72\eta^2 \sum_{j=rK}^t \|u_j^m - u_j^n\|_{H_r^{-1}}^2 \\ &\quad + 195\eta^2 K \frac{(1 - \beta_2)^2 B_1^2}{(1 - \beta_1)^3} \|u_r\|_{H_r^{-1}}^2 + 48\eta^2 K \left(\frac{1 - \beta_2}{1 - \beta_1} \right)^2 B_1^2 \sum_{j=rK}^t \|\widehat{g}_j^m\|_{H_r^{-1}}^2 + \frac{24\eta^2 K^2}{\lambda} \cdot \frac{\|2\sigma\|_{2\alpha}^{2\alpha}}{\rho^{2(\alpha-1)}} \\ &\stackrel{(\text{D.95})}{\leq} \textcircled{1} + \sum_{j=rK}^t (1 + \frac{1}{K})^{t-j} \left[-2\eta \langle x_j^m - x_j^n, \nabla f(x_j^m) - \nabla f(x_j^n) \rangle + \frac{2\eta}{L} \|\nabla f(x_j^m) - \nabla f(x_j^n)\|^2 \right] \\ &\quad + 144 \sum_{j=rK}^t \eta^2 \|\widehat{g}_j^m - \widehat{g}_j^n - \nabla f(x_j^m) + \nabla f(x_j^n)\|_{H_r^{-1}}^2 + 195\eta^2 K \frac{(1 - \beta_2)^2 B_1^2}{(1 - \beta_1)^3} \|u_r\|_{H_r^{-1}}^2 \\ &\quad + 48\eta^2 K \left(\frac{1 - \beta_2}{1 - \beta_1} \right)^2 B_1^2 \sum_{j=rK}^t \|\widehat{g}_j^m\|_{H_r^{-1}}^2 + \frac{24\eta^2 K^2}{\lambda} \cdot \frac{\|2\sigma\|_{2\alpha}^{2\alpha}}{\rho^{2(\alpha-1)}}. \end{aligned} \quad (\text{D.96})$$

906 Note that by definition, $u_r = (1 - \beta_1) \sum_{j=1}^K \beta_1^{j-1} \mathbb{E}_m \widehat{g}_{rK-j}^m + \beta_1^K u_{r-1}$. By Cauchy-Schwarz inequality,

$$\|u_r\| \leq \beta_1^K \|u_{r-1}\| + \sqrt{\sum_{j=1}^K \|\mathbb{E}_m \widehat{g}_{rK-j}^m\|^2 \sum_{j=1}^K (1 - \beta_1)^2 \beta_1^{2(j-1)}}. \quad (\text{D.97})$$

907 Therefore, event $E_{t,2}$ implies

$$\|u_r\|^2 \leq \frac{(1-\beta_1)^2 \sigma^2 A}{2^{12}(1-\beta_2)^2 B_1^2} \cdot \frac{1-\beta_1}{1-\beta_1^K} \leq \frac{(1-\beta_1)^3 \sigma^2 A}{2^{11}(1-\beta_2)^2 B_1^2}. \quad (\text{D.98})$$

908 By Lemma D.5, and $\|\nabla f(x_j^m)\| \leq G$,

$$\begin{aligned} \|z_{t+1}^m - z_{t+1}^n\|_{H_r}^2 &\leq \textcircled{1} \stackrel{\text{Lemma D.5}}{+} 6\eta\tau K \cdot \frac{\eta^2 \sigma^2}{\lambda^2} KA \\ &\quad + \frac{288\eta^2}{\lambda} \sum_{j=rK}^t \left[\|\widehat{g}_j^m - \nabla f(x_j^m)\|^2 + \|\widehat{g}_j^n - \nabla f(x_j^n)\|^2 \right] \\ &\quad + 96\eta^2 K \left(\frac{1-\beta_2}{1-\beta_1} \right)^2 \frac{B_1^2}{\lambda} \sum_{j=rK}^t \left(\|\widehat{g}_j^m - \nabla f(x_j^m)\|^2 + G^2 \right) \\ &\stackrel{(\text{D.98})}{+} \frac{\eta^2 \sigma^2 KA}{10\lambda} + \frac{24\eta^2 K^2}{\lambda} \cdot \frac{\|2\sigma\|_{2\alpha}^2}{\rho^{2(\alpha-1)}} \\ &\leq \textcircled{1} + 6\eta\tau K \cdot \frac{\eta^2 \sigma^2}{\lambda^2} KA \stackrel{\text{Lemma B.2}}{+} \frac{2^{10}\eta^2}{\lambda} K\sigma^2 \\ &\quad + \frac{2^{10}\eta^2}{\lambda} \max_{s \in [M]} \sum_{j=rK}^t \underbrace{\left[\|\widehat{g}_j^s - \nabla f(x_j^s)\|^2 - \mathbb{E}_j[\|\widehat{g}_j^s - \nabla f(x_j^s)\|^2] \right]}_{\textcircled{2}: \text{martingale}} \\ &\quad + 96\eta^2 K^2 \left(\frac{1-\beta_2}{1-\beta_1} \right)^2 \frac{B_1^2}{\lambda} G^2 + \frac{\eta^2 \sigma^2 KA}{10\lambda} + \frac{24\eta^2 K^2}{\lambda} \cdot \frac{\|2\sigma\|_{2\alpha}^2}{\rho^{2(\alpha-1)}}. \end{aligned} \quad (\text{D.99})$$

909 Define

$$\zeta_j^{m,n} = \begin{cases} -2\eta(1 + \frac{1}{K})^{t-j} \left\langle z_j^m - z_j^n, \widehat{g}_j^m - \widehat{g}_j^n - \mathbb{E}_j[\widehat{g}_j^m - \widehat{g}_j^n] \right\rangle, & \text{if event } E_j \text{ holds,} \\ 0, & \text{otherwise.} \end{cases} \quad (\text{D.100})$$

910

$$\theta_j^m = \begin{cases} \|\widehat{g}_j^m - \nabla f(x_j^m)\|^2 - \mathbb{E}_j[\|\widehat{g}_j^m - \nabla f(x_j^m)\|^2], & \text{if event } E_j \text{ holds,} \\ 0, & \text{otherwise.} \end{cases} \quad (\text{D.101})$$

911 Then (D.99) implies $\|z_{t+1}^m - z_{t+1}^n\|_{H_r}^2 \leq \frac{\eta^2 \sigma^2}{2\lambda} KA + \sum_{j=rK}^t \zeta_j^{m,n} + \frac{2^{10}\eta^2}{\lambda} \max_{s \in [M]} \sum_{j=rK}^t \theta_j^s$. Note that
912 by Lemma B.2,

$$|\theta_j^m| \leq 4\rho^2 d \stackrel{\text{def}}{=} c. \quad (\text{D.102})$$

913

$$\text{Var}_j(\theta_j^m) \leq \mathbb{E}_j[\|\widehat{g}_j^m - \nabla f(x_j^m)\|^4] \leq \sigma^4. \quad (\text{D.103})$$

914 Let $b = \frac{\sigma^2 KA}{2^{12}}$, $V = \sigma^4 K$. Then by Lemma B.1, $|\sum_{j=rK}^t \theta_j^m| \leq b$ with probability no less than

$$1 - 2 \exp\left(\frac{b^2}{2V + 2cb/3} \right) \geq 1 - \frac{\delta}{8MT}. \quad (\text{D.104})$$

915 This implies with probability no less than $1 - \frac{\delta}{8T}$,

$$\left| \sum_{j=rK}^t \theta_j^m \right| \leq \frac{\sigma^2 KA}{2^{12}}, \forall m \in [M]. \quad (\text{D.105})$$

916 Also note that

$$|\zeta_j^{m,n}| \leq 6\eta \cdot \frac{\eta\sigma}{\lambda} \sqrt{KA} \cdot 4\rho\sqrt{d} = \frac{24\eta^2 \sigma \rho \sqrt{d}}{\lambda} \sqrt{KA} \stackrel{\text{def}}{=} c. \quad (\text{D.106})$$

917

$$\text{Var}_j(\zeta_j^{m,n}) \leq \left(6\eta \cdot \frac{\eta\sigma}{\lambda} \sqrt{KA}\right)^2 \cdot 2\sigma^2 = \frac{72\eta^4\sigma^4}{\lambda^2} KA. \quad (\text{D.107})$$

918 Let $b = \frac{\eta^2\sigma^2}{4\lambda} KA, V = \frac{72\eta^4\sigma^4}{\lambda^2} K^2A$. Then by Lemma B.1, $|\sum_{j=rK}^t \zeta_j^{m,n}| \leq b$ with probability no

919 less than

$$1 - 2 \exp\left(\frac{b^2}{2V + 2cb/3}\right) \geq 1 - \frac{\delta}{8M^2T}. \quad (\text{D.108})$$

920 This implies with probability no less than $1 - \frac{\delta}{8T}$,

$$|\sum_{j=rK}^t \zeta_j^{m,n}| \leq \frac{\eta^2\sigma^2}{4\lambda} KA, \forall m, n \in [M]. \quad (\text{D.109})$$

921 We now turn to deal with $\sum_{j=rK}^t \|\widehat{g}_j^m\|^2$.

$$\begin{aligned} \sum_{j=rK}^t \|\widehat{g}_j^m\|^2 &\leq 2 \sum_{j=rK}^t [\|\widehat{g}_j^m - \nabla f(x_j^m)\|^2 + \|\nabla f(x_j^m)\|^2] \\ &\leq 2 \sum_{j=rK}^t \left[\|\widehat{g}_j^m - \nabla f(x_j^m)\|^2 - \mathbb{E}_j[\|\widehat{g}_j^m - \nabla f(x_j^m)\|^2] \right] + 2 \sum_{j=rK}^t \mathbb{E}_j[\|\widehat{g}_j^m - \nabla f(x_j^m)\|^2] + 2KG^2 \\ &\stackrel{\text{Lemma B.2}}{\leq} 2 \sum_{j=rK}^t \left[\|\widehat{g}_j^m - \nabla f(x_j^m)\|^2 - \mathbb{E}_j[\|\widehat{g}_j^m - \nabla f(x_j^m)\|^2] \right] + 2K(\sigma^2 + G^2). \end{aligned} \quad (\text{D.110})$$

