

CONVERGENCE OF DISTRIBUTED ADAPTIVE OPTIMIZATION WITH LOCAL UPDATES

Anonymous authors

Paper under double-blind review

ABSTRACT

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

1 INTRODUCTION

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

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

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

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

Despite the success of *Local SGD* in federated learning (McMahan et al., 2017), it may not exhibit good performance when training Transformer-based large language models (LLMs). Many empirical studies suggest that adaptive methods (e.g., Adam (Kingma & Ba, 2014)) are much better suited for natural language processing than vanilla SGD (Goodfellow et al., 2016; Zhang et al., 2020; Kunstner et al., 2023; Pan & Li, 2023). Furthermore, as shown in Zhang et al. (2019; 2020), language models tend to have unbounded global smoothness and heavy-tailed noise, which may also contribute to the worse performance of SGD. Parallelizing adaptive methods requires an even more expensive communication cost since additional terms, such as the momentum or the Adam denominator, need to be synchronized. Previous works on distributed adaptive optimization have utilized compression and quantization techniques to address this issue (Bernstein et al., 2018; Wangni et al., 2018; Wang et al., 2023). While Douillard et al. (2023) has shown the great empirical success of *Local Adam*, to the best of our knowledge, there are no theoretical results trying to improve training efficiency or adaptive methods from the perspective of intermittent communication.

In this paper, we investigate **distributed adaptive optimization algorithms in the homogeneous regime**, in order to establish theoretical guarantees for the benefits of local iterations in reduc-

054
 055
 056
 057
 058
 059
 060
 ing communication complexity. We focus on the convex or weakly convex setting, because in the
 non-convex setting, without non-standard strong smoothness assumptions, we are not aware of any
 theoretical-proven advantages of local iterations, even for non-adaptive methods¹. Further, in the
 case of Adam, we consider the weakly convex setting (as opposed to the standard convex setting),
 since we are not aware of any results on the convergence rate of Adam which take advantage of con-
 vexity. To handle unbounded global smoothness and heavy-tailed noise, we use the coordinate-wise
 gradient clipping mechanism.

061
 062
 063
 We propose a distributed version of Adam, namely, *Local* Adam, with gradient clipping. Our al-
 gorithm also reduces to *Local* SGD with momentum (*Local* SGDM), with some specific hyper-
 parameter choices.

- 064
- In Theorem 2, we establish the first convergence guarantee for *Local* SGDM in the convex
 065 setting, which outperforms the convergence rate of *Minibatch* SGD. The rate we obtain
 066 is in line with the rate of *Local* SGD (Woodworth et al., 2020a) .
 - In Theorem 3, we establish a convergence rate for *Local* Adam in the weakly convex set-
 067 ting. We show that *Local* Adam can provably improve communication efficiency compared
 068 to its minibatch baseline.

069
 070
 071 For the first time, we are able to show the benefits of local iterations for the two commonly used
 072 algorithms, SGDM and Adam. This suggests that we may be able to improve the training efficiency
 073 of LLMs by using intermittent communication.

074
 075 Additionally, our results hold under generalized smoothness and heavy-tailed noise. Our result is
 076 the first high probability bound for distributed optimization algorithms with local updates, to the
 077 best of our knowledge. The conventional in-expectation rate seems fail to capture some important
 078 properties like heavy/light tailed noise distribution. The high probability convergence guarantee can
 079 sometimes be more informative and useful in practice (Gorbunov et al., 2020).

080 As for technical contribution, we use a **novel technique to prove contraction for adaptive meth-**
 081 **ods**, which bounds the consensus error between the iterates at different workers. This is a key step
 082 in proving benefits of local updates. Different from *Local* SGD, our update direction involves mo-
 083 mentum or even distorted momentum due to the denominator in *Local* Adam, making it challenging
 084 to disentangle these accumulated stochastic gradients. To address this issue, we define and analyze
 085 an auxiliary sequence which is conditionally independent of the latest stochastic gradient and thus
 086 can construct a martingale. We will introduce the technique in more details in Section 5.

087 1.1 ORGANIZATION

088
 089 Section 2 provides the most related work to ours. Section 3 provides the problem setup, assumptions
 090 and the *Local* Adam algorithm. We then show our main results for *Local* SGDM in Section 4.1 and
 091 *Local* Adam in Section 4.2. Finally, in Section 5, we present the proof sketch of *Local* Adam,
 092 highlighting the technical challenges and our solution.

093 1.2 NOTATION

094 Let $\|\cdot\|$ be the standard Euclidean norm of a vector or the spectral norm of a matrix. For any
 095 $x, y \in \mathbb{R}^d$, the expressions $x + y$, $x \odot y$, $\frac{x}{y}$ stand for coordinate-wise sum, product and division,
 096 respectively. And $x \preceq y$ means each coordinate of $x - y$ is no greater than 0. Furthermore, we use
 097 $x^2, \sqrt{x}, |x|$ to denote the coordinate-wise square, square root and absolute value. We use $\mathbb{E}_m[X_m]$
 098 to denote the average $\frac{1}{M} \sum_{m=1}^M X_m$. The coordinate-wise clipping operator $\text{clip}(\cdot, \rho) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is
 099 defined as $[\text{clip}(X, \rho)]_i = \text{sgn}([X]_i) \cdot \min\{|X_i|, \rho\}$. We use $[N]$ to denote the set $\{1, 2, \dots, N\}$. For

100
 101
 102
 103
 104
 105
 106
 107
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).

108 a subset $\Omega_0 \subset \mathbb{R}^d$, let $\text{conv}(\cdot)$ denote the convex hull of Ω_0 and $\mathbf{B}_{R_0}(\Omega_0)$ denote the neighborhood
 109 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
 110 omit logarithmic factors.
 111

112 2 RELATED WORK

113 **Theoretical benefits of local updates in distributed optimization.** Algorithms with local up-
 114 dates have been used among practitioners for a long time to reduce communication complexity
 115 ([McMahan et al., 2017](#)). In the homogeneous and convex setting, *Local SGD* and its variants have
 116 been shown to outperform the minibatch baseline, for a fixed amount of gradient computations and
 117 communication rounds. [Woodworth et al. \(2020a\)](#) is the first to show that *Local SGD* can prov-
 118 ably outperform *Minibatch SGD*. [Yuan & Ma \(2020\)](#) develops FedAC to further accelerate *Local*
 119 *SGD*. In the heterogeneous case, [Woodworth et al. \(2020b\)](#) demonstrates the advantages of *Local*
 120 *SGD* when heterogeneity is very low. Algorithms with local updates have also been studied in the
 121 non-convex setting ([Karimireddy et al., 2020b](#); [Yang et al., 2021](#); [Glasgow et al., 2022](#)), includ-
 122 ing momentum-based and adaptive methods ([Reddi et al., 2020](#); [Karimireddy et al., 2020a](#)), though
 123 no advantage of local iterations over minibatch has been shown, without non-standard assumptions
 124 such as 3rd-order smoothness. Notably, [Liu et al. \(2022\)](#) is one closely related work to ours, which
 125 considers *Local SGD* with gradient clipping in homogeneous and non-convex setting and claims
 126 that the convergence guarantee is better than naive parallel of centralized clipped-SGD. However, it
 127 still cannot outperform minibatch baseline (with batch size K for each worker in each round) and
 128 thus fails to demonstrate the benefits of local iterations.
 129

130 **Convergence of centralized Adam.** Adam was first proposed by [Kingma & Ba \(2014\)](#) with con-
 131 vergence guarantee in online convex optimization. However, [Reddi et al. \(2019\)](#) found a gap in the
 132 original analysis of Adam and constructed a counter example to show its divergence. Since then,
 133 many works have developed convergence analyses of Adam with various assumptions and hyper-
 134 parameter settings. [Guo et al. \(2021\)](#) assumed the denominator is bounded from below and above
 135 by two constants, which typically requires a bounded gradient assumption or the AdaBound variant
 136 ([Luo et al., 2019](#)). [Défossez et al. \(2020\)](#) assumed a bounded gradient and their convergence guar-
 137 aintee depends on $\text{poly}(d)$. [Zhang et al. \(2022b\)](#); [Wang et al. \(2022\)](#) considered a finite sum setting and
 138 showed that Adam converges to the neighborhood of stationary points. One closely related work to
 139 ours is [Li et al. \(2024c\)](#), which established a high probability bound without a bounded gradient as-
 140 sumption. However they assumed that noise is bounded almost surely. Another recent work ([Wang](#)
 141 [et al., 2024](#)) provided a guarantee of $\mathcal{O}(1/\varepsilon^4)$ with dependence on $\text{poly}(d)$. Beyond the guarantees
 142 on gradient norm given by non-convex analyses, no stronger bounds (e.g., on function error) are
 143 known for Adam in the convex case.

144 **Convergence of distributed adaptive algorithms.** In the federated learning literature, [Reddi et al.](#)
 145 [\(2020\)](#) introduced a framework, FedOPT, to leverage both worker optimizer and server optimizer.
 146 Many works explored adaptive server optimizer while fixing worker side as vanilla SGD. The theo-
 147 retical results of local adaptive algorithms are much fewer. Some works have studied *Local Adam*
 148 and *Local AMSGrad* with fixed momentum state during local iterations ([Karimireddy et al., 2020a](#);
 149 [Chen et al., 2020](#); [Zhao et al., 2022](#)). They also needed stringent assumptions such as a huge batch
 150 size depending on the inverse of target error, bounded stochastic gradients, vanishing difference
 151 between denominator, etc., which are not standard. [Wang et al. \(2021\)](#) explored adaptive worker
 152 optimizer based on centralized algorithm, where the state of worker optimizer changes in local up-
 153 dates. However, their analysis relied on an explicit assumptions ([Wang et al., 2021](#), Assumption
 154 1) on the contraction property of worker optimizer. Some recent works ([Li et al., 2024a](#); [Anyszka](#)
 155 [et al., 2024](#)) discussed Polyak stepsizes with an exact local proximal operator, which is inaccessible
 156 in most cases by gradient-based optimizers. To the best of our knowledge, there is no end-to-end
 157 convergence guarantee for distributed adaptive algorithms with local iterations.

158 3 PROBLEM SETUP

159 Consider the distributed optimization problem
 160

$$\min_{x \in \mathbb{R}^d} f(x) := \mathbb{E}_{\xi \sim \mathcal{D}}[F(x; \xi)]. \quad (3.1)$$

162 Here \mathcal{D} is the data distribution and f is the population loss function. We consider a setting with
 163 M parallel workers, and a budget of R total communication rounds, and T total gradient computa-
 164 tions at each worker. We will describe the implementation of the *local* and *minibatch* versions of
 165 a centralized algorithm \mathcal{A} , which uses a single stochastic gradient in each iteration. And these are
 166 illustrated in Figure 1.

167

168

169

170

171

172

173

174

175

176

177

178

179

180

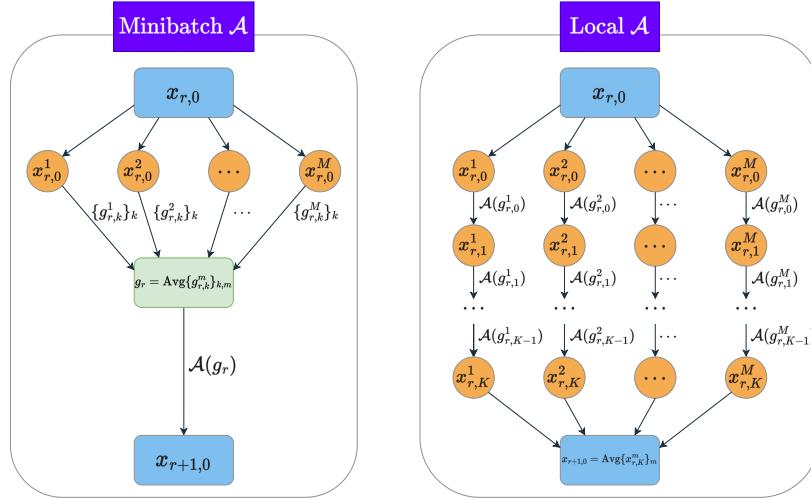
181

182

183

184

185



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

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

196 *Local* Adam is shown in Algorithm 1, which is a natural extension of centralized Adam (Kingma &
 197 Ba, 2014). The stochastic gradient is clipped by an coordinate-wise clipping operator with threshold
 198 ρ . After K steps of local updates, all the workers average their current iterates x_t^m , their first order
 199 momentum u_t^m , and their second order momentum v_t^m . These averaged quantities become the values
 200 used at the beginning of the next local round. Note that there are two slight differences from original
 201 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
 202 $1 - \beta_2^t$, respectively. Second, λ in the denominator is in the square root, while it is outside of the
 203 denominator in original Adam. These modifications do not harm the spirit of Adam and are made
 204 for the convenience of analysis.

205 3.1 ASSUMPTIONS

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

207 **Assumption 1** (Lower-boundedness). f is closed, twice continuously differentiable and
 208 $\inf_{x \in \mathbb{R}^d} f(x) = f(x_*) = f_* > -\infty$.

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

$$210 \| \nabla f(x) - \nabla f(y) \| \leq L \| x - y \|, \quad (3.2)$$

$$211 \| \nabla f(x) \|^2 \leq 2L(f(x) - f_*). \quad (3.3)$$

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

Similar to [Sadiev et al. \(2023\)](#), we only require some properties of f on a subset $\Omega \subset \mathbb{R}^d$, since we can prove that all the iterates will not leave this subset with high probability. In contrast, the typical smoothness assumption requires (3.2) on the entire domain.

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

$$\|\nabla^2 f(x)\| \leq L_0 + L_1 \|\nabla f(x)\|, \quad \forall x \in \mathbb{R}^d. \quad (3.4)$$

[Li et al. \(2024b\)](#) considers an extension called ℓ -smoothness, which replaces the linear function of $\|\nabla f\|$ in the right hand side of (3.4) with a sub-quadratic function $\ell(\cdot)$. As pointed out in [Li et al. \(2024b, Corollary 3.6\)](#), all of these will induce Assumption 2 if Ω is some level-set of the objective function². Therefore, we directly use this more general assumption to get cleaner results.

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

$$\mathbb{E}_{\xi \sim \mathcal{D}} |\nabla F(x; \xi) - \nabla f(x)|^\alpha \preceq \sigma^\alpha. \quad (3.5)$$

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

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

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

$$\mathbb{E}_{\{\xi_i\}^{i.i.d.}_{\mathcal{D}}} \left| \frac{1}{N} \sum_{i=1}^N \nabla F(x; \xi_i) - \nabla f(x) \right|^{\alpha} \preceq c(\alpha) (\sigma / \sqrt{N})^{\alpha}, \quad (3.6)$$

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

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

270 **4 MAIN RESULTS**
 271

272 In this section, we provide our main results for *Local* Adam and its simplified version: *Local* SGDM.
 273 For the first time, we will be able to show the benefits of local iterations for the two algorithms,
 274 compared with their minibatch baselines in certain regime of M, K, R .
 275

276 **4.1 LOCAL SGDM**
 277

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

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

$$\langle \nabla f(x) - \nabla f(y), x - y \rangle \geq \mu \|x - y\|^2, \quad (4.1)$$

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\mu}{2} \|x - y\|^2. \quad (4.2)$$

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

288 **Theorem 1** (Strongly convex, full version see Theorem C.4). *Let Assumption 1, 2, 3, 4 hold for*
 289 $\Omega := \{\|x - x_*\| \leq \sqrt{3}D_0\}$ *and $\mu > 0$. Further assume that $K \gtrsim \log \frac{MKR}{\delta}$, $1 - \beta_1 = \Omega(1)$ and*
 290 $\|\sigma\|_{2\alpha} d^{\frac{1}{2} - \frac{1}{2\alpha}} = \mathcal{O}(\sigma)$. *Then with probability no less than $1 - \delta$, Local SGDM yields*

$$293 f(\hat{x}) - f_* \leq \exp\left(-\Theta\left(\frac{\mu KR}{L}\right)\right) + \tilde{\mathcal{O}}\left(\frac{\sigma^2}{\mu MKR} + \frac{L\sigma^2}{\mu^2 KR^2} + \frac{\sigma^2}{\mu} \left(\frac{L^{\frac{1}{2}}}{\mu^{\frac{1}{2}} KR}\right)^{\frac{2(\alpha-1)}{\alpha}}\right). \quad (4.3)$$

296 **Theorem 2** (Convex, full version see Theorem C.5). *Let Assumption 1, 2, 3, 4 hold for $\Omega :=$*
 297 $\{\|x - x_*\| \leq \sqrt{3}D_0\}$ *and $\mu = 0$. Further assume that $K \gtrsim \log \frac{MKR}{\delta}$, $1 - \beta_1 = \Omega(1)$ and*
 298 $\|\sigma\|_{2\alpha} d^{\frac{1}{2} - \frac{1}{2\alpha}} = \mathcal{O}(\sigma)$. *Then with probability no less than $1 - \delta$, Local SGDM yields*

$$300 f(\hat{x}) - f_* \leq \tilde{\mathcal{O}}\left(\frac{LD_0^2}{KR} + \frac{\sigma D_0}{\sqrt{MKR}} + \frac{L^{\frac{1}{2}} \sigma^{\frac{2}{3}} D_0^{\frac{4}{3}}}{K^{\frac{1}{3}} R^{\frac{2}{3}}} + D_0 \left(\frac{(LD_0)^{\frac{1}{2}} \sigma^{\frac{\alpha}{\alpha-1}}}{KR}\right)^{\frac{2(\alpha-1)}{3\alpha-1}}\right). \quad (4.4)$$

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

306 Our method can also be applied to *Minibath* SGDM (by substituting M, K with 1 and σ with
 307 $\frac{\sigma}{\sqrt{MK}}$; see Remark 2), whose convergence guarantee is

$$309 f(\hat{x}) - f_* \lesssim \begin{cases} \exp\left(-\Theta\left(\frac{\mu R}{L}\right)\right) + \tilde{\mathcal{O}}\left(\frac{\sigma^2}{\mu MKR}\right), & \text{if } \mu > 0, \\ \tilde{\mathcal{O}}\left(\frac{LD_0^2}{R} + \frac{\sigma D_0}{\sqrt{MKR}}\right), & \text{otherwise.} \end{cases} \quad (4.5)$$

313 This rate matches the well-known in-expectation lower bound on the convergence rate of *Minibatch*
 314 SGD (up to logarithmic factors). In fact, our analysis improves the state-of-the-art rate for strongly-
 315 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
 316 convex setting, our rate is consistent with the state-of-the-art centralized in-expectation bound of
 317 SGDM in Sebbouh et al. (2021). Further notice that the last term in both (4.3) and (4.4) is due to the
 318 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
 319 this case, our guarantee for *Local* SGDM is aligned with the rate of *Local* SGD in Woodworth et al.
 320 (2020a); Khaled et al. (2020) up to logarithmic factor. Therefore, we can see the benefits of local
 321 iterations in the large M and large K regime compared to minibatch baseline.

323 We defer the complete version and detailed proof to Appendix C.

324 4.2 LOCAL ADAM
325

326 The convergence of Adam is much more difficult to prove. Reddi et al. (2019) pointed out that
327 the original proof in Kingma & Ba (2014) in centralized convex setting was incorrect. Therefore,
328 the convergence of Adam in for convex function is of independent interest and beyond our scope.
329 Instead, we turn to consider Adam in the weakly convex setting.

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

$$\langle \nabla f(x) - \nabla f(y), x - y \rangle \geq -\tau \|x - y\|^2, \quad (4.6)$$

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle - \frac{\tau}{2} \|x - y\|^2, \quad \nabla^2 f(x) \succeq -\tau I_d. \quad (4.7)$$

335 Note that L -smoothness implies that Assumption 5 always holds with $\tau = L$. Also note that here we
336 assume the weak convexity holds in \mathbb{R}^d for technical simplicity. Let $H_r = \text{diag}(\sqrt{v_r + \lambda^2}) \succeq \lambda I_d$
337 and $\Delta := f(x_0) - f_*$. Furthermore, inspired by Liu et al. (2020b), define an auxiliary sequence
338 $\{z_{r,k}^m\}$ as:

$$z_{r,k+1}^m = \begin{cases} (x_{r,k+1}^m - \beta_1 x_{r,k}^m)/(1 - \beta_1) & \text{if } k \neq K - 1, \\ (x_{r,k+1}^m - \beta_1 \bar{x}_{r,k})/(1 - \beta_1) & \text{otherwise.} \end{cases} \quad (4.8)$$

342 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
343 more general results on Moreau envelope).

344 **Theorem 3** (Full version see Theorem D.3). *Let Assumption 1, 2, 3, 5 hold for $\Omega = \text{conv}(\mathbf{B}_{R_0}(\Omega_0))$,
345 where $\Omega_0 := \{f(x) - f_* \leq 4\Delta\}$ and $R_0 = \sqrt{\Delta/(80L)}$. Further assume $K \gtrsim \log(MKR/\delta)$,
346 $1 - \beta_1 = \Omega(1)$, $\|\sigma\|_{2\alpha} d^{\frac{1}{2} - \frac{1}{2\alpha}} = \mathcal{O}(\sigma)$ and $1 - \beta_2 = \tilde{\mathcal{O}}(K^{-3/2} R^{-1/2})$. Then with probability no
347 less than $1 - \delta$, Local Adam yields*

$$\begin{aligned} & \frac{\lambda}{KR} \sum_{r=0}^{R-1} \sum_{k=0}^{K-1} \|\nabla f(\bar{z}_{r,k})\|_{H_r^{-1}}^2 \\ &= \tilde{\mathcal{O}}\left(\frac{\tau\Delta}{R} + \frac{L\Delta}{KR} + \sqrt{\frac{L\Delta\sigma^2}{MKR}} + \frac{(L\Delta\sigma)^{\frac{2}{3}}}{K^{\frac{1}{3}} R^{\frac{2}{3}}} + \left(\frac{L\Delta\sigma^{\frac{\alpha}{\alpha-1}}}{KR}\right)^{\frac{2(\alpha-1)}{3\alpha-2}}\right). \end{aligned} \quad (4.9)$$

355 The RHS of (4.9) consists of four parts. The first part is $\frac{\tau\Delta}{R} + \frac{L\Delta}{KR}$, which is the optimization term

356 and determined by the upper bound of learning rate η . The second term is $\sqrt{\frac{L\Delta\sigma^2}{MKR}}$, corresponding
357 to the standard statistical lower bound from MKR stochastic gradients (Arjevani et al., 2023).

358 The third 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
359 iterations. And the last one, $\left(\frac{L\Delta\sigma^{\frac{\alpha}{\alpha-1}}}{KR}\right)^{\frac{2(\alpha-1)}{3\alpha-2}}$, is induced by the bias of clipped stochastic gradient
360 and can be dominated when $K^{\frac{3\alpha-4}{2}} \gtrsim \sigma^2 R/(L\Delta)$.

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

$$\tilde{\mathcal{O}}\left(\frac{L\Delta}{R} + \sqrt{\frac{L\Delta\sigma^2}{MKR}}\right), \quad (4.10)$$

363 aligned with (up to logarithmic factor) the state-of-the-art convergence guarantees for smooth
364 weakly convex functions (Davis & Drusvyatskiy, 2019; Deng & Gao, 2021). Suppose $K^{\frac{3\alpha-4}{2}} \gtrsim$
365 $\sigma^2 R/(L\Delta)$ and hence the last term in (4.9) would be dominated and negligible. Now we can ob-
366 serve the benefits of local iterations. Note that both (4.9) and (4.10) have the statistical lower bound
367 $1/\sqrt{MKR}$. Hence when the statistical term dominates, both algorithms have similar worst-case
368 rate. Once we leave the noise-dominated regime, then *Local* Adam converges faster than *Minibatch*
369 Adam whenever $K \gtrsim \sigma^2 R/(L\Delta)$. And the gap will increase as K grows until $K \asymp L/\tau$.