922 Then $\sum_{j=rK}^t \|\widehat{g}_j^m\|^2 \leq 2 \sum_{j=rK}^t \theta_j^m + 2K(\sigma^2 + G^2)$ under event E_t . Therefore, by (D.105),

$$\sum_{j=rK}^t \|\widehat{g}_j^m\|^2 \leq \frac{\sigma^2 KA}{2^{11}} + 2K(\sigma^2 + G^2) \leq \frac{(1 - \beta_1)^2 \sigma^2 A}{2^{12}(1 - \beta_2)^2 B_1^2}. \quad (\text{D.111})$$

923 In conclusion, combining (D.105), (D.109), (D.111), we have

$$\mathbb{P} \left\{ E_{t,2} \text{ and } \|z_{t+1}^m - z_{t+1}^n\|_{H_r}^2 \leq \frac{\eta^2 \sigma^2 KA}{\lambda}, \sum_{j=rK}^t \|\widehat{g}_j^m\|^2 \leq \frac{(1 - \beta_1)^2 \sigma^2 A}{2^{12}(1 - \beta_2)^2 B_1^2} \text{ for all } m, n \right\} \geq \mathbb{P}(E_{t,2}) - \frac{\delta}{4T}. \quad (\text{D.112})$$

924

□

925 D.4 Proof of Descent Lemma

926 After laying all the groundwork above, we are now in the position of showing the main descent
927 lemma.

928 **Lemma D.11.** Assume that $\rho \geq \max\{3\sigma_\infty, 2G_\infty\}$ and

$$\begin{aligned} \frac{\eta\sigma^2}{\lambda M} \log \frac{T}{\delta} &\lesssim \Delta, \quad \frac{\eta\rho\sqrt{d}}{(1 - \beta_1)\sqrt{\gamma\lambda}} \log^{\frac{1}{2}} \frac{T}{\delta} \lesssim \sqrt{\Delta}, \quad \frac{\left(\frac{\eta L}{\lambda}\right)^3 \log \frac{T}{\delta}}{(1 - \beta_1)(\sqrt{\beta_2} - \beta_1)} \lesssim \frac{L\Delta}{\rho^2 d}, \\ \left(\frac{\eta L}{\lambda}\right)^3 \sigma^2 KA &\lesssim \frac{L\Delta}{T}, \quad \frac{\eta^2 \sigma^2}{\lambda \gamma M} \lesssim \frac{\Delta}{T}, \quad \frac{\eta \|2\sigma\|_{2\alpha}^{2\alpha}}{\lambda \rho^{2(\alpha-1)}} \lesssim \frac{\Delta}{T}, \end{aligned} \quad (\text{D.113})$$

929 and

$$(1 - \beta_2)B \leq \frac{\eta}{4\gamma} \leq \frac{\eta L}{4\lambda}, \quad \frac{\eta L}{\lambda} \leq \frac{(1 - \beta_1)^2}{2^6}. \quad (\text{D.114})$$

930 Then the following holds:

$$\mathbb{P}(E_{t+1}) \geq \mathbb{P}(E_{t,3}) - \frac{\delta}{4T}. \quad (\text{D.115})$$

931 *Proof.* For any $x \in \mathbb{R}^d$, since $\nabla^2 f(\cdot) \succeq -\tau I_d$ and $H_r \succeq \lambda I_d$, $y \mapsto f(y) + \frac{1}{2\gamma} \|x - y\|_{H_r}^2$
 932 is $(\frac{1}{\gamma} - \frac{\tau}{\lambda})$ -convex with respect to $\|\cdot\|_{H_r}$. Note that under event E_t , $\bar{z}_t \in \Omega_0$. Let $y_t :=$
 933 $\arg \min_y f(y) + \frac{1}{2\gamma} \|\bar{z}_t - y\|_{H_r}^2$ and by Lemma D.4, $y_t \in \Omega_0$. Then

$$f(y_t) + \frac{1}{2\gamma} \|y_t - \bar{z}_t\|_{H_r}^2 \leq f(\bar{z}_{t+1}) + \frac{1}{2\gamma} \|\bar{z}_{t+1} - \bar{z}_t\|_{H_r}^2 - \frac{1}{2} \left(\frac{1}{\gamma} - \frac{\tau}{\lambda}\right) \|\bar{z}_{t+1} - y_t\|_{H_r}^2. \quad (\text{D.116})$$

934 Recall that the definition of $\{z_t^m\}$ implies

$$\begin{aligned} z_{t+1}^m - z_t^m &= -\frac{\eta(H_t^m)^{-1}u_t^m}{1 - \beta_1} + \frac{\eta\beta_1(H_{t-1}^m)^{-1}u_{t-1}^m}{1 - \beta_1} \\ &= -\frac{\eta\beta_1}{1 - \beta_1} [(H_t^m)^{-1} - (H_{t-1}^m)^{-1}]u_{t-1}^m - \eta(H_t^m)^{-1}\widehat{g}_t^m \\ &= -\eta(H_t^m)^{-1}(\widehat{g}_t^m + e_t^m). \end{aligned} \quad (\text{D.117})$$

935 Here $e_t^m = \frac{\beta_1}{1 - \beta_1} (I_d - H_t^m(H_{t-1}^m)^{-1})u_{t-1}^m$.

936 Also, since $\|\bar{z}_{t+1} - \bar{z}_t\| \leq \frac{(1 + \beta_1)\eta\rho\sqrt{d}}{(1 - \beta_1)\lambda} \leq \sqrt{\frac{\Delta\gamma}{160\lambda}} = R_0$, we have $\bar{z}_{t+1} \in \Omega$ and

$$\begin{aligned} f(\bar{z}_{t+1}) - f(y_t) &\leq f(\bar{z}_t) + \langle \nabla f(\bar{z}_t), \bar{z}_{t+1} - \bar{z}_t \rangle + \frac{L}{2} \|\bar{z}_{t+1} - \bar{z}_t\|^2 - f(y_t) \\ &\leq \langle \nabla f(\bar{z}_t), \bar{z}_{t+1} - y_t \rangle + \frac{\tau}{2} \|\bar{z}_t - y_t\|^2 + \frac{L}{2} \|\bar{z}_{t+1} - \bar{z}_t\|^2 \\ &\leq \langle \nabla f(\bar{z}_t), \bar{z}_{t+1} - y_t \rangle + \frac{\tau}{2\lambda} \|\bar{z}_t - y_t\|_{H_r}^2 + \frac{L}{2\lambda} \|\bar{z}_{t+1} - \bar{z}_t\|_{H_r}^2. \end{aligned} \quad (\text{D.118})$$

937 Combine this with (D.116),

$$\begin{aligned} &\frac{\frac{1}{\eta} + \frac{1}{\gamma} - \frac{\tau}{\lambda}}{2} \|\bar{z}_{t+1} - y_t\|_{H_r}^2 - \frac{\frac{1}{\eta} - \frac{1}{\gamma} + \frac{\tau}{\lambda}}{2} \|\bar{z}_t - y_t\|_{H_r}^2 + \frac{\frac{1}{\eta} + \frac{1}{\gamma} - \frac{L}{\lambda}}{2} \|\bar{z}_{t+1} - \bar{z}_t\|_{H_r}^2 \\ &\leq \left\langle \bar{z}_{t+1} - y_t, \nabla f(\bar{z}_t) + \frac{H_r(\bar{z}_{t+1} - \bar{z}_t)}{\eta} \right\rangle \\ &= \left\langle \bar{z}_t - \eta\mathbb{E}_m[(H_t^m)^{-1}(\widehat{g}_t^m + e_t^m)] - y_t, \nabla f(\bar{z}_t) - H_r\mathbb{E}_m[(H_t^m)^{-1}(\widehat{g}_t^m + e_t^m)] \right\rangle \\ &= \left\langle \bar{z}_t - \eta H_r^{-1} \nabla f(\bar{z}_t) - y_t, \nabla f(\bar{z}_t) - H_r\mathbb{E}_m[(H_t^m)^{-1}(\widehat{g}_t^m + e_t^m)] \right\rangle \\ &\quad + \eta \|\nabla f(\bar{z}_t) - H_r\mathbb{E}_m[(H_t^m)^{-1}(\widehat{g}_t^m + e_t^m)]\|_{H_r^{-1}}^2 \\ &\leq \left\langle \bar{z}_t - \eta H_r^{-1} \nabla f(\bar{z}_t) - y_t, \nabla f(\bar{z}_t) - H_r\mathbb{E}_m[(H_t^m)^{-1}(\widehat{g}_t^m + e_t^m)] \right\rangle \\ &\quad + 4\eta \|\nabla f(\bar{z}_t) - \mathbb{E}_m[\nabla f(x_t^m)]\|_{H_r^{-1}}^2 + 4\eta \|\mathbb{E}_m[\nabla f(x_t^m) - \widehat{g}_t^m]\|_{H_r^{-1}}^2 \\ &\quad + 4\eta \|\mathbb{E}_m[(H_r(H_t^m)^{-1} - I_d)\widehat{g}_t^m]\|_{H_r^{-1}}^2 + 4\eta \|\mathbb{E}_m[(H_t^m)^{-1}e_t^m]\|_{H_r}^2. \end{aligned} \quad (\text{D.119})$$