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

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

5 PROOF SKETCH

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

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

$$f_\gamma^H(x) := \min_{y \in \mathbb{R}^d} f(y) + \frac{1}{2\gamma} \|x - y\|_H^2. \quad (5.1)$$

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

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 some $H \succeq \lambda I_d$ and $L/\lambda \geq \gamma^{-1} \geq 2\tau/\lambda$. Then*

$$\nabla f_\gamma^H(z) = \nabla f(y) = H(z - y)/\gamma, \quad \|\nabla f(z)\|_{H^{-1}} \leq 2\gamma L \|\nabla f_\gamma^H(z)\|_{H^{-1}}/\lambda. \quad (5.2)$$

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

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 $k(t) \leq K-1$. We will omit the dependence on t and let $r = r(t), k = k(t)$ if not causing confusion. Further define

$$x_t^m := x_{r,k}^m, g_t^m := g_{r,k}^m, \widehat{g}_t^m := \widehat{g}_{r,k}^m, u_t^m = u_{r,k}^m, v_t^m = v_{r,k}^m, H_t^m := \text{diag}(\sqrt{v_t^m + \lambda^2}) \quad (5.3)$$

Then Algorithm 1 is equivalent to the following update rule:

$$x_{t+1}^m = \begin{cases} x_t^m - \eta(H_t^m)^{-1} u_t^m & \text{if } t \bmod K \not\equiv -1, \\ \bar{x}_t - \eta \mathbb{E}_m[(H_t^m)^{-1} u_t^m] & \text{otherwise.} \end{cases} \quad (5.4)$$

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

$$z_{t+1}^m = \begin{cases} (x_{t+1}^m - \beta_1 x_t^m)/(1 - \beta_1) & \text{if } t \bmod K \not\equiv -1, \\ (x_{t+1}^m - \beta_1 \bar{x}_t)/(1 - \beta_1) & \text{otherwise.} \end{cases} \quad (5.5)$$

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

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

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

$$E_t = \{\mathcal{A}_{j,i} \text{ holds for all } j \leq t-1, i \in \{1, 2, 3, 4\}\}, \quad (5.6)$$

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

432 5.1 BOUNDING THE TRAJECTORY WITH HIGH PROBABILITY
 433

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

$$\begin{aligned}
 f_{\gamma}^{H_r(t+1)}(\bar{z}_{t+1}) &\leq f_{\gamma}^{\lambda}(x_0) - \Omega(\eta) \sum_{j=0}^t \|\nabla f_{\gamma}^{H_r(j)}(\bar{z}_j)\|_{H_{r(j)}^{-1}}^2 + \underbrace{\mathcal{O}(\eta^2) \sum_{j=0}^t \|\mathbb{E}_m[\nabla f(x_j^m) - \widehat{g}_j^m]\|^2}_{\text{stochastic noise}} \\
 &\quad + \underbrace{\mathcal{O}(\eta) \sum_{j=0}^t \|\nabla f(\bar{z}_j) - \mathbb{E}_m[\nabla f(x_j^m)]\|^2}_{\text{discrepancy}} \\
 &\quad + \underbrace{\mathcal{O}(\eta) \sum_{j=0}^t \left\langle \bar{z}_j - \eta H_{r(j)}^{-1} \nabla f(\bar{z}_j) - y_j, \mathbb{E}_m[\mathbb{E}_j[\widehat{g}_j^m] - \widehat{g}_j^m] \right\rangle}_{\text{martingale}} \\
 &\quad + \text{higher order terms.}
 \end{aligned} \tag{5.7}$$

453 To see that the last term is a martingale, note that $H_{r(j)}$ is independent of \widehat{g}_j^m since the stochastic
 454 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
 455 thanks to clipping. Now (5.7) allows us to inductively bound $f_{\gamma}^{H_r(j)}(\bar{z}_j)$ and thus bound $\|\bar{z}_j - \eta H_{r(j)}^{-1} \nabla f(\bar{z}_j) - y_j\|$. After these preliminaries, we are able to apply Bernstein's inequality (Bennett,
 456 1962; Freedman, 1975) to control this martingale. Hence the Moreau envelope at \bar{z}_{t+1} can be
 457 bounded by a constant with high probability. Combining this with contraction results below, we can
 458 show that all the iterates stay in Ω with high probability.
 459

460 5.2 CONTRACTION
 461

462 Next, we aim to show contraction, i.e., $\|x_t^m - x_t^n\|$ will not diverge during local iterations with
 463 high probability. This property is crucial for showing the benefits of local updates in distributed
 464 optimization. However, different from Woodworth et al. (2020a); Khaled et al. (2020), the update of
 465 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
 466 moving average (EMA) and coordinate-wise product. Thus, the weak monotonicity (4.6) can not be
 467 directly applied as in standard analysis of gradient descent. This will further impede contraction.
 468

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

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

$$\mathcal{A}_{t,1} := \left\{ \beta_2^{K/2} \preceq H_{r(t)}^{-1} H_t^m \preceq 1 + (1 - \beta_2)B \text{ and for all } m \in [M] \right\}, \tag{5.8}$$

$$\mathcal{A}_{t,2} := \left\{ \|H_{r(t)}((H_t^m)^{-1} - (H_t^n)^{-1})\| \leq (1 - \beta_2)B_1 \text{ for all } m, n \in [M] \right\}, \tag{5.9}$$

472 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 $B = \tilde{\mathcal{O}}(K), B_1 = \tilde{\mathcal{O}}(K)$, it holds that $\mathbb{P}(E_{t,1}) \geq \mathbb{P}(E_t) - \delta/(4T), \mathbb{P}(E_{t,2}) \geq \mathbb{P}(E_{t,1}) - \delta/(4T)$.

473 Event $\mathcal{A}_{t,1}$ implies the denominator of each worker during local iterations tends to be stagnant
 474 and close to the averaged one after communication. Event $\mathcal{A}_{t,2}$ suggests the denominator at each
 475 worker is close to each other. Note that when there is no noise, all the workers will be exactly
 476 the same and then event $\mathcal{A}_{t,2}$ will always hold. Therefore, although $\mathcal{A}_{t,2}$ seems to be implied by
 477 $\mathcal{A}_{t,1}$, we will be able to take $B_1 \ll B$ as long as $\sigma \ll 1$ by handling them separately. The key
 478 idea to prove Lemma 5 is to control the magnitude of the EMA of squared stochastic gradients, i.e.,
 479

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

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

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

$$\|x_{t+1}^m - x_{t+1}^n\|_{H_r}^2 \approx \|x_t^m - x_t^n\|_{H_r}^2 - 2\eta \langle x_t^m - x_t^n, u_t^m - u_t^n \rangle + \mathcal{O}(\eta^2). \quad (5.10)$$

By the definition of u_t^m , however, it would be influenced by noises from past stochastic gradients. In 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 and apply Bernstein's inequality. This is the reason why we introduce the auxiliary sequence $\{z_t^m\}$ defined in (5.5). Fortunately, noticing that $x_t^m - x_t^n \in \text{conv}(\{z_j^m - z_j^n\}_{j \leq t})$, it suffices to show that $\|z_t^m - z_t^n\|$ will not get too large with high probability.

Lemma 6 (Informal). Define probabilistic event

$$\mathcal{A}_{t,3} := \left\{ \|z_{t+1}^m - z_{t+1}^n\|_{H_r}^2 \leq \frac{\eta^2 \sigma^2}{\lambda} KA, \sum_{j=rK}^t \|\widehat{g}_j^m\|^2 \leq \frac{(1-\beta_1)^2 \sigma^2 A}{2^{12}(1-\beta_2)^2 B_1^2} \text{ for all } m, n \in [M] \right\}, \quad (5.11)$$

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

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

$$\|z_{t+1}^m - z_{t+1}^n\|_{H_r}^2 \approx \|z_t^m - z_t^n\|_{H_r}^2 - 2\eta \langle z_t^m - z_t^n, \widehat{g}_t^m - \widehat{g}_t^n \rangle + \mathcal{O}(\eta^2). \quad (5.12)$$

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 $\nabla f(z_t^m) - \nabla f(z_t^n)$ 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 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 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 Bernstein's inequality.

5.3 DESCENT

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

$$\mathcal{A}_{t,4} := \left\{ f_\gamma^{H_{r(t+1)}}(\bar{z}_{t+1}) - f_* + \frac{\eta}{12} \sum_{j=0}^t \|\nabla f_\gamma^{H_{r(j)}}(\bar{z}_j)\|_{H_{r(j)}^{-1}}^2 \leq 2\Delta \right\}. \quad (5.13)$$

We proceed with (5.7) and control the stochastic noise term by subtracting its expectation to construct a martingale and apply Bernstein's inequality. Its expectation can be bounded by properties of clipping operator and variance bound. As for the discrepancy overhead, we apply the upper bound of $\|x_j^m - x_j^n\|^2$, which is induced by event E_t and utilize the $\mathcal{O}(\eta^2)$ bound on $\|\bar{z}_j - \bar{x}_j\|^2$. Therefore, thanks to all the foundations beforehand, we are able to bound each of these terms.

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

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 carefully choosing the learning rate η , we complete the proof of Theorem 3.

6 CONCLUSION

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

540 REFERENCES
541

- 542 Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and
543 Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC
544 conference on computer and communications security*, pp. 308–318, 2016.
- 545 Ahmet Alacaoglu, Yura Malitsky, and Volkan Cevher. Convergence of adaptive algorithms for
546 weakly convex constrained optimization. *arXiv preprint arXiv:2006.06650*, 2020.
- 547 Wojciech Anyszka, Kaja Gruntkowska, Alexander Tyurin, and Peter Richtárik. Tighter performance
548 theory of fedexprox. *arXiv preprint arXiv:2410.15368*, 2024.
- 549
- 550 Yossi Arjevani, Yair Carmon, John C Duchi, Dylan J Foster, Nathan Srebro, and Blake Woodworth.
551 Lower bounds for non-convex stochastic optimization. *Mathematical Programming*, 199(1-2):
552 165–214, 2023.
- 553 George Bennett. Probability inequalities for the sum of independent random variables. *Journal of
554 the American Statistical Association*, 57(297):33–45, 1962.
- 555
- 556 Jeremy Bernstein, Yu-Xiang Wang, Kamyar Azizzadenesheli, and Animashree Anandkumar.
557 signsdg: Compressed optimisation for non-convex problems. In *International Conference on
558 Machine Learning*, pp. 560–569. PMLR, 2018.
- 559
- 560 Xiangyi Chen, Xiaoyun Li, and Ping Li. Toward communication efficient adaptive gradient method.
561 In *Proceedings of the 2020 ACM-IMS on Foundations of Data Science Conference*, pp. 119–128,
562 2020.
- 563
- 564 Ziheng Cheng, Xinmeng Huang, Pengfei Wu, and Kun Yuan. Momentum benefits non-iid federated
learning simply and provably. *arXiv preprint arXiv:2306.16504*, 2023.
- 565
- 566 Michael Crawshaw, Mingrui Liu, Francesco Orabona, Wei Zhang, and Zhenxun Zhuang. Robustness
567 to unbounded smoothness of generalized signsdg. *Advances in Neural Information Processing
568 Systems*, 35:9955–9968, 2022.
- 569
- 570 Ashok Cutkosky and Harsh Mehta. High-probability bounds for non-convex stochastic optimization
with heavy tails. *Advances in Neural Information Processing Systems*, 34:4883–4895, 2021.
- 571
- 572 Damek Davis and Dmitriy Drusvyatskiy. Stochastic model-based minimization of weakly convex
functions. *SIAM Journal on Optimization*, 29(1):207–239, 2019.
- 573
- 574 Alexandre Défossez, Léon Bottou, Francis Bach, and Nicolas Usunier. A simple convergence proof
of adam and adagrad. *arXiv preprint arXiv:2003.02395*, 2020.
- 575
- 576 Qi Deng and Wenzhi Gao. Minibatch and momentum model-based methods for stochastic weakly
577 convex optimization. *Advances in Neural Information Processing Systems*, 34:23115–23127,
578 2021.
- 579
- 580 Arthur Douillard, Qixuan Feng, Andrei A Rusu, Rachita Chhaparia, Yani Donchev, Adhiguna
Kuncoro, Marc’Aurelio Ranzato, Arthur Szlam, and Jiajun Shen. Diloco: Distributed low-
581 communication training of language models. *arXiv preprint arXiv:2311.08105*, 2023.
- 582
- 583 Matthew Faw, Litu Rout, Constantine Caramanis, and Sanjay Shakkottai. Beyond uniform smooth-
584 ness: A stopped analysis of adaptive sgd. In *The Thirty Sixth Annual Conference on Learning
585 Theory*, pp. 89–160. PMLR, 2023.
- 586
- 587 David A Freedman. On tail probabilities for martingales. *the Annals of Probability*, pp. 100–118,
1975.
- 588
- 589 Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, and Yann N Dauphin. Convolutional
590 sequence to sequence learning. In *International conference on machine learning*, pp. 1243–1252.
PMLR, 2017.
- 591
- 592 Margalit R Glasgow, Honglin Yuan, and Tengyu Ma. Sharp bounds for federated averaging (lo-
593 cal sgd) and continuous perspective. In *International Conference on Artificial Intelligence and
Statistics*, pp. 9050–9090. PMLR, 2022.

- 594 Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016.
 595
- 596 Eduard Gorbunov, Marina Danilova, and Alexander Gasnikov. Stochastic optimization with heavy-
 597 tailed noise via accelerated gradient clipping. *Advances in Neural Information Processing Sys-*
 598 *tems*, 33:15042–15053, 2020.
- 599 Eduard Gorbunov, Abdurakhmon Sadiev, Marina Danilova, Samuel Horváth, Gauthier Gidel, Pavel
 600 Dvurechensky, Alexander Gasnikov, and Peter Richtárik. High-probability convergence for com-
 601 posite and distributed stochastic minimization and variational inequalities with heavy-tailed noise.
 602 *arXiv preprint arXiv:2310.01860*, 2023.
 603
- 604 Zhishuai Guo, Yi Xu, Wotao Yin, Rong Jin, and Tianbao Yang. A novel convergence analysis for
 605 algorithms of the adam family. *arXiv preprint arXiv:2112.03459*, 2021.
- 606 Nicholas JA Harvey, Christopher Liaw, and Sikander Randhawa. Simple and optimal
 607 high-probability bounds for strongly-convex stochastic gradient descent. *arXiv preprint*
 608 *arXiv:1909.00843*, 2019.
 609
- 610 Sai Praneeth Karimireddy, Martin Jaggi, Satyen Kale, Mehryar Mohri, Sashank J Reddi, Sebas-
 611 tian U Stich, and Ananda Theertha Suresh. Mime: Mimicking centralized stochastic algorithms
 612 in federated learning. *arXiv preprint arXiv:2008.03606*, 2020a.
- 613 Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and
 614 Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In
 615 *International conference on machine learning*, pp. 5132–5143. PMLR, 2020b.
 616
- 617 Sai Praneeth Karimireddy, Lie He, and Martin Jaggi. Learning from history for byzantine robust
 618 optimization. In *International Conference on Machine Learning*, pp. 5311–5319. PMLR, 2021.
- 619 Ahmed Khaled, Konstantin Mishchenko, and Peter Richtárik. Tighter theory for local sgd on identi-
 620 cal and heterogeneous data. In *International Conference on Artificial Intelligence and Statistics*,
 621 pp. 4519–4529. PMLR, 2020.
- 622 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
 623 *arXiv:1412.6980*, 2014.
 624
- 625 Frederik Kunstner, Jacques Chen, Jonathan Wilder Lavington, and Mark Schmidt. Noise is not the
 626 main factor behind the gap between sgd and adam on transformers, but sign descent might be.
 627 *arXiv preprint arXiv:2304.13960*, 2023.
 628
- 629 Hanmin Li, Kirill Acharya, and Peter Richtarik. The power of extrapolation in federated learning.
 630 *arXiv preprint arXiv:2405.13766*, 2024a.
- 631 Haochuan Li, Jian Qian, Yi Tian, Alexander Rakhlin, and Ali Jadbabaie. Convex and non-convex
 632 optimization under generalized smoothness. *Advances in Neural Information Processing Systems*,
 633 36, 2024b.
 634
- 635 Haochuan Li, Alexander Rakhlin, and Ali Jadbabaie. Convergence of adam under relaxed assump-
 636 tions. *Advances in Neural Information Processing Systems*, 36, 2024c.
- 637 Xiaoyu Li and Francesco Orabona. A high probability analysis of adaptive sgd with momentum.
 638 *arXiv preprint arXiv:2007.14294*, 2020.
 639
- 640 Hong Liu, Zhiyuan Li, David Hall, Percy Liang, and Tengyu Ma. Sophia: A scalable stochastic
 641 second-order optimizer for language model pre-training. *arXiv preprint arXiv:2305.14342*, 2023.
- 642 Mingrui Liu, Zhenxun Zhuang, Yunwen Lei, and Chunyang Liao. A communication-efficient dis-
 643 tributed gradient clipping algorithm for training deep neural networks. *Advances in Neural Infor-*
 644 *mation Processing Systems*, 35:26204–26217, 2022.
 645
- 646 Wei Liu, Li Chen, Yunfei Chen, and Wenyi Zhang. Accelerating federated learning via momentum
 647 gradient descent. *IEEE Transactions on Parallel and Distributed Systems*, 31(8):1754–1766,
 2020a.

- 648 Yanli Liu, Yuan Gao, and Wotao Yin. An improved analysis of stochastic gradient descent with
 649 momentum. *Advances in Neural Information Processing Systems*, 33:18261–18271, 2020b.
 650
- 651 Liangchen Luo, Yuanhao Xiong, Yan Liu, and Xu Sun. Adaptive gradient methods with dynamic
 652 bound of learning rate. *arXiv preprint arXiv:1902.09843*, 2019.
- 653 Vien Mai and Mikael Johansson. Convergence of a stochastic gradient method with momentum
 654 for non-smooth non-convex optimization. In *International conference on machine learning*, pp.
 655 6630–6639. PMLR, 2020.
- 656
- 657 Vien V Mai and Mikael Johansson. Stability and convergence of stochastic gradient clipping: Be-
 658 yond lipschitz continuity and smoothness. In *International Conference on Machine Learning*, pp.
 659 7325–7335. PMLR, 2021.
- 660 Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas.
 661 Communication-efficient learning of deep networks from decentralized data. In *Artificial intelli-
 662 gence and statistics*, pp. 1273–1282. PMLR, 2017.
- 663
- 664 Stephen Merity, Nitish Shirish Keskar, and Richard Socher. Regularizing and optimizing lstm lan-
 665 guage models. *arXiv preprint arXiv:1708.02182*, 2017.
- 666
- 667 Yan Pan and Yuanzhi Li. Toward understanding why adam converges faster than sgd for transform-
 668 ers. *arXiv preprint arXiv:2306.00204*, 2023.
- 669
- 670 Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent neural
 671 networks. In *International conference on machine learning*, pp. 1310–1318. Pmlr, 2013.
- 672
- 673 Kumar Kshitij Patel, Lingxiao Wang, Blake E Woodworth, Brian Bullins, and Nati Srebro. Towards
 674 optimal communication complexity in distributed non-convex optimization. *Advances in Neural
 675 Information Processing Systems*, 35:13316–13328, 2022.
- 676
- 677 V. V. Petrov. Moments of sums of independent random variables. *Journal of Soviet Mathematics*,
 678 61(1):1905–1906, Aug 1992. ISSN 1573-8795. doi: 10.1007/BF01362802. URL <https://doi.org/10.1007/BF01362802>.
- 679
- 680 Sashank Reddi, Zachary Charles, Manzil Zaheer, Zachary Garrett, Keith Rush, Jakub Konečný,
 681 Sanjiv Kumar, and H Brendan McMahan. Adaptive federated optimization. *arXiv preprint
 682 arXiv:2003.00295*, 2020.
- 683
- 684 Sashank J Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. *arXiv
 685 preprint arXiv:1904.09237*, 2019.
- 686
- 687 Abdurakhmon Sadiev, Marina Danilova, Eduard Gorbunov, Samuel Horváth, Gauthier Gidel, Pavel
 688 Dvurechensky, Alexander Gasnikov, and Peter Richtárik. High-probability bounds for stochas-
 689 tic optimization and variational inequalities: the case of unbounded variance. In *Proceed-
 690 ings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings
 691 of Machine Learning Research*, pp. 29563–29648. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/sadiev23a.html>.
- 692
- 693 Othmane Sebbouh, Robert M Gower, and Aaron Defazio. Almost sure convergence rates for stochas-
 694 tic gradient descent and stochastic heavy ball. In *Conference on Learning Theory*, pp. 3935–3971.
 695 PMLR, 2021.
- 696
- 697 Naichen Shi, Dawei Li, Mingyi Hong, and Ruoyu Sun. Rmsprop converges with proper hyper-
 698 parameter. In *International Conference on Learning Representations*, 2020.
- 699
- 700 Bohan Wang, Yushun Zhang, Huishuai Zhang, Qi Meng, Zhi-Ming Ma, Tie-Yan Liu, and Wei Chen.
 701 Provable adaptivity in adam. *arXiv preprint arXiv:2208.09900*, 2022.
- 702
- 703 Bohan Wang, Jingwen Fu, Huishuai Zhang, Nanning Zheng, and Wei Chen. Closing the gap between
 704 the upper bound and lower bound of adam’s iteration complexity. *Advances in Neural Information
 705 Processing Systems*, 36, 2024.

- 702 Jianyu Wang, Vinayak Tantia, Nicolas Ballas, and Michael Rabbat. Slowmo: Improving
 703 communication-efficient distributed sgd with slow momentum. *arXiv preprint arXiv:1910.00643*,
 704 2019.
- 705 Jianyu Wang, Zheng Xu, Zachary Garrett, Zachary Charles, Luyang Liu, and Gauri Joshi. Local
 706 adaptivity in federated learning: Convergence and consistency. *arXiv preprint arXiv:2106.02305*,
 707 2021.
- 708 Jue Wang, Yucheng Lu, Binhang Yuan, Beidi Chen, Percy Liang, Christopher De Sa, Christopher
 709 Re, and Ce Zhang. Cocktailsd: fine-tuning foundation models over 500mbps networks. In
 710 *International Conference on Machine Learning*, pp. 36058–36076. PMLR, 2023.
- 711 Jianqiao Wangni, Jialei Wang, Ji Liu, and Tong Zhang. Gradient sparsification for communication-
 712 efficient distributed optimization. *Advances in Neural Information Processing Systems*, 31, 2018.
- 713 Blake Woodworth, Kumar Kshitij Patel, Sebastian Stich, Zhen Dai, Brian Bullins, Brendan Mcma-
 714 han, Ohad Shamir, and Nathan Srebro. Is local sgd better than minibatch sgd? In *International
 715 Conference on Machine Learning*, pp. 10334–10343. PMLR, 2020a.
- 716 Blake E Woodworth, Kumar Kshitij Patel, and Nati Srebro. Minibatch vs local sgd for heterogeneous
 717 distributed learning. *Advances in Neural Information Processing Systems*, 33:6281–6292, 2020b.
- 718 Jing Xu, Sen Wang, Liwei Wang, and Andrew Chi-Chih Yao. Fedcm: Federated learning with
 719 client-level momentum. *arXiv preprint arXiv:2106.10874*, 2021.
- 720 Haibo Yang, Minghong Fang, and Jia Liu. Achieving linear speedup with partial worker participa-
 721 tion in non-iid federated learning. *arXiv preprint arXiv:2101.11203*, 2021.
- 722 Hao Yu, Rong Jin, and Sen Yang. On the linear speedup analysis of communication efficient mo-
 723 mentum sgd for distributed non-convex optimization. In *International Conference on Machine
 724 Learning*, pp. 7184–7193. PMLR, 2019.
- 725 Honglin Yuan and Tengyu Ma. Federated accelerated stochastic gradient descent. *Advances in
 726 Neural Information Processing Systems*, 33:5332–5344, 2020.
- 727 Jingzhao Zhang, Tianxing He, Suvrit Sra, and Ali Jadbabaie. Why gradient clipping accelerates
 728 training: A theoretical justification for adaptivity. *arXiv preprint arXiv:1905.11881*, 2019.
- 729 Jingzhao Zhang, Sai Praneeth Karimireddy, Andreas Veit, Seungyeon Kim, Sashank Reddi, Sanjiv
 730 Kumar, and Suvrit Sra. Why are adaptive methods good for attention models? *Advances in
 731 Neural Information Processing Systems*, 33:15383–15393, 2020.
- 732 Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuhui Chen, Christo-
 733 pher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer
 734 language models. *arXiv preprint arXiv:2205.01068*, 2022a.
- 735 Yushun Zhang, Congliang Chen, Naichen Shi, Ruoyu Sun, and Zhi-Quan Luo. Adam can converge
 736 without any modification on update rules. *Advances in Neural Information Processing Systems*,
 737 35:28386–28399, 2022b.
- 738 Weijie Zhao, Xuewu Jiao, Mingqing Hu, Xiaoyun Li, Xiangyu Zhang, and Ping Li. Communication-
 739 efficient terabyte-scale model training framework for online advertising. *arXiv preprint
 740 arXiv:2201.05500*, 2022.
- 741 Fangyu Zou, Li Shen, Zequn Jie, Weizhong Zhang, and Wei Liu. A sufficient condition for conver-
 742 gences of adam and rmsprop. In *Proceedings of the IEEE/CVF Conference on computer vision
 743 and pattern recognition*, pp. 11127–11135, 2019.
- 744
- 745
- 746
- 747
- 748
- 749
- 750
- 751
- 752
- 753
- 754
- 755

756 **A ADDITIONAL RELATED WORK**
 757

758 **Gradient clipping.** Pascanu et al. (2013) first proposed gradient clipping technique to address the
 759 issue of exploding gradient problem of deep neural networks. Since then, it has become standard
 760 practice in the training of language models (Gehring et al., 2017; Merity et al., 2017; Zhang et al.,
 761 2022a; Liu et al., 2023). Furthermore, from theoretical perspective, gradient clipping is also used
 762 for multiple purposes, including differential privacy (Abadi et al., 2016), distributed optimization
 763 (Karimireddy et al., 2021; Liu et al., 2022), heavy-tailed noise (Zhang et al., 2020).

764 **Generalized smoothness.** The generalized smoothness condition was initially proposed by
 765 (Zhang et al., 2019) to justify gradient clipping, and was called (L_0, L_1) -smoothness. The empirical
 766 evidence therein illustrated that the norm of Hessian matrix of language models depends linearly
 767 on the magnitude of gradient, contradicting the standard L -smoothness. A recent work (Li et al.,
 768 2024b) further generalized this condition to ℓ -smoothness and proved convergence of classical SGD
 769 in this setting. Apart from bounding the Hessian through gradient, Sadiev et al. (2023) proposed to
 770 assume that the norm of Hessian is uniformly bounded in certain subset of whole space, in order to
 771 get high probability bounds for (accelerated) clipped-SGD. Gorbunov et al. (2023) further extended
 772 this setting to composite and distributed optimization without local updates. Here we follow the
 773 setting of (Sadiev et al., 2023) since (L_0, L_1) -smoothness would reduce to it in most cases. See
 774 Section 3.1 for details.