938 By Lemma D.4, we have

$$\begin{aligned}
& \left\langle \bar{z}_t - \eta H_r^{-1} \nabla f(\bar{z}_t) - y_t, \nabla f(\bar{z}_t) - H_r \mathbb{E}_m[(H_t^m)^{-1} \widehat{g}_t^m] \right\rangle \\
&= \left\langle \bar{z}_t - \eta H_r^{-1} \nabla f(\bar{z}_t) - y_t, \nabla f(\bar{z}_t) - \mathbb{E}_m[\nabla f(x_t^m)] \right\rangle \\
&\quad + \left\langle \bar{z}_t - \eta H_r^{-1} \nabla f(\bar{z}_t) - y_t, \mathbb{E}_m[\nabla f(x_t^m) - \widehat{g}_t^m] \right\rangle \\
&\quad + \left\langle \bar{z}_t - \eta H_r^{-1} \nabla f(\bar{z}_t) - y_t, \mathbb{E}_m[(I_d - H_r(H_t^m)^{-1}) \widehat{g}_t^m] \right\rangle \\
&\stackrel{(D.44)}{\leq} \frac{\gamma}{16} \|\nabla f(y_t)\|_{H_r^{-1}}^2 + 8\gamma \|\nabla f(\bar{z}_t) - \mathbb{E}_m[\nabla f(x_t^m)]\|_{H_r^{-1}}^2 + 8\gamma \left\| \mathbb{E}_m[(H_r(H_t^m)^{-1} - I_d) \widehat{g}_t^m] \right\|_{H_r^{-1}}^2 \\
&\quad + \left\langle \bar{z}_t - \eta H_r^{-1} \nabla f(\bar{z}_t) - y_t, \mathbb{E}_m[\nabla f(x_t^m) - \widehat{g}_t^m] \right\rangle.
\end{aligned} \tag{D.120}$$

939 Also,

$$\left\langle \bar{z}_t - \eta H_r^{-1} \nabla f(\bar{z}_t) - y_t, -H_r \mathbb{E}_m[(H_t^m)^{-1} e_t^m] \right\rangle \leq \frac{\gamma}{16} \|\nabla f(y_t)\|_{H_r^{-1}}^2 + 4\gamma \left\| \mathbb{E}_m[(H_t^m)^{-1} e_t^m] \right\|_{H_r}^2 \tag{D.121}$$

940 Further noticing that $\eta \leq \frac{\gamma}{4}$ and by AM-GM inequality, we conclude that

LHS of (D.119)

$$\begin{aligned}
& \leq \frac{\gamma}{8} \|\nabla f(y_t)\|_{H_r^{-1}}^2 + 9\gamma \|\nabla f(\bar{z}_t) - \mathbb{E}_m[\nabla f(x_t^m)]\|_{H_r^{-1}}^2 + 9\gamma \left\| \mathbb{E}_m[(H_r(H_t^m)^{-1} - I_d) \widehat{g}_t^m] \right\|_{H_r^{-1}}^2 \\
&\quad + 4\eta \left\| \mathbb{E}_m[\nabla f(x_t^m) - \widehat{g}_t^m] \right\|_{H_r^{-1}}^2 + 5\gamma \left\| \mathbb{E}_m[(H_t^m)^{-1} e_t^m] \right\|_{H_r}^2 \\
&\quad + \left\langle \bar{z}_t - \eta H_r^{-1} \nabla f(\bar{z}_t) - y_t, \mathbb{E}_m[\nabla f(x_t^m) - \widehat{g}_t^m] \right\rangle.
\end{aligned} \tag{D.122}$$

941 If $t \bmod K \equiv -1$, then $r(t+1) = r(t) + 1 = r + 1$ and event $E_{t,1}$ implies

$$H_r^{-1} H_{r+1} \leq 1 + (1 - \beta_2) B \leq 1 + \frac{\eta}{4\gamma}, \tag{D.123}$$

942

$$\begin{aligned}
f_\gamma^{H_{r+1}}(\bar{z}_{t+1}) &\leq f(y_t) + \frac{1}{2\gamma} \|\bar{z}_{t+1} - y_t\|_{H_{r+1}}^2 \\
&\leq f(y_t) + \frac{1 + \eta/4\gamma}{2\gamma} \|\bar{z}_{t+1} - y_t\|_{H_r}^2.
\end{aligned} \tag{D.124}$$

943 On the other hand, if $t \bmod K \not\equiv -1$, then $r(t+1) = r(t) = r$,

$$f_\gamma^{H_{r(t+1)}}(\bar{z}_{t+1}) \leq f(y_t) + \frac{1}{2\gamma} \|\bar{z}_{t+1} - y_t\|_{H_r}^2. \tag{D.125}$$

944 Hence the following always holds:

$$\begin{aligned}
f_\gamma^{H_r(t+1)}(\bar{z}_{t+1}) &\leq f_\gamma^{H_r}(\bar{z}_t) - \frac{1}{2\gamma} \|\bar{z}_t - y_t\|_{H_r}^2 + \frac{1 + \eta/4\gamma}{2\gamma} \|\bar{z}_{t+1} - y_t\|_{H_r}^2 \\
&\stackrel{(D.122)}{\leq} f_\gamma^{H_r}(\bar{z}_t) - \frac{7\gamma^{-1}}{8\gamma(\eta^{-1} + \gamma^{-1})} \|\bar{z}_t - y_t\|_{H_r}^2 \\
&\quad + \frac{(1 + \eta/4\gamma) \left[\frac{1}{8} \|\nabla f(y_t)\|_{H_r^{-1}}^2 + 9 \|\nabla f(\bar{z}_t) - \mathbb{E}_m[\nabla f(x_t^m)]\|_{H_r^{-1}}^2 + 9 \left\| \mathbb{E}_m[(H_r(H_t^m)^{-1} - I_d)\widehat{g}_t^m] \right\|_{H_r^{-1}}^2 \right]}{\eta^{-1} + \gamma^{-1} - \tau/\lambda} \\
&\quad + \frac{(1 + \eta/4\gamma) \left[4\eta \left\| \mathbb{E}_m[\nabla f(x_t^m) - \widehat{g}_t^m] \right\|_{H_r^{-1}}^2 + 5\gamma \left\| \mathbb{E}_m[(H_t^m)^{-1} e_t^m] \right\|_{H_r}^2 \right]}{\gamma(\eta^{-1} + \gamma^{-1} - \tau/\lambda)} \\
&\quad + \frac{(1 + \eta/4\gamma) \left\langle \bar{z}_t - \eta H_r^{-1} \nabla f(\bar{z}_t) - y_t, \mathbb{E}_m[\nabla f(x_t^m) - \widehat{g}_t^m] \right\rangle}{\gamma(\eta^{-1} + \gamma^{-1} - \tau/\lambda)} \\
&\stackrel{(D.42)}{\leq} f_\gamma^{H_r}(\bar{z}_t) - \frac{\eta}{8} \|\nabla f(y_t)\|_{H_r^{-1}}^2 + \frac{5\eta^2}{\lambda\gamma} \left\| \mathbb{E}_m[\nabla f(x_t^m) - \widehat{g}_t^m] \right\|^2 + 6\eta \left\| \mathbb{E}_m[(H_t^m)^{-1} e_t^m] \right\|_{H_r}^2 \\
&\quad + \frac{10\eta}{\lambda} \|\nabla f(\bar{z}_t) - \mathbb{E}_m[\nabla f(x_t^m)]\|^2 + 10\eta \left\| \mathbb{E}_m[(H_r(H_t^m)^{-1} - I_d)\widehat{g}_t^m] \right\|_{H_r^{-1}}^2 \\
&\quad + \frac{1 + \eta/4\gamma}{\gamma(\eta^{-1} + \gamma^{-1} - \tau/\lambda)} \left\langle \bar{z}_t - \eta H_r^{-1} \nabla f(\bar{z}_t) - y_t, \mathbb{E}_m[\nabla f(x_t^m) - \widehat{g}_t^m] \right\rangle.
\end{aligned} \tag{D.126}$$

945 Sum over t and we get

$$\begin{aligned}
f_\gamma^{H_r(t+1)}(\bar{z}_{t+1}) &\leq f_\gamma^\lambda(x_0) - \frac{\eta}{8} \sum_{j=0}^t \|\nabla f(y_j)\|_{H_r^{-1}(j)}^2 + \frac{5\eta^2}{\lambda\gamma} \sum_{j=0}^t \left\| \mathbb{E}_m[\nabla f(x_j^m) - \widehat{g}_j^m] \right\|^2 + 6\eta \sum_{j=0}^t \left\| \mathbb{E}_m[(H_j^m)^{-1} e_j^m] \right\|_{H_r(j)}^2 \\
&\quad + \frac{10\eta}{\lambda} \sum_{j=0}^t \|\nabla f(\bar{z}_j) - \mathbb{E}_m[\nabla f(x_j^m)]\|^2 + 10\eta \sum_{j=0}^t \left\| \mathbb{E}_m[(H_{r(j)}(H_j^m)^{-1} - I_d)\widehat{g}_j^m] \right\|_{H_r^{-1}(j)}^2 \\
&\quad + \underbrace{\frac{1 + \eta/4\gamma}{\gamma(\eta^{-1} + \gamma^{-1} - \tau/\lambda)} \sum_{j=0}^t \left\langle \bar{z}_j - \eta H_{r(j)}^{-1} \nabla f(\bar{z}_j) - y_j, \mathbb{E}_m[\nabla f(x_j^m) - \widehat{g}_j^m] \right\rangle}_{(*)}.
\end{aligned} \tag{D.127}$$

946 By AM-GM inequality and notice that $\bar{x}_t, \bar{z}_t \in \Omega$,

$$\begin{aligned}
&\|\nabla f(\bar{z}_t) - \mathbb{E}_m[\nabla f(x_t^m)]\|^2 \\
&\leq 2\|\nabla f(\bar{z}_t) - \nabla f(\bar{x}_t)\|^2 + 2\|\nabla f(\bar{x}_t) - \mathbb{E}_m[\nabla f(x_t^m)]\|^2 \\
&\leq 2L^2\|\bar{z}_t - \bar{x}_t\|^2 + 2\|\nabla f(\bar{x}_t) - \mathbb{E}_m[\nabla f(x_t^m)]\|^2.
\end{aligned} \tag{D.128}$$

947 Under event $E_{t,3}$,

$$\left\| \mathbb{E}_m[(H_r(H_t^m)^{-1} - I_d)\widehat{g}_t^m] \right\|_{H_r^{-1}}^2 \leq (1 - \beta_2)^2 B^2 \mathbb{E}_m \left[\|\widehat{g}_t^m\|_{H_r^{-1}}^2 \right]. \tag{D.129}$$