775 **B TECHNICAL LEMMAS**
 776

777 **Lemma B.1** ((Bennett, 1962; Freedman, 1975)). *Let the sequence of random variables $\{X_i\}_{i \geq 1}$ 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 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 exists deterministic constant $c > 0$ such that $|X_i| \leq c$ almost surely for all $i \geq 1$. Then for all $b > 0$, $V > 0$ and $n \geq 1$,*

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

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

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

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

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

794 *Proof.* The first claim is from (Sadiev et al., 2023) and we show the proof here for completeness.
 795 To start the proof, we introduce two indicator random variables. Let

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

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

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

803 we obtain

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

810 Since $1 - \frac{\rho}{|X|} \in (0, 1)$ when $\chi \neq 0$, we derive
 811

$$\begin{aligned} 812 \quad |\mathbb{E}[\tilde{X}] - x| &\leq \mathbb{E}[\chi|X|] \\ 813 \quad &\leq \mathbb{E}[\eta|X|] \\ 814 \quad &\leq \mathbb{E}[\eta|X - x| + \eta|x|] \\ 815 \quad &\leq (\mathbb{E}[|X - x|^\alpha])^{\frac{1}{\alpha}} (\mathbb{E}[\eta^{\frac{\alpha}{\alpha-1}}])^{\frac{\alpha-1}{\alpha}} + |x|\mathbb{E}[\eta] \\ 816 \quad &\stackrel{\eta \in \{0,1\}}{\leq} \sigma(\mathbb{E}[\eta])^{\frac{\alpha-1}{\alpha}} + \frac{\rho}{2}\mathbb{E}[\eta], \\ 817 \quad & \\ 818 \quad & \\ 819 \end{aligned} \tag{B.7}$$

820 By Markov's inequality,

$$\begin{aligned} 821 \quad \mathbb{E}[\eta] &= \mathbb{P}\left\{|X - x|^\alpha > \frac{\rho^\alpha}{2^\alpha}\right\} \\ 822 \quad &\leq \frac{2^\alpha}{\rho^\alpha} \mathbb{E}[|X - x|^\alpha] \\ 823 \quad &\leq \left(\frac{2\sigma}{\rho}\right)^\alpha. \\ 824 \quad & \\ 825 \quad & \\ 826 \quad & \\ 827 \end{aligned} \tag{B.8}$$

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

$$829 \quad |\mathbb{E}[\tilde{X}] - x| \leq \sigma\left(\frac{2\sigma}{\rho}\right)^{\alpha-1} + \frac{\rho}{2}\left(\frac{2\sigma}{\rho}\right)^\alpha = \frac{2^\alpha\sigma^\alpha}{\rho^{\alpha-1}}. \tag{B.9}$$

830 For the second part, since

$$831 \quad |\tilde{X} - x| = |\text{clip}(X, \rho) - \text{clip}(x, \rho)| \leq |X - x|, \tag{B.10}$$

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

$$\begin{aligned} 833 \quad \mathbb{E}|\tilde{X} - \mathbb{E}[\tilde{X}]|^{\alpha} &\leq q^{1-\alpha}\mathbb{E}|\tilde{X} - x|^\alpha + (1-q)^{1-\alpha}|\mathbb{E}[\tilde{X}] - x|^\alpha \\ 834 \quad &\leq q^{1-\alpha}\sigma^\alpha + (1-q)^{1-\alpha}\left(\frac{(2\sigma)^\alpha}{\rho^{\alpha-1}}\right)^\alpha. \\ 835 \quad & \\ 836 \quad & \\ 837 \quad & \\ 838 \quad & \\ 839 \end{aligned} \tag{B.11}$$

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

$$841 \quad \mathbb{E}|\tilde{X} - \mathbb{E}[\tilde{X}]|^{\alpha} \leq \left(\sigma + \frac{(2\sigma)^\alpha}{\rho^{\alpha-1}}\right)^\alpha \leq (2\sigma)^\alpha. \tag{B.12}$$

842 This completes the proof. □

843 **Lemma B.3.** For M independent random vectors $X_1, \dots, X_M \in \mathbb{R}^d$ such that $\mathbb{E}[X_m] = 0$,
 844 $\mathbb{E}[\|X_m\|^4] \leq \sigma^4$, the following holds

$$845 \quad \mathbb{E}[\|\mathbb{E}_m X_m\|^2]^2 \leq \frac{4\sigma^4}{M^2}. \tag{B.13}$$

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

$$\begin{aligned} 847 \quad \mathbb{E}[\|\mathbb{E}_m X_m\|^2]^2 &\leq \mathbb{E}\left[\frac{1}{M^2} \sum_m \|X_m\|^2 + \frac{2}{M^2} \sum_{m < n} \langle X_m, X_n \rangle\right]^2 \\ 848 \quad &= \mathbb{E}\left[\frac{1}{M^2} \sum_m \|X_m\|^2\right]^2 + \mathbb{E}\left[\frac{2}{M^2} \sum_{m < n} \langle X_m, X_n \rangle\right]^2 \\ 849 \quad &\leq \frac{\sigma^4}{M^2} + \frac{4}{M^4} \mathbb{E} \sum_{m < n} \langle X_m, X_n \rangle^2 \\ 850 \quad &\leq \frac{4\sigma^4}{M^2}. \\ 851 \quad & \\ 852 \quad & \\ 853 \quad & \\ 854 \quad & \\ 855 \quad & \\ 856 \quad & \\ 857 \quad & \\ 858 \quad & \\ 859 \quad & \\ 860 \quad & \\ 861 \quad & \\ 862 \quad & \\ 863 \end{aligned} \tag{B.14}$$

□

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

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

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

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

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

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

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

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

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

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

C PROOF OF LOCAL SGDM

We restate the *Local SGDM* algorithm here.

Algorithm 2 Local SGDM

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

Require: initial model x_0 , learning rate η , momentum
 Set $x_{0,0}^m \equiv x_0$, $y_{0,-1}^m \equiv 0$ 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

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

$$g_{r,k}^m = \nabla F(x_{r,k}^m; \xi_{r,k}^m), \quad \widehat{g_{r,k}^m} = \text{clip}(g_{r,k}^m, \rho) \quad \triangleright \text{Compute clipped stochastic gradient}$$

$$\widehat{\omega}_{r,k}^m = \varrho \cdot \widehat{\omega}_{r,k-1}^m + (1 - \varrho) \cdot \widehat{g_{r,k}^m} \quad \triangleright \text{Update momentum}$$

$$u_{r,k}^m = \beta_1 u_{r,k-1}^m + (1 - \beta_1) g_{r,k}^m$$

▷ Update model

nd for

for [View Details](#) [Edit](#) [Delete](#)

C.1 OVERVIEW AND MAIN THEOREM

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 $k(t) \leq K-1$. We omit the dependence on t and let $r = r(t), k = k(t)$ through out the proof if not causing confusion. 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 following update rule:

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

$$x_{t+1}^m = \begin{cases} x_t^m - \eta u_t^m & \text{if } t \bmod K \not\equiv -1, \\ \bar{x}_t - \eta \bar{u}_t & \text{otherwise.} \end{cases} \quad (\text{C.2})$$

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

$$z_{t+1}^m = \begin{cases} \frac{1}{1-\beta_1} x_{t+1}^m - \frac{\beta_1}{1-\beta_1} x_t^m & \text{if } t \bmod K \not\equiv -1, \\ \frac{1}{1-\beta_1} x_{t+1}^m - \frac{\beta_1}{1-\beta_1} \bar{x}_t & \text{otherwise.} \end{cases} \quad (\text{C.3})$$

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

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

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

Besides, let

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

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

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 $\eta \leq \min \left\{ \frac{(1-\beta_1)^2}{2L}, \frac{D_0}{4\sigma \sqrt{KA}} \right\}$ and $\rho \geq \max\{3\sigma_\infty, 2G_\infty\}$, then the following holds:

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

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

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

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

If

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

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

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

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

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

if

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

Here

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

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

When $t = 0$, this is trivial. Assume that the statement is true for some $t \leq T - 1$. We aim to prove 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 (C.12) hold. Hence we have

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

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

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

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

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

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

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

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

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

□

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

1026 **Theorem C.4** (Complete version of Theorem 1). *Under the conditions of Theorem C.3 and $\mu > 0$,
1027 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 $K \gtrsim \log \frac{MT}{\delta} \left(\frac{\|\boldsymbol{\sigma}\|_{2\alpha} d^{\frac{1}{2} - \frac{1}{2\alpha}}}{\sigma}\right)^{\frac{2\alpha}{\alpha-2}}$.
1028 Then with probability no less than $1 - \delta$, Local SGDM with optimal η, ρ yields $f(\hat{x}) - f_* \leq \varepsilon$, if
1029*

$$1030 T \gtrsim \log \frac{\mu D_0^2}{\varepsilon} \left[\frac{L}{\mu} + \frac{\sigma^2}{\mu M \varepsilon} \log \frac{T}{\delta} + \sqrt{\frac{L \sigma^2 K \log \frac{MT}{\delta}}{\mu^2 \varepsilon}} + \sqrt{\frac{Ld}{\mu^2 \varepsilon} \log \frac{MT}{\delta} \left(\frac{\|\boldsymbol{\sigma}\|_{2\alpha}^{2\alpha}}{\mu\varepsilon}\right)^{\frac{1}{2(\alpha-1)}}} \right]. \quad (C.18)$$

1031 And equivalently, let $\kappa := L/\mu$,

$$1032 f(\hat{x}) - f_* \lesssim \exp \left(-\Theta \left(\frac{\mu K R}{L} \right) \right) + \frac{\sigma^2 \log(MKR)}{\mu MKR} \log \frac{KR}{\delta} \\ 1033 + \frac{L \sigma^2 \log^2(KR)}{\mu^2 KR^2} \log \frac{MKR}{\delta} + \frac{\|\boldsymbol{\sigma}\|_{2\alpha}^2 (\kappa d)^{\frac{\alpha-1}{\alpha}}}{\mu} \left(\frac{\log \frac{MKR}{\delta}}{KR} \right)^{\frac{2(\alpha-1)}{\alpha}}. \quad (C.19)$$

1044 *Proof.* Plug the definition of A in (C.11),

$$1045 T \gtrsim \log \frac{\mu D_0^2}{\varepsilon} \left[\frac{L}{\mu} + \frac{\sigma^2}{\mu M \varepsilon} \log \frac{T}{\delta} + \sqrt{\frac{L \sigma^2 K \log \frac{MT}{\delta}}{\mu^2 \varepsilon}} + \frac{\rho \sqrt{d}}{\sqrt{\mu \varepsilon}} \log \frac{T}{\delta} \right] \\ 1046 + \log \frac{\mu D_0^2}{\varepsilon} \sqrt{\frac{LK}{\mu^2 \varepsilon}} \sqrt{\frac{\rho^2 d}{K} \log^2 \frac{MT}{\delta} + \frac{K \|2\boldsymbol{\sigma}\|_{2\alpha}^{2\alpha}}{\rho^{2(\alpha-1)}}} \\ 1047 \asymp \log \frac{\mu D_0^2}{\varepsilon} \left[\frac{L}{\mu} + \frac{\sigma^2}{\mu M \varepsilon} \log \frac{T}{\delta} + \sqrt{\frac{L \sigma^2 K \log \frac{MT}{\delta}}{\mu^2 \varepsilon}} \right] \\ 1048 + \log \frac{\mu D_0^2}{\varepsilon} \sqrt{\frac{LK}{\mu^2 \varepsilon}} \sqrt{\frac{\rho^2 d}{K} \log^2 \frac{MT}{\delta} + \frac{K \|2\boldsymbol{\sigma}\|_{2\alpha}^{2\alpha}}{\rho^{2(\alpha-1)}}}. \quad (C.20)$$

1058 Hence the optimal ρ is given by

$$1059 \rho \asymp \max \left\{ \|\boldsymbol{\sigma}\|_{2\alpha} \left(\frac{K}{\sqrt{d} \log \frac{MT}{\delta}} \right)^{1/\alpha}, \left(\frac{\|\boldsymbol{\sigma}\|_{2\alpha}^{2\alpha}}{\mu \varepsilon} \right)^{\frac{1}{2(\alpha-1)}}, \sigma_\infty, G_\infty \right\}. \quad (C.21)$$

1060 Note that $\left(\frac{\|\boldsymbol{\sigma}\|_{2\alpha}^{2\alpha}}{\mu \varepsilon} \right)^{\frac{1}{2(\alpha-1)}} \gtrsim G_\infty \vee \sigma_\infty$ and this implies

$$1061 T \gtrsim \log \frac{\mu D_0^2}{\varepsilon} \left[\frac{L}{\mu} + \frac{\sigma^2}{\mu M \varepsilon} \log \frac{T}{\delta} + \sqrt{\frac{L \sigma^2 K \log \frac{MT}{\delta}}{\mu^2 \varepsilon}} \right] \\ 1062 + \log \frac{\mu D_0^2}{\varepsilon} \sqrt{\frac{L}{\mu^2 \varepsilon} \cdot \left[\|\boldsymbol{\sigma}\|_{2\alpha}^2 K^{\frac{2}{\alpha}} \left(d \log^2 \frac{MT}{\delta} \right)^{1-\frac{1}{\alpha}} + \left(\frac{\|\boldsymbol{\sigma}\|_{2\alpha}^{2\alpha}}{\mu \varepsilon} \right)^{\frac{1}{2(\alpha-1)}} d \log^2 \frac{MT}{\delta} \right]} \\ 1063 \asymp \log \frac{\mu D_0^2}{\varepsilon} \left[\frac{L}{\mu} + \frac{\sigma^2}{\mu M \varepsilon} \log \frac{T}{\delta} + \sqrt{\frac{L \sigma^2 K \log \frac{MT}{\delta}}{\mu^2 \varepsilon}} + \sqrt{\frac{Ld}{\mu^2 \varepsilon} \log \frac{MT}{\delta} \left(\frac{\|\boldsymbol{\sigma}\|_{2\alpha}^{2\alpha}}{\mu \varepsilon} \right)^{\frac{1}{2(\alpha-1)}}} \right]. \quad (C.22)$$

1064 In the last equation we use $K \gtrsim \log \frac{MT}{\delta} \left(\frac{\|\boldsymbol{\sigma}\|_{2\alpha} d^{\frac{1}{2} - \frac{1}{2\alpha}}}{\sigma} \right)^{\frac{2\alpha}{\alpha-2}}$. This completes the proof. \square

1080
 1081 **Theorem C.5** (Complete version of Theorem 2). *Under the conditions of Theorem C.3 and $\mu = 0$,*
 1082 *assume $1 - \beta_1 = \Omega(1)$, $\left(\frac{D_0\|\boldsymbol{\sigma}\|_{2\alpha}^\alpha}{\varepsilon}\right)^{\frac{1}{\alpha-1}} \gtrsim G_\infty \vee \sigma_\infty$, and $K \gtrsim \log \frac{MT}{\delta} \left(\frac{\|\boldsymbol{\sigma}\|_{2\alpha} d^{\frac{1}{2}-\frac{1}{2\alpha}}}{\sigma}\right)^{\frac{2\alpha}{\alpha-2}}$.*
 1083 *Then with probability no less than $1 - \delta$, Local SGDM with optimal η, ρ yields $f(\hat{x}) - f_* \leq \varepsilon$ if*

1084
 1085
$$T \gtrsim \frac{D_0^2}{\varepsilon} \left[L + \frac{\sigma^2}{M\varepsilon} \log \frac{T}{\delta} + \sqrt{\frac{L\sigma^2 K \log \frac{MT}{\delta}}{\varepsilon}} + \sqrt{\frac{dL}{\varepsilon}} \left(\frac{D_0\|\boldsymbol{\sigma}\|_{2\alpha}^\alpha}{\varepsilon}\right)^{\frac{1}{\alpha-1}} \log \frac{MT}{\delta} \right]. \quad (\text{C.23})$$

1086
 1087 *And equivalently,*

1088
 1089
$$\begin{aligned} f(\hat{x}) - f_* &\lesssim \frac{LD_0^2}{KR} + \frac{\sigma D_0}{\sqrt{MKR}} \log^{\frac{1}{2}} \frac{KR}{\delta} \\ &+ \frac{L^{\frac{1}{3}} \sigma^{\frac{2}{3}} D_0^{\frac{4}{3}}}{K^{\frac{1}{3}} R^{\frac{2}{3}}} \log^{\frac{1}{3}} \frac{MKR}{\delta} + \left(\|\boldsymbol{\sigma}\|_{2\alpha}^{\frac{2\alpha}{\alpha-1}} dLD_0\right)^{\frac{\alpha-1}{3\alpha-1}} D_0 \left(\frac{\log \frac{MKR}{\delta}}{KR}\right)^{\frac{2(\alpha-1)}{3\alpha-1}}. \end{aligned} \quad (\text{C.24})$$

1090
 1091 *Proof.* Plug the definition of A in (C.11),

1092
 1093
$$\begin{aligned} T &\gtrsim \frac{D_0^2}{\varepsilon} \left[L + \frac{\sigma^2}{M\varepsilon} \log \frac{T}{\delta} + \sqrt{\frac{L\sigma^2 K \log \frac{MT}{\delta}}{\varepsilon}} + \frac{\rho\sqrt{d}}{D_0} \log \frac{T}{\delta} \right] \\ &+ \frac{D_0^2}{\varepsilon} \sqrt{\frac{LK}{\varepsilon}} \sqrt{\frac{\rho^2 d}{K} \log^2 \frac{MT}{\delta} + \frac{K\|2\boldsymbol{\sigma}\|_{2\alpha}^{2\alpha}}{\rho^{2(\alpha-1)}}} \\ &\asymp \frac{D_0^2}{\varepsilon} \left[L + \frac{\sigma^2}{M\varepsilon} \log \frac{T}{\delta} + \sqrt{\frac{L\sigma^2 K \log \frac{MT}{\delta}}{\varepsilon}} + \sqrt{\frac{LK}{\varepsilon}} \sqrt{\frac{\rho^2 d}{K} \log^2 \frac{MT}{\delta} + \frac{K\|2\boldsymbol{\sigma}\|_{2\alpha}^{2\alpha}}{\rho^{2(\alpha-1)}}} \right]. \end{aligned} \quad (\text{C.25})$$

1094 Hence the optimal ρ is given by

1095
 1096
$$\rho \asymp \max \left\{ \|\boldsymbol{\sigma}\|_{2\alpha} \left(\frac{K}{\sqrt{d} \log \frac{MT}{\delta}} \right)^{1/\alpha}, \left(\frac{D_0\|\boldsymbol{\sigma}\|_{2\alpha}^\alpha}{\varepsilon} \right)^{\frac{1}{\alpha-1}}, \sigma_\infty, G_\infty \right\}. \quad (\text{C.26})$$

1097 Note that $\left(\frac{D_0\|\boldsymbol{\sigma}\|_{2\alpha}^\alpha}{\varepsilon}\right)^{\frac{1}{\alpha-1}} \gtrsim G_\infty \vee \sigma_\infty$ and this implies

1098
 1099
$$\begin{aligned} T &\gtrsim \frac{D_0^2}{\varepsilon} \left[L + \frac{\sigma^2}{M\varepsilon} \log \frac{T}{\delta} + \sqrt{\frac{L\sigma^2 K \log \frac{MT}{\delta}}{\varepsilon}} \right] \\ &+ \frac{D_0^2}{\varepsilon} \sqrt{\frac{L}{\varepsilon}} \cdot \left[\|\boldsymbol{\sigma}\|_{2\alpha}^2 K^{\frac{2}{\alpha}} \left(d \log^2 \frac{MT}{\delta} \right)^{1-\frac{1}{\alpha}} + \left(\frac{D_0\|\boldsymbol{\sigma}\|_{2\alpha}^\alpha}{\varepsilon} \right)^{\frac{2}{\alpha-1}} d \log^2 \frac{MT}{\delta} \right] \\ &\asymp \frac{D_0^2}{\varepsilon} \left[L + \frac{\sigma^2}{M\varepsilon} \log \frac{T}{\delta} + \sqrt{\frac{L\sigma^2 K \log \frac{MT}{\delta}}{\varepsilon}} + \sqrt{\frac{dL}{\varepsilon}} \left(\frac{D_0\|\boldsymbol{\sigma}\|_{2\alpha}^\alpha}{\varepsilon} \right)^{\frac{1}{\alpha-1}} \log \frac{MT}{\delta} \right]. \end{aligned} \quad (\text{C.27})$$

1100 In the last equation we use $K \gtrsim \log \frac{MT}{\delta} \left(\frac{\|\boldsymbol{\sigma}\|_{2\alpha} d^{\frac{1}{2}-\frac{1}{2\alpha}}}{\sigma} \right)^{\frac{2\alpha}{\alpha-2}}$. Solve ε and we get the upper bound of $f(\hat{x}) - f_*$. This completes the proof. \square

C.2 PRELIMINARIES

1101 In this subsection, we show that event E_t implies all the iterates remain in certain area, so that we
 1102 can apply all kinds of properties of f afterwards.

1134 **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
 1135 $x_j^m, \bar{x}_j, z_j^m, \bar{z}_j \in \Omega$. And $\|x_j^m - x_j^n\| \leq \eta\sigma\sqrt{KA}$ for all m, n .

1137 *Proof.* Event E_t implies that for all $j \leq t$,

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

1140 Hence $\bar{z}_j \in \Omega$, $\|z_j^m - x_*\| \leq \sqrt{3}D_0$ and $z_j^m \in \Omega$. Also, notice that $\bar{x}_j \in \text{conv}\{\bar{z}_i\}_{i \leq j}$ and
 1141 $x_j^m - x_j^n \in \text{conv}\{z_i^m - z_i^n\}_{i \leq j}$. We have

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

1144 Therefore $x_j^m, \bar{x}_j \in \Omega$. This completes the proof. \square

C.3 PROOF OF CONTRACTION LEMMA C.1

1148 In this subsection, we aim to show contraction, *i.e.*, $\|x_t^m - x_t^n\|$ won't be too large during local
 1149 iterations with high probability. This property is crucial for showing the benefits of local updates in
 1150 distributed optimization. However, different from (Woodworth et al., 2020a; Khaled et al., 2020),
 1151 the update of x_t^m is in the direction of momentum u_t^m , which incorporates information from all past
 1152 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
 1153 that $x_t^m - x_t^n \in \text{conv}(\{z_j^m - z_j^n\}_{j \leq t})$, it suffices to show that $\|z_t^m - z_t^n\|$ won't get too large with
 1154 high probability. Besides, the update rule of z_t^m is much easier to handle.

1155 *Proof.* First note that by the upper bound of η , Lemma C.6 holds. Since $z_{t+1}^m = z_t^m - \eta\widehat{g}_t^m$,

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

1163 Event E_t implies $z_t^m, x_t^m \in \Omega$ and thus by $\forall x, y \in \Omega, \langle x - y, \nabla f(x) - \nabla f(y) \rangle \geq \frac{1}{L} \|\nabla f(x) - \nabla f(y)\|^2$,

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

1172 Therefore, for the second and third term in the RHS of (C.30),

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

1178 By the update rule,

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

1188
1189 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
1190
1191
1192 LHS of (C.32) $\leq -\frac{\eta}{L}(S_t - \beta_1 S_{t-1}) + \frac{4\eta L(\eta\beta_1)^2}{1-\beta_1} \left[S_{t-1} + \sum_{j=rK}^{t-1} \beta_1^{t-j-1} [\|\widehat{g}_j^m - \widehat{g}_j^n - \nabla f(x_j^m) + \nabla f(x_j^n)\|^2] \right]$
1193
1194 $= -\frac{\eta}{L}(S_t - S_{t-1}) + \frac{4\eta L(\eta\beta_1)^2}{1-\beta_1} \left[\sum_{j=rK}^{t-1} \beta_1^{t-j-1} [\|\widehat{g}_j^m - \widehat{g}_j^n - \nabla f(x_j^m) + \nabla f(x_j^n)\|^2] \right]$
1195
1196
1197
1198 (C.34)

1199 Then plug in (C.30),
1200

1201 $\|z_{t+1}^m - z_{t+1}^n\|^2 \leq \|z_t^m - z_t^n\|^2 - \frac{\eta}{L}(S_t - S_{t-1})$
1202
1203 $+ \frac{4\eta L(\eta\beta_1)^2}{1-\beta_1} \left[\sum_{j=rK}^{t-1} \beta_1^{t-j-1} [\|\widehat{g}_j^m - \widehat{g}_j^n - \nabla f(x_j^m) + \nabla f(x_j^n)\|^2] \right]$
1204
1205
1206 $+ 2\eta \langle z_t^m - z_t^n, \nabla f(x_t^m) - \nabla f(x_t^n) - \widehat{g}_t^m + \widehat{g}_t^n \rangle + 2\eta^2 \|\widehat{g}_t^m - \widehat{g}_t^n - \nabla f(x_t^m) + \nabla f(x_t^n)\|^2.$
1207 (C.35)

1208 Notice that this recursive bound holds for any $rK \leq i \leq t$. Unroll it and recalculate the coefficients
1209 using $\eta L \leq (1 - \beta_1)^2 / 2$,