948

$$\left\| \mathbb{E}_m[(H_t^m)^{-1} e_t^m] \right\|_{H_r}^2 \leq 4 \left(\frac{\beta_1(1 - \beta_2)}{1 - \beta_1} \right)^2 B^2 \mathbb{E}_m \left[\|u_{t-1}^m\|_{H_r^{-1}}^2 \right]. \tag{D.130}$$

949 By the definition of u_{t-1}^m , we have

$$\begin{aligned}
\mathbb{E}_m \left[\|u_{t-1}^m\|_{H_r^{-1}}^2 \right] &\leq (1 - \beta_1) \sum_{j=0}^{t-1} \beta_1^{t-j-1} \mathbb{E}_m \left[\|\widehat{g}_j^m\|_{H_r^{-1}}^2 \right] \\
&\leq \frac{(1 - \beta_1)}{\beta_2^{K/2}} \sum_{j=0}^{t-1} (\beta_1/\sqrt{\beta_2})^{t-j-1} \mathbb{E}_m \left[\|\widehat{g}_j^m\|_{H_r^{-1}(j)}^2 \right].
\end{aligned} \tag{D.131}$$

950 Plug these inequalities above in (D.127),

$$\begin{aligned}
f_\gamma^{H_r(t+1)}(\bar{z}_{t+1}) &\leq f_\gamma^\lambda(x_0) - \frac{\eta}{8} \sum_{j=0}^t \|\nabla f(y_j)\|_{H_r(j)}^2 + \frac{5\eta^2}{\lambda\gamma} \sum_{j=0}^t \|\mathbb{E}_m[\nabla f(x_j^m) - \widehat{g}_j^m]\|^2 \\
&\stackrel{(D.128)}{+} \frac{20\eta}{\lambda} \sum_{j=0}^t [L^2 \|\bar{z}_j - \bar{x}_j\|^2 + \|\nabla f(\bar{x}_j) - \mathbb{E}_m[\nabla f(x_j^m)]\|^2] \\
&\stackrel{(D.129)-(D.131)}{+} \eta \left(\frac{48\beta_1^2}{(1-\beta_1)(\sqrt{\beta_2}-\beta_1)} + 10 \right) (1-\beta_2)^2 B^2 \sum_{j=0}^t \mathbb{E}_m \left[\|\widehat{g}_j^m\|_{H_r(j)}^2 \right] + (*).
\end{aligned} \tag{D.132}$$

951 By AM-GM inequality and Lemma D.4,

$$\begin{aligned}
\mathbb{E}_m \left[\|\widehat{g}_t^m\|_{H_r}^2 \right] &\leq 4\mathbb{E}_m \left[\|\widehat{g}_t^m - \nabla f(x_t^m)\|_{H_r}^2 + \|\nabla f(x_t^m) - \nabla f(\bar{x}_t)\|_{H_r}^2 \right. \\
&\quad \left. + \|\nabla f(\bar{x}_t) - \nabla f(\bar{z}_t)\|_{H_r}^2 + \|\nabla f(\bar{z}_t)\|_{H_r}^2 \right] \\
&\leq \frac{4}{\lambda} \left[\mathbb{E}_m \|\widehat{g}_t^m - \nabla f(x_t^m)\|^2 + L^2 \mathbb{E}_m [\|x_t^m - \bar{x}_t\|^2] + L^2 \|\bar{z}_t - \bar{x}_t\|^2 \right] + \frac{16(\gamma L)^2}{\lambda^2} \|\nabla f^{H_r}(\bar{z}_t)\|_{H_r}^2.
\end{aligned} \tag{D.133}$$

952 Therefore, we achieve that

$$\begin{aligned}
f_\gamma^{H_r(t+1)}(\bar{z}_{t+1}) &\leq f_\gamma^{H_0}(x_0) - \frac{\eta}{9} \sum_{j=0}^t \|\nabla f(y_j)\|_{H_r(j)}^2 + \frac{5\eta^2}{\lambda\gamma} \sum_{j=0}^t \|\mathbb{E}_m[\nabla f(x_j^m) - \widehat{g}_j^m]\|^2 \\
&\quad + \frac{40\eta}{\lambda} \sum_{j=0}^t [L^2 \|\bar{z}_j - \bar{x}_j\|^2 + \|\nabla f(\bar{x}_j) - \mathbb{E}_m[\nabla f(x_j^m)]\|^2] \\
&\quad + \frac{160\eta(1-\beta_2)^2 B^2}{\lambda(1-\beta_1)(\sqrt{\beta_2}-\beta_1)} \sum_{j=0}^t \left[\mathbb{E}_m \|\widehat{g}_j^m - \nabla f(x_j^m)\|^2 + L^2 \mathbb{E}_m [\|x_j^m - \bar{x}_j\|^2] \right] + (*).
\end{aligned} \tag{D.134}$$

953 By (D.160), (D.164) in Lemma D.12, under event $E_{t,3}$,

$$\begin{aligned}
\|\bar{z}_j - \bar{x}_j\|^2 &\leq \left(\frac{\beta_1}{1-\beta_1} \right)^2 \left[64\eta^2 \left(\|\nabla f(\bar{z}_j)\|_{H_r(j)}^2 + \frac{L^2}{\lambda^2} \Lambda_{j-1} \right) \right. \\
&\quad \left. + \frac{36\eta^2}{\lambda^2} (1-\beta_1) \sum_{i=r(j)K}^{j-1} \beta_1^{j-i-1} \left[\frac{\eta^2 L^2 \sigma^2}{\lambda^2} KA + \mathbb{E}_m \|\widehat{g}_i^m - \nabla f(x_i^m)\|^2 \right] \right].
\end{aligned} \tag{D.135}$$

954 Hence

$$\begin{aligned}
\sum_{j=0}^t \|\bar{z}_j - \bar{x}_j\|^2 &\leq \left(\frac{\beta_1}{1-\beta_1} \right)^2 \left[64\eta^2 \sum_{j=0}^t \left(\|\nabla f(\bar{z}_j)\|_{H_r(j)}^2 + \frac{L^2}{\lambda^2} \Lambda_{j-1} \right) \right. \\
&\quad \left. + \frac{36\eta^2}{\lambda^2} \sum_{j=0}^{t-1} \left[\frac{\eta^2 L^2 \sigma^2}{\lambda^2} KA + \mathbb{E}_m \|\widehat{g}_j^m - \nabla f(x_j^m)\|^2 \right] \right].
\end{aligned} \tag{D.136}$$

955 Additionally by Lemma D.12,

$$\begin{aligned}
\Lambda_t + \frac{(1-\beta_1)^2}{2} \sum_{j=0}^{t-1} \Lambda_j &\leq \frac{64\eta^2}{1-\beta_1} \sum_{j=0}^t \|\nabla f(\bar{z}_j)\|_{H_r(j)}^2 \\
&\quad + \frac{36\eta^2}{\lambda^2} (1-\beta_1) \sum_{j=0}^{t-1} \left[\frac{\eta^2 L^2 \sigma^2}{\lambda^2} KA + \mathbb{E}_m \|\widehat{g}_j^m - \nabla f(x_j^m)\|^2 \right].
\end{aligned} \tag{D.137}$$

956 Therefore, by noticing that $\Lambda_t \geq 0$ and $\frac{\eta L}{\lambda} \leq \frac{(1 - \beta_1)^2}{16}$,

$$\sum_{j=0}^t \|\bar{z}_j - \bar{x}_j\|^2 \leq 2 \left(\frac{\eta \beta_1}{1 - \beta_1} \right)^2 \left[64 \sum_{j=0}^t \|\nabla f(\bar{z}_j)\|_{H_{r(j)}^{-2}}^2 + \frac{36}{\lambda^2} \sum_{j=0}^{t-1} \left[\frac{\eta^2 L^2 \sigma^2}{\lambda^2} KA + \mathbb{E}_m \|g_j^m - \nabla f(x_j^m)\|^2 \right] \right] \quad (\text{D.138})$$

957 For the third term of RHS of (D.130),

$$\begin{aligned} \frac{5\eta^2}{\lambda\gamma} \sum_{j=0}^t \|\mathbb{E}_m[\nabla f(x_j^m) - \widehat{g}_j^m]\|^2 &\leq \frac{10\eta^2}{\lambda\gamma} \sum_{j=0}^t \left[\|\mathbb{E}_m[\widehat{g}_j^m - \mathbb{E}_j[\widehat{g}_j^m]]\|^2 + \|\mathbb{E}_m[\nabla f(x_j^m) - \mathbb{E}_j[\widehat{g}_j^m]]\|^2 \right] \\ &\stackrel{\text{Lemma B.2}}{\leq} \frac{10\eta^2}{\lambda\gamma} \sum_{j=0}^t \left[\|\mathbb{E}_m[\widehat{g}_j^m - \mathbb{E}_j[\widehat{g}_j^m]]\|^2 + \frac{\|2\sigma\|_{2\alpha}^{2\alpha}}{\rho^{2(\alpha-1)}} \right] \\ &\leq \frac{10\eta^2}{\lambda\gamma} \sum_{j=0}^t \underbrace{\left[\|\mathbb{E}_m[\widehat{g}_j^m - \mathbb{E}_j[\widehat{g}_j^m]]\|^2 - \mathbb{E}_j \left[\|\mathbb{E}_m[\widehat{g}_j^m - \mathbb{E}_j[\widehat{g}_j^m]]\|^2 \right] \right]}_{\text{①: martingale}} \\ &\quad + \frac{10\eta^2 T}{\lambda\gamma} \left[\frac{\|2\sigma\|_{2\alpha}^{2\alpha}}{\rho^{2(\alpha-1)}} + \frac{\sigma^2}{M} \right] \end{aligned} \quad (\text{D.139})$$

958 For the (*) term of RHS of (D.130),

$$\begin{aligned} &\frac{1 + \eta/4\gamma}{\gamma(\eta^{-1} + \gamma^{-1} - \tau/\lambda)} \sum_{j=0}^t \left\langle \bar{z}_j - \eta H_{r(j)}^{-1} \nabla f(\bar{z}_j) - y_j, \mathbb{E}_m[\nabla f(x_j^m) - \widehat{g}_j^m] \right\rangle \\ &= \frac{1 + \eta/4\gamma}{\gamma(\eta^{-1} + \gamma^{-1} - \tau/\lambda)} \sum_{j=0}^t \left\langle \bar{z}_j - \eta H_{r(j)}^{-1} \nabla f(\bar{z}_j) - y_j, \mathbb{E}_m[\nabla f(x_j^m) - \mathbb{E}_j[\widehat{g}_j^m]] \right\rangle \\ &\quad + \frac{1 + \eta/4\gamma}{\gamma(\eta^{-1} + \gamma^{-1} - \tau/\lambda)} \sum_{j=0}^t \underbrace{\left\langle \bar{z}_j - \eta H_{r(j)}^{-1} \nabla f(\bar{z}_j) - y_j, \mathbb{E}_m[\mathbb{E}_j[\widehat{g}_j^m] - \widehat{g}_j^m] \right\rangle}_{\text{②: martingale}} \end{aligned} \quad (\text{D.140})$$

$$\begin{aligned} &\stackrel{\text{AM-GM}}{\leq} \frac{2\eta}{\gamma} \sum_{j=0}^t \left[\frac{1}{120\gamma} \|H_{r(j)}(\bar{z}_j - y_j) - \eta \nabla f(\bar{z}_j)\|_{H_{r(j)}^{-1}}^2 + 30\gamma \frac{\|2\sigma\|_{2\alpha}^{2\alpha}}{\lambda \rho^{2(\alpha-1)}} \right] + \text{②} \\ &\stackrel{(\text{D.44})}{\leq} \frac{\eta}{60} \sum_{j=0}^t \|\nabla f(y_j)\|_{H_{r(j)}^{-1}}^2 + \frac{60\eta T}{\lambda} \frac{\|2\sigma\|_{2\alpha}^{2\alpha}}{\rho^{2(\alpha-1)}} + \text{②} \end{aligned}$$

959 Here we remark that ② is a martingale because $H_{r(j)}$ only depends on stochastic gradients drawn
960 strictly before round $r(j)$ and thus independent of \widehat{g}_j^m , which is drawn during round $r(j)$.

961 Plug (D.138),(D.139), (D.140) in (D.130),

$$\begin{aligned}
f_\gamma^{H_r(t+1)}(\bar{z}_{t+1}) &\leq f_\gamma^\lambda(x_0) - \frac{\eta}{12} \sum_{j=0}^t \|\nabla f(y_j)\|_{H_r^{-1}(j)}^2 + \textcircled{1} + \frac{10\eta^2 T}{\lambda\gamma} \left[\frac{\|2\sigma\|_{2\alpha}^{2\alpha}}{\rho^{2(\alpha-1)}} + \frac{\sigma^2}{M} \right] \\
&\quad + \frac{40\eta}{\lambda} \sum_{j=0}^t \left[\frac{72(\eta L\beta_1)^2}{(\lambda(1-\beta_1))^2} \left[\frac{\eta^2 L^2 \sigma^2}{\lambda^2} K A + \mathbb{E}_m \|\widehat{g}_j^m - \nabla f(x_j^m)\|^2 \right] + \frac{\eta^2 L^2 \sigma^2}{\lambda^2} K A \right] \\
&\quad + \frac{160\eta(1-\beta_2)^2 B^2}{\lambda(1-\beta_1)(\sqrt{\beta_2}-\beta_1)} \sum_{j=0}^t \left[\mathbb{E}_m \|\widehat{g}_j^m - \nabla f(x_j^m)\|^2 + \frac{\eta^2 L^2 \sigma^2}{\lambda^2} K A \right] \\
&\quad + \frac{60\eta T}{\lambda} \frac{\|2\sigma\|_{2\alpha}^{2\alpha}}{\rho^{2(\alpha-1)}} + \textcircled{2} \\
&\leq f_\gamma^\lambda(x_0) - \frac{\eta}{12} \sum_{j=0}^t \|\nabla f(y_j)\|_{H_r^{-1}(j)}^2 + \textcircled{1} + \frac{10\eta^2 T}{\lambda\gamma} \left[\frac{\|2\sigma\|_{2\alpha}^{2\alpha}}{\rho^{2(\alpha-1)}} + \frac{\sigma^2}{M} \right] \\
&\quad + \frac{160\eta}{\lambda} \frac{[18(\frac{\eta L\beta_1}{\lambda})^2 + (1-\beta_2)^2 B^2]}{(1-\beta_1)(\sqrt{\beta_2}-\beta_1)} \sum_{j=0}^t \left[\mathbb{E}_m \|\widehat{g}_j^m - \nabla f(x_j^m)\|^2 \right] \\
&\quad + \frac{160\eta T}{\lambda} \cdot \left[\frac{1}{4} + \frac{18(\frac{\eta L\beta_1}{\lambda})^2 + (1-\beta_2)^2 B^2}{(1-\beta_1)(\sqrt{\beta_2}-\beta_1)} \right] \cdot \frac{\eta^2 L^2 \sigma^2}{\lambda^2} K A \\
&\quad + \frac{60\eta T}{\lambda} \frac{\|2\sigma\|_{2\alpha}^{2\alpha}}{\rho^{2(\alpha-1)}} + \textcircled{2} \\
&\leq f_\gamma^\lambda(x_0) - \frac{\eta}{12} \sum_{j=0}^t \|\nabla f(y_j)\|_{H_r^{-1}(j)}^2 + \textcircled{1} + \frac{10\eta^2 T}{\lambda\gamma} \left[\frac{\|2\sigma\|_{2\alpha}^{2\alpha}}{\rho^{2(\alpha-1)}} + \frac{\sigma^2}{M} \right] \\
&\quad + \underbrace{\frac{160\eta}{\lambda} \frac{20(\frac{\eta L}{\lambda})^2}{(1-\beta_1)(\sqrt{\beta_2}-\beta_1)} \sum_{j=0}^t \mathbb{E}_m \left[\|\widehat{g}_j^m - \nabla f(x_j^m)\|^2 - \mathbb{E}_j \left[\|\widehat{g}_j^m - \nabla f(x_j^m)\|^2 \right] \right]}_{\textcircled{3}: \text{martingale}} \\
&\quad + \frac{50\eta T}{\lambda} \cdot \frac{\eta^2 L^2 \sigma^2}{\lambda^2} \left(K A + \frac{64}{(1-\beta_1)(\sqrt{\beta_2}-\beta_1)} \right) + \frac{60\eta T}{\lambda} \frac{\|2\sigma\|_{2\alpha}^{2\alpha}}{\rho^{2(\alpha-1)}} + \textcircled{2} \\
&\leq f_\gamma^\lambda(x_0) - \frac{\eta}{12} \sum_{j=0}^t \|\nabla f(y_j)\|_{H_r^{-1}(j)}^2 + \frac{10\eta^2 \sigma^2}{\lambda\gamma M} T + \frac{60\eta T}{\lambda} \cdot \frac{\eta^2 L^2 \sigma^2}{\lambda^2} K A + \frac{60\eta T}{\lambda} \frac{\|2\sigma\|_{2\alpha}^{2\alpha}}{\rho^{2(\alpha-1)}} \\
&\quad + \textcircled{1} + \textcircled{2} + \textcircled{3}. \tag{D.141}
\end{aligned}$$

962 where in the third inequality, we apply $(1-\beta_2)B \leq \frac{\eta L}{\lambda}$.

963 For $\textcircled{1}$, define

$$\theta_j = \begin{cases} \frac{10\eta^2}{\lambda\gamma} \left[\left\| \mathbb{E}_m[\widehat{g}_j^m - \mathbb{E}_j[\widehat{g}_j^m]] \right\|^2 - \mathbb{E}_j \left[\left\| \mathbb{E}_m[\widehat{g}_j^m - \mathbb{E}_j[\widehat{g}_j^m]] \right\|^2 \right] \right], & \text{if event } E_j \text{ holds,} \\ 0, & \text{otherwise.} \end{cases} \tag{D.142}$$

964 Then event E_t implies $\textcircled{1} = \sum_{j=0}^t \theta_j$ and notice that

$$|\theta_j| \leq \frac{10\eta^2}{\lambda\gamma} \cdot 4\rho^2 d = \frac{40\eta^2 \rho^2 d}{\lambda\gamma} \stackrel{\text{def}}{=} c, \tag{D.143}$$

965

$$\text{Var}_j(\theta_j) \leq \left(\frac{10\eta^2}{\lambda\gamma} \right)^2 \mathbb{E}_j \left[\left\| \mathbb{E}_m[\widehat{g}_j^m - \mathbb{E}_j[\widehat{g}_j^m]] \right\|^2 \right] \stackrel{\text{Lemma B.3}}{\leq} 1600 \left(\frac{\eta^2 \sigma^2}{\lambda\gamma M} \right)^2. \tag{D.144}$$

966 Let $b = \Delta/4$, $V = 1600T \left(\frac{\eta^2 \sigma^2}{\lambda \gamma M} \right)^2$. Then by Lemma B.1, $|\sum_{j=0}^t \theta_j| \leq b$ with probability no less

967 than

$$1 - 2 \exp\left(-\frac{b^2}{2V + 2cb/3}\right) \geq 1 - \frac{\delta}{12T}. \quad (\text{D.145})$$

968 For ③, define

$$\xi_j = \begin{cases} \frac{160\eta}{\lambda} \frac{20(\frac{\eta L}{\lambda})^2}{(1-\beta_1)(\sqrt{\beta_2}-\beta_1)} \left(\mathbb{E}_m \left[\|\widehat{g}_j^m - \nabla f(x_j^m)\|^2 - \mathbb{E}_j[\|\widehat{g}_j^m - \nabla f(x_j^m)\|^2] \right] \right), & \text{if event } E_j \text{ holds,} \\ 0, & \text{otherwise.} \end{cases} \quad (\text{D.146})$$

969 Note that

$$|\xi_j| \leq \frac{160\eta}{\lambda} \frac{20(\frac{\eta L}{\lambda})^2}{(1-\beta_1)(\sqrt{\beta_2}-\beta_1)} \cdot 4\rho^2 d \stackrel{def}{=} c \quad (\text{D.147})$$