1210
1211 $\|z_{t+1}^m - z_{t+1}^n\|^2 + \frac{\eta}{L} S_t \leq \sum_{j=rK}^t 2\eta \langle z_j^m - z_j^n, \nabla f(x_j^m) - \nabla f(x_j^n) - \widehat{g}_j^m + \widehat{g}_j^n \rangle$
1212
1213 $+ \sum_{j=rK}^t 4\eta^2 \|\nabla f(x_j^m) - \nabla f(x_j^n) - \widehat{g}_j^m + \widehat{g}_j^n\|^2$
1214
1215
1216 $\leq \underbrace{\sum_{j=rK}^t 2\eta \langle z_j^m - z_j^n, \mathbb{E}_j[\widehat{g}_j^m - \widehat{g}_j^n] - [\widehat{g}_j^m - \widehat{g}_j^n] \rangle}_{\textcircled{1}: \text{martingale}}$
1217
1218
1219
1220
1221 $+ \underbrace{\sum_{j=rK}^t 2\eta \langle z_j^m - z_j^n, \nabla f(x_j^m) - \nabla f(x_j^n) - \mathbb{E}_j[\widehat{g}_j^m - \widehat{g}_j^n] \rangle}_{\textcircled{2}: \text{clipping bias}}$
1222
1223
1224
1225 $+ \underbrace{\sum_{j=rK}^t 4\eta^2 [\|\nabla f(x_j^m) - \nabla f(x_j^n) - \widehat{g}_j^m + \widehat{g}_j^n\|^2 - \mathbb{E}_j[\|\nabla f(x_j^m) - \nabla f(x_j^n) - [\widehat{g}_j^m - \widehat{g}_j^n]\|^2]}_{\textcircled{3}: \text{martingale}}$
1226
1227
1228
1229
1230
1231 (C.36)

1232 For ①, define
1233

1234 $\zeta_j^{m,n} = \begin{cases} 2\eta \langle z_j^m - z_j^n, \mathbb{E}_j[\widehat{g}_j^m - \widehat{g}_j^n] - [\widehat{g}_j^m - \widehat{g}_j^n] \rangle, & \text{if event } E_j \text{ holds,} \\ 0, & \text{otherwise.} \end{cases}$ (C.37)
1235
1236

1237 Then since event E_j implies $\|z_j^m - z_j^n\| \leq \eta\sigma\sqrt{KA}$,
1238

1239 $|\zeta_j^{m,n}| \leq 2\eta \cdot \eta\sigma\sqrt{KA} \cdot 2\rho\sqrt{d} = 4\eta^2\sigma\rho\sqrt{dKA} \stackrel{\text{def}}{=} c,$ (C.38)
1240

1241 $\text{Var}_j(\zeta_j^{m,n}) \leq 4\eta^2 \cdot \eta^2\sigma^2 KA \cdot 2\sigma^2 = 8\eta^4\sigma^4 KA.$ (C.39)

1242 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
1243
1244
1245

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

1246 For ②,

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

1247 For ③, define

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

1248 Then,

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

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

1249 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
1250
1251
1252

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

1253 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
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403
1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565
1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835
1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889
1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943
1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051
2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105
2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159
2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213
2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263
2264
2265
2266
2267
2268
2269
2270
2271
2272
2273
2274
2275
2276
2277
2278
2279
2280
2281
2282
2283
2284
2285
2286
2287
2288
2289
2290
2291
2292
2293
2294
2295
2296
2297
2298
2299
2300
2301
2302
2303
2304
2305
2306
2307
2308
2309
2310
2311
2312
2313
2314
2315
2316
2317
2318
2319
2320
2321
2322
2323
2324
2325
2326
2327
2328
2329
2330
2331
2332
2333
2334
2335
2336
2337
2338
2339
2340
2341
2342
2343
2344
2345
2346
2347
2348
2349
2350
2351
2352
2353
2354
2355
2356
2357
2358
2359
2360
2361
2362
2363
2364
2365
2366
2367
2368
2369
2370
2371
2372
2373
2374
2375
2376
2377
2378
2379
2380
2381
2382
2383
2384
2385
2386
2387
2388
2389
2390
2391
2392
2393
2394
2395
2396
2397
2398
2399
2400
2401
2402
2403
2404
2405
2406
2407
2408
2409
2410
2411
2412
2413
2414
2415
2416
2417
2418
2419
2420
2421
2422
2423
2424
2425
2426
2427
2428
2429
2430
2431
2432
2433
2434
2435
2436
2437
2438
2439
2440
2441
2442
2443
2444
2445
2446
2447
2448
2449
2450
2451
2452
2453
2454
2455
2456
2457
2458
2459
2460
2461
2462
2463
2464
2465
2466
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
2483
2484
2485
2486
2487
2488
2489
2490
2491
2492
2493
2494
2495
2496
2497
2498
2499
2500
2501
2502
2503
2504
2505
2506
2507
2508
2509
2510
2511
2512
2513
2514
2515
2516
2517
2518
2519
2520
2521
2522
2523
2524
2525
2526
2527
2528
2529
2530
2531
2532
2533
2534
2535
2536
2537
2538
2539
2540
2541
2542
2543
2544
2545
2546
2547
2548
2549
2550
2551
2552
2553
2554
2555
2556
2557
2558
2559
2550
2551
2552
2553
2554
2555
2556
2557
2558
2559
2560
2561
2562
2563
2564
2565
2566
2567
2568
2569
2570
2571
2572
2573
2574
2575
2576
2577
2578
2579
2580
2581
2582
2583
2584
2585
2586
2587
2588
2589
2590
2591
2592
2593
2594
2595
2596
2597
2598
2599
2590
2591
2592
2593
2594
2595
2596
2597
2598
2599
2600
2601
2602
2603
2604
2605
2606
2607
2608
2609
2610
2611
2612
2613
2614
2615
2616
2617
2618
2619
2620
2621
2622
2623
2624
2625
2626
2627
2628
2629
2620
2621
2622
2623
2624
2625
2626
2627
2628
2629
2630
2631
2632
2633
2634
2635
2636
2637
2638
2639
2630
2631
2632
2633
2634
2635
2636
2637
2638
2639
2640
2641
2642
2643
2644
2645
2646
2647
2648
2649
2640
2641
2642
2643
2644
2645
2646
2647
2648
2649
2650
2651
2652
2653
2654
2655
2656
2657
2658
2659
2660
2661
2662
2663
2664
2665
2666
2667
2668
2669
2670
2671
2672
2673
2674
2675
2676
2677
2678
2679
2680
2681
2682
2683
2684
2685
2686
2687
2688
2689
2690
2691
2692
2693
2694
2695
2696
2697
2698
2699
2690
2691
2692
2693
2694
2695
2696
2697
2698
2699
2700
2701
2702
2703
2704
2705
2706
2707
2708
2709
27010
27011
27012
27013
27014
27015
27016
27017
27018
27019
27020
27021
27022
27023
27024
27025
27026
27027
27028
27029
27030
27031
27032
27033
27034
27035
27036
27037
27038
27039
27040
27041
27042
27043
27044
27045
27046
27047
27048
27049
27050
27051
27052
27053
27054
27055
27056
27057
27058
27059
27060
27061
27062
27063
27064
27065
27066
27067
27068
27069
27070
27071
27072
27073
27074
27075
27076
27077
27078
27079
27080
27081
27082
27083
27084
27085
27086
27087
27088
27089
27090
27091
27092
27093
27094
27095
27096
27097
27098
27099
270100
270101
270102
270103
270104
270105
270106
270107
270108
270109
270110
270111
270112
270113
270114
270115
270116
270117
270118
270119
270120
270121
270122
270123
270124
270125
270126
270127
270128
270129
270130
270131
270132
270133
270134
270135
270136
270137
270138
270139
270140
270141
270142
270143
270144
270145
270146

1296 Therefore,

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

For the last term in (C.46),

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

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

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

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,

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

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

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

1350 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
 1351 contraction results implied by event $E_{t,1}$.
 1352

1353 Unroll this recursive bound and re-calculate the coefficients,

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

1364 Simplify Φ_{t+1} term,

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

1381 For the last term,

$$\begin{aligned}
& 4\eta^2 \sum_{j=0}^t (1 - \frac{\eta\mu}{2})^{t-j} \|\mathbb{E}_m[\widehat{g_j^m} - \nabla f(x_j^m)]\|^2 \leq 8\eta^2 \underbrace{\sum_{j=0}^t (1 - \frac{\eta\mu}{2})^{t-j} \left[\|\mathbb{E}_m[\widehat{g_j^m} - \mathbb{E}_j[\widehat{g_j^m}]]\|^2 - \mathbb{E}_j[\|\mathbb{E}_m[\widehat{g_j^m} - \mathbb{E}_j[\widehat{g_j^m}]]\|^2] \right]}_{\textcircled{3}: \text{ martingale}} \\
& + 8\eta^2 \underbrace{\sum_{j=0}^t (1 - \frac{\eta\mu}{2})^{t-j} \mathbb{E}_j[\|\mathbb{E}_m[\widehat{g_j^m} - \mathbb{E}_j[\widehat{g_j^m}]]\|^2]}_{\textcircled{4}: \text{ Lemma B.2}} \\
& + 8\eta^2 \underbrace{\sum_{j=0}^t (1 - \frac{\eta\mu}{2})^{t-j} \|\mathbb{E}_m[\mathbb{E}_j[\widehat{g_j^m}] - \nabla f(x_j^m)]\|^2}_{\textcircled{5}: \text{ clipping bias}}, \tag{C.57}
\end{aligned}$$

1396 we finally get

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

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

1404 For ①, define
 1405

$$1406 \zeta_j = \begin{cases} -2\eta(1 - \frac{\eta\mu}{2})^{t-j} \langle \bar{z}_j - x_*, \mathbb{E}_m[\widehat{g}_j^m] - \mathbb{E}_j[\widehat{g}_j^m] \rangle, & \text{if event } E_j \text{ holds,} \\ 1407 0, & \text{otherwise.} \end{cases} \quad (\text{C.59})$$

1409 Then since event E_j implies $\|\bar{z}_j - x_*\| \leq \sqrt{2}(1 - \frac{\eta\mu}{2})^{j/2}D_0$,

$$1410 |\zeta_j| \leq 2\eta \cdot \sqrt{2}(1 - \frac{\eta\mu}{2})^{t/2}D_0 \cdot 2\rho\sqrt{d} = 4(1 - \frac{\eta\mu}{2})^{t/2}\eta\rho\sqrt{2d}D_0 \stackrel{\text{def}}{=} c, \quad (\text{C.60})$$

$$1414 \text{Var}_j(\zeta_j) \leq 4\eta^2(1 - \frac{\eta\mu}{2})^{2(t-j)} \cdot 2(1 - \frac{\eta\mu}{2})^j D_0^2 \cdot \frac{\sigma^2}{M} = 8(1 - \frac{\eta\mu}{2})^{2t-j} \frac{\eta^2 D_0^2 \sigma^2}{M}. \quad (\text{C.61})$$

1417 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
 1418 no less than

$$1420 1 - 2 \exp\left(\frac{b^2}{2V + 2cb/3}\right) \geq 1 - \frac{\delta}{4T}. \quad (\text{C.62})$$

1423 For ②, since by Lemma B.2,

$$1424 \|\mathbb{E}_j[\widehat{g}_j^m - \nabla f(x_j^m)]\|^2 \leq \frac{\|2\sigma\|_{2\alpha}^{2\alpha}}{\rho^{2(\alpha-1)}}, \quad (\text{C.63})$$

1427 event E_t implies that

$$\begin{aligned} 1429 |②| &\leq 2\eta \sum_{j=0}^t (1 - \frac{\eta\mu}{2})^{t-j} \cdot \sqrt{2}(1 - \frac{\eta\mu}{2})^{j/2}D_0 \cdot \frac{\|2\sigma\|_{2\alpha}^\alpha}{\rho^{\alpha-1}} \\ 1430 &\leq 4\sqrt{2}(1 - \frac{\eta\mu}{2})^{t/2} \frac{D_0 \|2\sigma\|_{2\alpha}^\alpha}{\mu \rho^{\alpha-1}} \\ 1431 &\leq \frac{(1 - \frac{\eta\mu}{2})^{t+1} D_0^2}{5}. \end{aligned} \quad (\text{C.64})$$

1437 Here we use the definition of η and conditions of ρ in (C.12).

1438 For ③, define
 1439

$$1440 \theta_j = \begin{cases} 8\eta^2(1 - \frac{\eta\mu}{2})^{t-j} [\|\mathbb{E}_m[\widehat{g}_j^m] - \mathbb{E}_j[\widehat{g}_j^m]\|^2 - \mathbb{E}_j[\|\mathbb{E}_m[\widehat{g}_j^m] - \mathbb{E}_j[\widehat{g}_j^m]\|^2]], & \text{if event } E_j \text{ holds,} \\ 1441 0, & \text{otherwise.} \end{cases} \quad (\text{C.65})$$

1443 Then

$$1444 |\theta_j| \leq 8\eta^2 \cdot 4\rho^2 d = 32\eta^2 \rho^2 d \stackrel{\text{def}}{=} c, \quad (\text{C.66})$$

$$1446 \text{Var}_j(\theta_j) \leq 64\eta^4(1 - \frac{\eta\mu}{2})^{2(t-j)} \cdot \mathbb{E}_j[\|\mathbb{E}_m[\widehat{g}_j^m] - \mathbb{E}_j[\widehat{g}_j^m]\|^2]^2 \stackrel{\text{Lemma B.3}}{\leq} 64\eta^4(1 - \frac{\eta\mu}{2})^{2(t-j)} \cdot \frac{4(2\sigma)^4}{M^2}. \quad (\text{C.67})$$