970

$$\begin{aligned} \text{Var}_j(\xi_j) &\leq \left(\frac{160\eta}{\lambda} \frac{20(\frac{\eta L}{\lambda})^2}{(1-\beta_1)(\sqrt{\beta_2}-\beta_1)} \right)^2 \frac{\mathbb{E}_j \mathbb{E}_m \|\widehat{g}_j^m - \nabla f(x_j^m)\|^4}{M} \\ &\leq \left(\frac{160\eta}{\lambda} \frac{20(\frac{\eta L}{\lambda})^2}{(1-\beta_1)(\sqrt{\beta_2}-\beta_1)} \right)^2 \frac{\sigma^4}{M}. \end{aligned} \quad (\text{D.148})$$

971 Let $b = \Delta/4$, $V = \left(\frac{160\eta}{\lambda} \frac{20(\frac{\eta L}{\lambda})^2}{(1-\beta_1)(\sqrt{\beta_2}-\beta_1)} \right)^2 \frac{T\sigma^4}{M}$. Then by Lemma B.1, $|\sum_{j=0}^t \xi_j| \leq b$ with

972 probability no less than

$$1 - 2 \exp\left(-\frac{b^2}{2V + 2cb/3}\right) \geq 1 - \frac{\delta}{12T}. \quad (\text{D.149})$$

973 For ②, define

$$\zeta_j = \begin{cases} \frac{1 + \eta/4\gamma}{\gamma(\eta^{-1} + \gamma^{-1} - \tau/\lambda)} \left\langle \bar{z}_j - \eta H_{r(j)}^{-1} \nabla f(\bar{z}_j) - y_j, \mathbb{E}_m[\mathbb{E}_j[\widehat{g}_j^m] - \widehat{g}_j^m] \right\rangle, & \text{if event } E_j \text{ holds,} \\ 0, & \text{otherwise.} \end{cases} \quad (\text{D.150})$$

974 Then event E_t implies ② = $\sum_{j=0}^t \zeta_j$ and notice that by Lemma D.4,

$$\begin{aligned} \|\bar{z}_j - \eta H_{r(j)}^{-1} \nabla f(\bar{z}_j) - y_j\|^2 &\leq \frac{\|H_{r(j)}(\bar{z}_j - y_j) - \eta \nabla f(\bar{z}_j)\|_{H_{r(j)}^{-1}}^2}{\lambda} \\ &\leq \frac{\gamma^2 \|\nabla f_{\gamma}^{H_{r(j)}}(\bar{z}_j)\|_{H_{r(j)}^{-1}}^2}{\lambda} \\ &\leq \frac{2\gamma\Delta}{\lambda}. \end{aligned} \quad (\text{D.151})$$

975 Therefore,

$$|\zeta_j| \leq \frac{2\eta}{\gamma} \cdot \sqrt{\frac{2\gamma\Delta}{\lambda}} \cdot 2\rho\sqrt{d} = 4\eta\rho\sqrt{\frac{2\Delta d}{\gamma\lambda}} \stackrel{def}{=} c, \quad (\text{D.152})$$

976

$$\text{Var}_j(\zeta_j) \leq \left(\frac{2\eta}{\gamma} \right)^2 \cdot \frac{\gamma^2}{\lambda} \|\nabla f(y_j)\|_{H_{r(j)}^{-1}}^2 \cdot \frac{\sigma^2}{M} \leq \frac{4\eta^2 \sigma^2}{\lambda M} \|\nabla f(y_j)\|_{H_{r(j)}^{-1}}^2. \quad (\text{D.153})$$

977 Let $b = \Delta/4$, $V = \frac{100\eta\sigma^2\Delta}{\lambda M}$. Then by Lemma B.1,

$$\mathbb{P} \left\{ \left| \sum_{j=0}^t \zeta_j \right| > b \text{ and } \sum_{j=0}^t \text{Var}_j(\zeta_j) \leq V \right\} \leq 2 \exp\left(-\frac{b^2}{2V + 2cb/3}\right) \leq \frac{\delta}{12T}. \quad (\text{D.154})$$

978 Note that by Lemma D.4 and event E_t ,

$$\|\nabla f(y_t)\|_{H_r^{-1}}^2 \leq \frac{2}{\gamma} (f_\gamma^{H_r(t)}(\bar{z}_t) - \min f_\gamma^\lambda) \leq \frac{4\Delta}{\gamma}. \quad (\text{D.155})$$

979

$$\sum_{j=0}^t \text{Var}_j(\zeta_j) \leq \frac{4\eta^2\sigma^2}{\lambda M} \sum_{j=0}^t \|\nabla f(y_j)\|_{H_r^{-1}}^2 \leq \frac{4\eta^2\sigma^2}{\lambda M} \cdot \left(\frac{24\Delta}{\eta} + \frac{4\Delta}{\gamma}\right) \leq V. \quad (\text{D.156})$$

980 Therefore, combining ①, ②, ③, with probability no less than $\mathbb{P}(E_{t,3}) - 3 \cdot \frac{\delta}{12T}$, event $E_{t,3}$ holds

981 and $|\sum_{j=0}^t \zeta_j| \leq \frac{\Delta}{4}$, $|\sum_{j=0}^t \theta_j| \leq \frac{\Delta}{4}$, $|\sum_{j=0}^t \xi_j| \leq \frac{\Delta}{4}$. These implies

$$\begin{aligned} f_\gamma^{H_r(t+1)}(\bar{z}_{t+1}) - \min f_\gamma^\lambda &\leq \frac{7}{4}\Delta - \frac{\eta}{12} \sum_{j=0}^t \|\nabla f(y_j)\|_{H_r^{-1}}^2 + \frac{10\eta^2\sigma^2}{\lambda\gamma M}T + \frac{60\eta T}{\lambda} \cdot \frac{\eta^2 L^2 \sigma^2}{\lambda^2} KA + \frac{60\eta T}{\lambda} \frac{\|2\sigma\|_{2\alpha}^{2\alpha}}{\rho^{2(\alpha-1)}} \\ &\leq 2\Delta - \frac{\eta}{12} \sum_{j=0}^t \|\nabla f_\gamma^{H_r(j)}(\bar{z}_j)\|_{H_r^{-1}}^2. \end{aligned} \quad (\text{D.157})$$

982 In the last inequality, we apply

$$\frac{10\eta^2\sigma^2}{\lambda\gamma M}T \leq \frac{\Delta}{12}, \quad \frac{60\eta T}{\lambda} \cdot \frac{\eta^2 L^2 \sigma^2}{\lambda^2} KA \leq \frac{\Delta}{12}, \quad \frac{60\eta T}{\lambda} \frac{\|2\sigma\|_{2\alpha}^{2\alpha}}{\rho^{2(\alpha-1)}} \leq \frac{\Delta}{12} \quad (\text{D.158})$$

983 Therefore, we can conclude that $\mathbb{P}(E_{t+1}) \geq \mathbb{P}(E_{t,3}) - \frac{\delta}{4T}$. \square

984 **Lemma D.12.** Define $\Lambda_t := \sum_{j=0}^{t-1} a_{t,j} \|\bar{x}_j - \bar{x}_{j+1}\|^2$ where $a_{t,j} := \beta_1^{t-j-1} (t-j + \frac{\beta_1}{1-\beta_1})$. Under
985 the same conditions in Lemma D.11, event $E_{t,3}$ implies

$$\begin{aligned} \Lambda_t &\leq \left(1 - \frac{(1-\beta_1)^2}{2}\right) \Lambda_{t-1} + \frac{64\eta^2}{1-\beta_1} \|\nabla f(\bar{z}_t)\|_{H_r^{-2}}^2 \\ &\quad + \frac{36\eta^2}{\lambda^2} (1-\beta_1) \sum_{j=rK}^{t-1} \beta_1^{t-j-1} \left[\frac{\eta^2 L^2 \sigma^2}{\lambda^2} KA + \mathbb{E}_m \|\widehat{g}_j^m - \nabla f(x_j^m)\|^2 \right]. \end{aligned} \quad (\text{D.159})$$

986 *Proof.* By the update rule, it always holds that

$$\|\bar{z}_t - \bar{x}_t\|^2 = \left(\frac{\beta_1}{1-\beta_1}\right)^2 \|\bar{x}_t - \bar{x}_{t-1}\|^2. \quad (\text{D.160})$$

987 By AM-GM inequality and event $E_{t,1}$,

$$\begin{aligned} \|\bar{x}_t - \bar{x}_{t-1}\|^2 &= \eta^2 \|\mathbb{E}_m(H_{t-1}^m)^{-1} u_{t-1}^m\|^2 \\ &\leq 2\eta^2 \|\mathbb{E}_m(H_{t-1}^m)^{-1} \bar{u}_{t-1}\|^2 + \frac{2\eta^2}{\lambda^2} \mathbb{E}_m \|u_{t-1}^m - \bar{u}_{t-1}\|^2 \\ &\leq 4\eta^2 \|\mathbb{E}_m H_r^{-1} \bar{u}_{t-1}\|^2 + \frac{2\eta^2}{\lambda^2} \mathbb{E}_m \|u_{t-1}^m - \bar{u}_{t-1}\|^2. \end{aligned} \quad (\text{D.161})$$

988 Event $E_{t,1}$ implies $z_j^m, x_j^m \in \mathbf{conv}(\mathbf{B}_{R_0}(\Omega))$ for all $j \leq t$ and thus

$$\begin{aligned}
\mathbb{E}_m \|u_{t-1}^m - \bar{u}_{t-1}\|^2 &\leq (1 - \beta_1) \sum_{j=rK}^{t-1} \beta_1^{t-j-1} \mathbb{E}_m [\|\widehat{g}_j^m - \bar{g}_j\|^2] \\
&\leq 2(1 - \beta_1) \sum_{j=rK}^{t-1} \beta_1^{t-j-1} \mathbb{E}_m \left[\|\widehat{g}_j^m - \nabla f(x_j^m)\|^2 + \|\nabla f(x_j^m) - \mathbb{E}_m \nabla f(x_j^m)\|^2 \right] \\
&\leq 2(1 - \beta_1) \sum_{j=rK}^{t-1} \beta_1^{t-j-1} \mathbb{E}_m \left[L^2 \|x_j^m - \bar{x}_j\|^2 + \|\widehat{g}_j^m - \nabla f(x_j^m)\|^2 \right] \\
&\leq 2(1 - \beta_1) \sum_{j=rK}^{t-1} \beta_1^{t-j-1} \left[\frac{\eta^2 L^2 \sigma^2}{\lambda^2} K A + \mathbb{E}_m \|\widehat{g}_j^m - \nabla f(x_j^m)\|^2 \right].
\end{aligned} \tag{D.162}$$

989

$$\begin{aligned}
\frac{1}{4} \|\bar{u}_{t-1}\|_{H_r^{-2}}^2 &\leq \left\| (1 - \beta_1) \sum_{j=0}^{t-1} \beta_1^{t-j-1} \nabla f(\bar{x}_t) \right\|_{H_r^{-2}}^2 + \left\| (1 - \beta_1) \sum_{j=0}^{t-1} \beta_1^{t-j-1} [\nabla f(\bar{x}_j) - \nabla f(\bar{x}_t)] \right\|_{H_r^{-2}}^2 \\
&\quad + \left\| (1 - \beta_1) \sum_{j=0}^{t-1} \beta_1^{t-j-1} \mathbb{E}_m [\nabla f(x_j^m) - \nabla f(\bar{x}_j)] \right\|_{H_r^{-2}}^2 + \left\| (1 - \beta_1) \sum_{j=0}^{t-1} \beta_1^{t-j-1} \mathbb{E}_m [\widehat{g}_j^m - \nabla f(x_j^m)] \right\|_{H_r^{-2}}^2 \\
&\leq \|\nabla f(\bar{x}_t)\|_{H_r^{-2}}^2 + \frac{(1 - \beta_1)}{\lambda^2} \sum_{j=0}^{t-1} \beta_1^{t-j-1} L^2 \|\bar{x}_j - \bar{x}_t\|^2 \\
&\quad + \frac{(1 - \beta_1)}{\lambda^2} \sum_{j=0}^{t-1} \beta_1^{t-j-1} \left[\|\mathbb{E}_m [\widehat{g}_j^m - \nabla f(x_j^m)]\|^2 + \|\mathbb{E}_m [\nabla f(x_j^m) - \nabla f(\bar{x}_j)]\|^2 \right] \\
&\leq 2 \|\nabla f(\bar{z}_t)\|_{H_r^{-2}}^2 + \frac{2L^2}{\lambda^2} \|\bar{z}_t - \bar{x}_t\|^2 + \frac{(1 - \beta_1)}{\lambda^2} \sum_{j=0}^{t-1} \beta_1^{t-j-1} L^2 (t - j) \sum_{i=j}^{t-1} \|\bar{x}_i - \bar{x}_{i+1}\|^2 \\
&\quad + \frac{(1 - \beta_1)}{\lambda^2} \sum_{j=0}^{t-1} \beta_1^{t-j-1} \left[\|\mathbb{E}_m [\widehat{g}_j^m - \nabla f(x_j^m)]\|^2 + \|\mathbb{E}_m [\nabla f(x_j^m) - \nabla f(\bar{x}_j)]\|^2 \right] \\
&\leq 2 \|\nabla f(\bar{z}_t)\|_{H_r^{-2}}^2 + \frac{2L^2}{\lambda^2} \|\bar{z}_t - \bar{x}_t\|^2 + \frac{L^2}{\lambda^2} \sum_{j=0}^{t-1} a_{t,j} \|\bar{x}_j - \bar{x}_{j+1}\|^2 \\
&\quad + \frac{(1 - \beta_1)}{\lambda^2} \sum_{j=0}^{t-1} \beta_1^{t-j-1} \left[\|\mathbb{E}_m [\widehat{g}_j^m - \nabla f(x_j^m)]\|^2 + \|\mathbb{E}_m [\nabla f(x_j^m) - \nabla f(\bar{x}_j)]\|^2 \right].
\end{aligned} \tag{D.163}$$

990 Here $a_{t,j} := \beta_1^{t-j-1}(t-j + \frac{\beta_1}{1-\beta_1})$. For $j \leq t-2$, we have $a_{t,j} \leq \beta_1(2-\beta_1)a_{t-1,j}$. Since

991 $\Lambda_t = \sum_{j=0}^{t-1} a_{t,j} \|\bar{x}_j - \bar{x}_{j+1}\|^2$, we conclude that

$$\begin{aligned} \|\bar{x}_t - \bar{x}_{t-1}\|^2 &\leq 64\eta^2 \left[\|\nabla f(\bar{z}_t)\|_{H_r^{-2}}^2 + \frac{L^2}{\lambda^2} \Lambda_{t-1} \right] + \frac{4\eta^2}{\lambda^2} (1-\beta_1) \sum_{j=rK}^{t-1} \beta_1^{t-j-1} \left[\frac{\eta^2 L^2 \sigma^2}{\lambda^2} KA + \mathbb{E}_m \|\widehat{g}_j^m - \nabla f(x_j^m)\|^2 \right] \\ &\quad + \frac{32\eta^2(1-\beta_1)}{\lambda^2} \sum_{j=0}^{t-1} \beta_1^{t-j-1} \left[\|\mathbb{E}_m[\widehat{g}_j^m - \nabla f(x_j^m)]\|^2 + \|\mathbb{E}_m[\nabla f(x_j^m) - \nabla f(\bar{x}_j)]\|^2 \right] \\ &\leq 64\eta^2 \left[\|\nabla f(\bar{z}_t)\|_{H_r^{-2}}^2 + \frac{L^2}{\lambda^2} \Lambda_{t-1} \right] \\ &\quad + \frac{36\eta^2}{\lambda^2} (1-\beta_1) \sum_{j=rK}^{t-1} \beta_1^{t-j-1} \left[\frac{\eta^2 L^2 \sigma^2}{\lambda^2} KA + \mathbb{E}_m \|\widehat{g}_j^m - \nabla f(x_j^m)\|^2 \right], \end{aligned} \quad (\text{D.164})$$

992 and

$$\Lambda_t \leq \beta_1(2-\beta_1)\Lambda_{t-1} + \frac{1}{1-\beta_1} \|\bar{x}_t - \bar{x}_{t-1}\|^2. \quad (\text{D.165})$$

993 This completes the proof. \square

994 D.5 Further Discussion

995 **Compared to other results under centralized weakly convex setting.** Theorem D.2 can reduce
996 to Minibatch Adam (by substituting M, K with 1 and σ with $\frac{\sigma}{\sqrt{MK}}$ in (D.27) [Petrov, 1992]), and
997 the convergence guarantee is

$$\frac{\lambda}{R} \sum_{r=0}^{R-1} \|\nabla f_{\gamma}^{H_r}(\bar{z}_r)\|_{H_r^{-1}}^2 = \tilde{\mathcal{O}} \left(\frac{L\Delta}{R} + \sqrt{\frac{\lambda\Delta\sigma^2}{\gamma MKR}} + \left(\frac{L\Delta\sigma^{\frac{\alpha}{\alpha-1}}}{(MK)^{\frac{2(\alpha-1)}{2(\alpha-1)}R}} \right)^{\frac{2(\alpha-1)}{3\alpha-2}} \right). \quad (\text{D.166})$$

998 Therefore, in centralized setting with iteration number R and batch size 1, our guarantee for squared
999 norm of gradient of Moreau envelope is

$$\tilde{\mathcal{O}} \left(\frac{L\Delta}{R} + \sqrt{\frac{\lambda\Delta\sigma^2}{\gamma R}} + \left(\frac{L\Delta\sigma^{\frac{\alpha}{\alpha-1}}}{R} \right)^{\frac{2(\alpha-1)}{3\alpha-2}} \right). \quad (\text{D.167})$$

1000 The last term is induced by the bias of clipped gradient. For simplicity, let $R \gtrsim \frac{L\Delta}{\sigma^2}$ so that the last
1001 term can be dominated by the first term. Then we obtain

$$\tilde{\mathcal{O}} \left(\frac{L\Delta}{R} + \sqrt{\frac{\lambda\Delta\sigma^2}{\gamma MKR}} \right). \quad (\text{D.168})$$

1002 In the previous literature of weakly convex function [Davis and Drusvyatskiy, 2019, Alacaoglu et al.,
1003 2020, Mai and Johansson, 2021], f is typically non-smooth and stochastic gradient is assumed to
1004 have bounded second order moment. This is weaker than the smoothness assumption but stronger
1005 than that of noise with bounded moment. There are a few existing results for smooth objective [Davis
1006 and Drusvyatskiy, 2019, Mai and Johansson, 2020, Deng and Gao, 2021], but they set $\tau = L$. Overall,
1007 our result is the first convergence guarantee for smooth weakly convex function with $\tau \ll L$ and
1008 bounded-moment noise.

1009 **Dependence on β_2 .** The default setting of β_2 in the Adam optimizer of PyTorch is 0.999, which is
1010 a constant close to 1. Adam with small β_2 has been shown to diverge in some examples [Reddi et al.,
1011 2019]. However, if it is too close to 1, e.g., $\beta_2 \geq 1 - \mathcal{O}(T^{-1})$, then the denominator would be too

1012 stagnant to provide adaptivity. Therefore, to derive a proper range for β_2 is crucial in the theoretical
 1013 analysis of Adam.

1014 On the other hand, β_2 is notoriously difficult to handle even under centralized setting. In finite sum
 1015 case, [Zou et al. \[2019\]](#) assumes $\beta_2 \geq 1 - \mathcal{O}(T^{-1})$. [Shi et al. \[2020\]](#) suggests that $\beta_2 \geq 1 - \mathcal{O}(n^{-3.5})$
 1016 suffices, where n is sample size. [Zhang et al. \[2022b\]](#) claims Adam can converge to the neighborhood
 1017 of stationary points with constant radius if $\beta_2 \geq 1 - \mathcal{O}(n^{-3})$. Further, [Wang et al. \[2022\]](#) shows
 1018 Adam can converge to stationary points if β_2 is sufficiently close to 1, but the explicit bound is
 1019 missing. In streaming data case, [Défossez et al. \[2020\]](#) shows β_2 can be a constant but relies on the
 1020 bounded gradient assumption. [\[Li et al., 2024b\]](#) suggests $\beta_2 \geq 1 - \tilde{\mathcal{O}}(T^{-\frac{1}{2}})$.

1021 As for distributed setting, works discussing the range of β_2 are much fewer. Our theory requires
 1022 $\beta_2 \geq 1 - \tilde{\mathcal{O}}(K^{-\frac{3}{2}}R^{-\frac{1}{2}})$. For distributed Adam, [Karimireddy et al. \[2020a\]](#), [Zhao et al. \[2022\]](#)
 1023 fixed the denominator during local iterations and thus did not discuss the range of β_2 . To the best
 1024 of our knowledge, our result is the first one to show the $\tilde{\mathcal{O}}(R^{-\frac{1}{2}})$ dependence with respect to R .
 1025 Nevertheless, it is an interesting question to improve the dependence on K . Since K is usually a
 1026 constant in practice, our results suggest $\beta_2 \geq 1 - \tilde{\mathcal{O}}(\mathcal{R}^{-\frac{1}{2}})$ in essence. Still, we believe that the
 1027 dependence on K has room for improvement. We leave this for future work.

1028 **Dependence on λ .** λ in the denominator of Adam is aimed to avoid numerical instability, and
 1029 usually a small constant in practice. Note $H_r = \mathbf{diag}(\sqrt{V_r + \lambda^2})$ and v_r is the EMA of squared
 1030 gradients. Informally, v_r vanishes as r grows and thus H_r would gradually reduce to λI_d . In the
 1031 worst case, H_r can be bounded by a constant. In conclusion, the LHS in (3.9) is roughly the averaged
 1032 squared gradient norm if λ is not too small. It is worth noting that λ can be arbitrarily small or even 0
 1033 in [\[Défossez et al., 2020, Wang et al., 2022, 2024\]](#). However, their results all depend on $\mathbf{poly}(d)$. It is
 1034 still an interesting question to get dimension-free result with small λ .

1035 **Dependence on β_1 .** The default setting of β_1 in PyTorch is 0.9, a constant away from 0 and 1. In
 1036 the centralized setting, [Li et al. \[2024b\]](#) requires $\beta_1 = 1 - \mathcal{O}(T^{-\frac{1}{2}})$ to converge, which is too large.
 1037 [Défossez et al. \[2020\]](#) shows $\mathcal{O}((1 - \beta_1)^{-1})$, which is the state of the art result to the best of our
 1038 knowledge. However, it relies on the bounded gradient assumption. Regarding the dependence on
 1039 β_1 , our convergence rate in [Theorem D.1](#) suggests $\mathcal{O}((1 - \beta_1)^{-2})$. Although it also supports any
 1040 constant choice of β_1 , we leave the exploration of better dependence for future work.

1041 E Failure of Standard SGD with Heavy-Tailed Noise

1042 The convergence of standard SGD in high probability is widely studied. If we assume the noises are
 1043 light-tailed, *e.g.*, sub-exponential, sub-gaussian, then SGD can get high probability bound depending
 1044 on $\log \frac{1}{\delta}$. However, if only finite variance is assumed, [Sadiev et al. \[2023\]](#) has shown that standard
 1045 SGD fails to get a high probability bound having logarithmic dependence on $\frac{1}{\delta}$. In fact, this claim is
 1046 still valid when the stochastic noises only have finite α th-moment, as shown in [Theorem E.1](#) below.
 1047 Therefore, gradient clipping is necessary to get the $\log \frac{1}{\delta}$ bound.

1048 **Theorem E.1.** *For any $\varepsilon > 0$, $\delta \in (0, 1)$, and SGD with the iteration number T and learning rate η ,*
 1049 *there exists an 1D-problem satisfying Assumption 1, 2, 3, 4, with $\Omega = \mathbb{R}$ and $L = \mu$, such that, if*
 1050 *$0 < \eta \leq 1/L$, then*

$$\mathbb{P}\{f(x_T) - f_* \geq \varepsilon\} \leq \delta \implies T = \tilde{\Omega}\left(\frac{\sigma}{\delta^{1/\alpha}} \sqrt{\frac{L}{\varepsilon}}\right). \quad (\text{E.1})$$

1051 *Proof.* We follow the construction of the counter example in [Sadiev et al. \[2023\]](#). To prove the above
 1052 theorem, we consider a simple 1D-problem $f(x) = Lx^2/2$. It is easy to see that the considered
 1053 problem is L -strongly convex, L -smooth, and has optimum at $x_* = 0$. We construct the noise in an
 1054 adversarial way with respect to the parameters of the SGD. Concretely, the noise depends on the
 1055 number of iterates t , learning rate η , target precision ε , the starting point x_0 , and the moment bound

1056 σ such that

$$\nabla F(x_t; \xi_t) = Lx_t - \sigma\xi_t, \quad (\text{E.2})$$

1057 where

$$\xi_t = \begin{cases} 0, & \text{if } t < T - 1 \text{ or } (1 - \eta L)^T |x_0| > \sqrt{\frac{2\varepsilon}{L}}, \\ \begin{cases} -A, & \text{with probability } \frac{1}{2A^\alpha}, \\ 0, & \text{with probability } 1 - \frac{1}{A^\alpha}, \\ A, & \text{with probability } \frac{1}{2A^\alpha}, \end{cases} & \text{otherwise} \end{cases} \quad (\text{E.3})$$

1058 where $A = \max \left\{ \frac{2\sqrt{\frac{2\varepsilon}{L}}}{\eta\sigma}, 1 \right\}$. We note that $\mathbb{E}[\xi_t] = 0$ and $\mathbb{E}[\nabla F(x_t; \xi_t)] = \nabla f(x_t)$. Furthermore,

$$\mathbb{E}[|\xi_t|^\alpha] \leq \frac{1}{2A^\alpha} A^\alpha + \frac{1}{2A^\alpha} A^\alpha = 1, \quad (\text{E.4})$$

1059 which implies that Assumption 3 holds.

1060 We are interested in the situation when

$$\mathbb{P}\{f(x_T) - f_* \geq \varepsilon\} \leq \delta, \quad (\text{E.5})$$

1061 for $\delta \in (0, 1)$. We first prove that this implies $(1 - \eta L)^T |x_0| \leq \sqrt{\frac{2\varepsilon}{L}}$. To do that we proceed by
1062 contradiction and assume that

$$(1 - \eta L)^T |x_0| > \sqrt{\frac{2\varepsilon}{L}}. \quad (\text{E.6})$$

1063 By construction, this implies that $\xi_t = 0, \forall t \in \{0, \dots, T - 1\}$. This, in turn, implies that $x_T =$
1064 $(1 - \eta L)^T x_0$, and further, by (E.6) that

$$\mathbb{P}\{f(x_T) - f_* \geq \varepsilon\} = \mathbb{P}\left\{|x_T| \geq \sqrt{\frac{2\varepsilon}{L}}\right\} = 1.$$

1065 Thus, the contradiction shows that $(1 - \eta L)^T |x_0| \leq \sqrt{\frac{2\varepsilon}{L}}$. Using (E.3), we obtain

$$f(x_T) - f_* = \frac{L}{2} [(1 - \eta L)^T x_0 + \eta\sigma\xi_{T-1}]^2. \quad (\text{E.7})$$

1066 Furthermore,

$$\begin{aligned} \mathbb{P}\{f(x_T) - f_* \geq \varepsilon\} &= \mathbb{P}\left\{|(1 - \eta L)^T x_0 + \eta\sigma\xi_{T-1}| \geq \sqrt{\frac{2\varepsilon}{L}}\right\} \\ &= \mathbb{P}\left\{|\eta\sigma\xi_{T-1}| \geq \sqrt{\frac{2\varepsilon}{L}} + (1 - \eta L)^T |x_0|\right\} \\ &\geq \mathbb{P}\left\{|\eta\sigma\xi_{T-1}| \geq 2\sqrt{\frac{2\varepsilon}{L}}\right\} \\ &= \mathbb{P}\left\{|\xi_{T-1}| \geq \frac{2\sqrt{\frac{2\varepsilon}{L}}}{\eta\sigma}\right\}. \end{aligned} \quad (\text{E.8})$$

1067 Now if $\frac{2\sqrt{\frac{2\varepsilon}{L}}}{\eta\sigma} < 1$ then $A = 1$. Therefore,

$$1 = \mathbb{P}\left\{|\xi_{T-1}| \geq \frac{2\sqrt{\frac{2\varepsilon}{L}}}{\eta\sigma}\right\} \leq \mathbb{P}\{f(x_T) - f_* > \varepsilon\} \leq \delta, \quad (\text{E.9})$$

1068 yielding contradiction, which implies that $\frac{2\sqrt{\frac{2\varepsilon}{L}}}{\eta\sigma} \geq 1$, *i.e.*, $\eta \leq 2\sqrt{\frac{2\varepsilon}{L\sigma^2}}$. In this case, $A = \frac{2\sqrt{\frac{2\varepsilon}{L}}}{\eta\sigma}$
 1069 and we have

$$\delta \geq \mathbb{P}\{f(x_T) - f_* \geq \varepsilon\} \geq \mathbb{P}\left\{|\xi_{T-1}| \geq \frac{2\sqrt{\frac{2\varepsilon}{L}}}{\eta\sigma}\right\} = \frac{1}{A^\alpha}. \quad (\text{E.10})$$

1070 This implies that $\eta \leq \frac{2\delta^{1/\alpha}}{\sigma} \sqrt{\frac{2\varepsilon}{L}}$. Combining this inequality with $T \geq \frac{1}{2\eta L} \log \frac{Lx_0^2}{2\varepsilon}$ yields

$$T = \Omega\left(\frac{\sigma}{\delta^{1/\alpha}} \sqrt{\frac{L}{\varepsilon}} \log \frac{Lx_0^2}{2\varepsilon}\right). \quad (\text{E.11})$$

1071 This concludes the proof. □