1449 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
 1450

$$1453 1 - 2 \exp\left(\frac{b^2}{2V + 2cb/3}\right) \geq 1 - \frac{\delta}{4T}. \quad (\text{C.68})$$

1455 For ④, by Lemma B.2,

$$1456 |④| \leq \frac{16\eta}{\mu} \cdot \frac{\|2\sigma\|_{2\alpha}^{2\alpha}}{\rho^{2(\alpha-1)}} \leq \frac{(1 - \frac{\eta\mu}{2})^{t+1} D_0^2}{5}. \quad (\text{C.69})$$

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

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

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

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

In this case, (C.58) reduces to

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

For $\textcircled{1}$, define

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

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

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

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

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

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

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

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

event E_t implies that

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

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

For $\textcircled{3}$, define

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

Then

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

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

1512 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
 1513
 1514

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

1515 For ④, by Lemma B.2,
 1516
 1517

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

1518 Combine the above claims, with probability no less than $\mathbb{P}(E_{t,1}) - 2 \cdot \frac{\delta}{4T}$, we have $|\textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4}| \leq$
 1519
 1520
 1521
 1522
 1523
 $\frac{4}{5}D_0^2$. By (C.58), these implies

$$\begin{aligned} 1524 & \frac{\eta}{2} \sum_{j=0}^t (f(\bar{z}_j) - f_*) + \|\bar{z}_{t+1} - x_*\|^2 \leq D_0^2 + 16 \left[\eta LKA + \frac{1}{M} \right] \eta^2 \sigma^2 (t+1) + \frac{4}{5} D_0^2 \\ 1525 & \leq 2D_0^2. \end{aligned} \quad (\text{C.83})$$

1526 Therefore, we conclude that $\mathbb{P}(E_{t+1}) \geq \mathbb{P}(E_{t,1}) - \frac{\delta}{2T}$.
 1527
 1528
 1529
 1530
 1531

□

1532
 1533 **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
 1534 conditions in Lemma C.2, then the following holds:
 1535

$$\begin{aligned} 1536 \Lambda_t & \leq \left(1 - \frac{(1-\beta_1)^2}{2}\right) \Lambda_{t-1} + \frac{32\eta^2}{1-\beta_1} \|\nabla f(\bar{z}_t)\|^2 \\ 1537 & + 4\eta^2 \sum_{j=0}^{t-1} \beta_1^{t-j-1} \left[2L^2 \mathbb{E}_m [\|x_j^m - \bar{x}_j\|^2] + \|\mathbb{E}_m[\widehat{g_j^m} - \nabla f(x_j^m)]\|^2 \right]. \end{aligned} \quad (\text{C.84})$$

$$\begin{aligned} 1541 \|\bar{z}_t - \bar{x}_t\|^2 & \leq \left(\frac{\eta\beta_1}{1-\beta_1}\right)^2 [16L^2\Lambda_{t-1} + 32\|\nabla f(\bar{z}_t)\|^2] \\ 1542 & + \frac{4(\eta\beta_1)^2}{1-\beta_1} \sum_{j=0}^{t-1} \beta_1^{t-j-1} \left[2L^2 \mathbb{E}_m [\|x_j^m - \bar{x}_j\|^2] + \|\mathbb{E}_m[\widehat{g_j^m} - \nabla f(x_j^m)]\|^2 \right]. \end{aligned} \quad (\text{C.85})$$

1548
 1549 *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
 1550

$$\begin{aligned} 1551 \|\bar{x}_t - \bar{x}_{t-1}\|^2 & = \eta^2 \|\bar{u}_{t-1}\|^2 \\ 1552 & = \eta^2 \left\| (1-\beta_1) \sum_{j=0}^{t-1} \beta_1^{t-j-1} \mathbb{E}_m[\widehat{g_j^m}] \right\|^2 \\ 1553 & \leq 2\eta^2 \left[\left\| (1-\beta_1) \sum_{j=0}^{t-1} \beta_1^{t-j-1} \mathbb{E}_m[\nabla f(x_j^m)] \right\|^2 + \left\| (1-\beta_1) \sum_{j=0}^{t-1} \beta_1^{t-j-1} \mathbb{E}_m[\widehat{g_j^m} - \nabla f(x_j^m)] \right\|^2 \right] \\ 1554 & \leq 4\eta^2 \left\| (1-\beta_1) \sum_{j=0}^{t-1} \beta_1^{t-j-1} \nabla f(\bar{x}_j) \right\|^2 \\ 1555 & + 2\eta^2 (1-\beta_1) \sum_{j=0}^{t-1} \beta_1^{t-j-1} \left[2L^2 \mathbb{E}_m [\|x_j^m - \bar{x}_j\|^2] + \|\mathbb{E}_m[\widehat{g_j^m} - \nabla f(x_j^m)]\|^2 \right]. \end{aligned} \quad (\text{C.86})$$

1566 Note that

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

1584 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

1586 $\Lambda_t = \sum_{j=0}^{t-1} a_{t,j} \|\bar{x}_j - \bar{x}_{j+1}\|^2$, we can conclude that

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

1594 which implies (C.85). We complete the proof by plugging the above inequality in

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

□

C.5 FURTHER DISCUSSION

1601 **Coordinate-wise clipping and global clipping.** Lemma B.2 can be easily extended to \mathbb{R}^d , similar
1602 to Sadiev et al. (2023, Lemma 5.1). Therefore, our results can be easily generalized to global
1603 clipping operator $\text{clip}_g(X, \rho_g) := \min \left\{ 1, \frac{\rho_g}{\|X\|} \right\} X$ with threshold $\rho_g := \rho\sqrt{d}$. We omit the
1604 details in this paper. Readers may also wonder why our Theorem C.4 and Theorem C.5 depend on
1605 $\text{poly}(d)$. 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
1606 our convergence guarantee will not depend on $\text{poly}(d)$ explicitly. Zhang et al. (2020, Corollary
1607 7) claims that coordinate-wise clipping has better dependence on dimension d . But they simply
1608 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
1609 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 $\mathcal{O}(d^{\frac{\alpha}{2}})$,
1610 then there is still no difference between coordinate-wise clipping and global clipping in their setting.
1611

1613 **Prior works on distributed SGDM with local updates.** There are many works on *Local* SGDM
1614 in distributed setting. Liu et al. (2020a) studies *Local* SGDM in convex setting and rely on some
1615 strong assumptions to show convergence. Xu et al. (2021) analyze *Local* SGDM with bounded
1616 gradient assumption and the use a global momentum parameter during local iterations. Yu et al.
1617 (2019) considers non-convex *Local* SGDM but is only able to prove linear speedup. Wang et al.
1618 (2019); Cheng et al. (2023) also study non-convex problem and use momentum to handle hetero-
1619 geneity in federated learning. All these works fail to show the benefits of local iterations compared
to minibatch baseline.

1620 **D PROOF OF LOCAL ADAM**
 1621

1622 **D.1 OVERVIEW AND MAIN THEOREM**
 1623

1624 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 $k(t) \leq K-1$.
 1625 We omit the dependence on t and let $r = r(t), k = k(t)$ through out the proof if not causing
 1626 confusion. 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
 1627 to the following update rule:
 1628

$$1629 \quad u_t^m = \begin{cases} \beta_1 u_{t-1}^m + (1 - \beta_1) \widehat{g}_t^m & \text{if } t \bmod K \not\equiv 0, \\ 1630 \quad \beta_1 \bar{u}_{t-1} + (1 - \beta_1) \widehat{g}_t^m & \text{otherwise,} \end{cases} \quad (\text{D.1})$$

$$1632 \quad v_t^m = \begin{cases} \beta_2 v_{t-1}^m + (1 - \beta_2) \widehat{g}_t^m & \text{if } t \bmod K \not\equiv 0, \\ 1633 \quad \beta_2 \bar{v}_{t-1} + (1 - \beta_2) \widehat{g}_t^m & \text{otherwise,} \end{cases} \quad (\text{D.2})$$

$$1635 \quad x_{t+1}^m = \begin{cases} x_t^m - \eta (H_t^m)^{-1} u_t^m & \text{if } t \bmod K \not\equiv -1, \\ 1636 \quad \bar{x}_t - \eta \mathbb{E}_m[(H_t^m)^{-1} u_t^m] & \text{otherwise.} \end{cases} \quad (\text{D.3})$$

1637 Define an auxiliary sequence $\{z_t^m\}$ as:
 1638

$$1639 \quad z_{t+1}^m = \begin{cases} \frac{1}{1 - \beta_1} x_{t+1}^m - \frac{\beta_1}{1 - \beta_1} x_t^m & \text{if } t \bmod K \not\equiv -1, \\ 1640 \quad \frac{1}{1 - \beta_1} x_{t+1}^m - \frac{\beta_1}{1 - \beta_1} \bar{x}_t & \text{otherwise.} \end{cases} \quad (\text{D.4})$$

1643 Let

$$1644 \quad e_t^m := \frac{\beta_1}{1 - \beta_1} (I_d - H_t^m (H_{t-1}^m)^{-1}) u_{t-1}^m. \quad (\text{D.5})$$

1647 Then the definition of $\{z_t^m\}$ implies

$$1648 \quad z_{t+1}^m - z_t^m = -\frac{\eta (H_t^m)^{-1} u_t^m}{1 - \beta_1} + \frac{\eta \beta_1 (H_{t-1}^m)^{-1} u_{t-1}^m}{1 - \beta_1} \\ 1649 \quad = -\frac{\eta \beta_1}{1 - \beta_1} [(H_t^m)^{-1} - (H_{t-1}^m)^{-1}] u_{t-1}^m - \eta (H_t^m)^{-1} \widehat{g}_t^m \\ 1650 \quad =: -\eta (H_t^m)^{-1} (\widehat{g}_t^m + e_t^m). \quad (\text{D.6})$$

1654 Finally, let $y_t := \arg \min_y f(y) + \frac{1}{2\gamma} \|y - \bar{z}_t\|_{H_{r(t)}}^2$.

1657 Define probabilistic events (see (D.15) for definition of some parameters)

$$1659 \quad \mathcal{A}_{t,1} := \left\{ \beta_2^{K/2} \preceq H_{r(t)}^{-1} H_t^m \preceq 1 + (1 - \beta_2) B \text{ and for all } m \in [M] \right\}, \quad (\text{D.7})$$

$$1661 \quad \mathcal{A}_{t,2} := \left\{ \|H_{r(t)}((H_t^m)^{-1} - (H_t^n)^{-1})\| \leq (1 - \beta_2) B_1 \text{ for all } m, n \in [M] \right\}, \quad (\text{D.8})$$

$$1663 \quad \mathcal{A}_{t,3} := \left\{ \|z_{t+1}^m - z_{t+1}^n\|_{H_r}^2 \leq \frac{\eta^2 \sigma^2}{\lambda} K A, \sum_{j=rK}^t \|\widehat{g}_j^m\|^2 \leq \frac{(1 - \beta_1)^2 \sigma^2 A}{2^{12}(1 - \beta_2)^2 B_1^2} \text{ for all } m, n \in [M] \right\}, \quad (\text{D.9})$$

$$1667 \quad \mathcal{A}_{t,4} := \left\{ f_\gamma^{H_{r(t+1)}}(\bar{z}_{t+1}) - \min f_\gamma^\lambda + \frac{\eta}{12} \sum_{j=0}^t \|\nabla f_\gamma^{H_{r(j)}}(\bar{z}_j)\|_{H_{r(j)}^{-1}}^2 \leq 2\Delta \right\}. \quad (\text{D.10})$$

1670 Here $\Delta := f_\gamma^\lambda(x_0) - \min f_\gamma^\lambda$. Besides, let
 1671

$$1672 \quad E_t := \{\mathcal{A}_{j,i} \text{ holds for all } j \leq t-1, i \in \{1, 2, 3, 4\}\}, \quad (\text{D.11})$$

$$1673 \quad E_{t,1} := E_t \cap \mathcal{A}_{t,1}, E_{t,2} := E_{t,1} \cap \mathcal{A}_{t,2}, E_{t,3} := E_{t,2} \cap \mathcal{A}_{t,3}. \quad (\text{D.12})$$

Theorem D.1. For $L/\lambda \geq \gamma^{-1} \geq 2\tau/\lambda$, let Assumption 1, 2, 3, 5 hold for $\Omega = \text{conv}(\mathbf{B}_{R_0}(\Omega_0))$, 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 that for any $x \in \Omega$, $\|\nabla f(x)\| \leq G$, $\|\nabla f(x)\|_\infty \leq G_\infty$, and

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

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

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

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

Here

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

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

When $t = 0$, this is trivial. Assume that the statement is true for some $t \leq T - 1$. We aim to prove 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

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

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

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

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

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

In Lemma D.9,

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

In Lemma D.10,

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

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

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

which can be ensured as long as

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

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

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

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

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

□

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

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

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

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

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

And equivalently,

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

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

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

Hence the optimal ρ is given by

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

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

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

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

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

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

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

This completes the proof.

Theorem D.3 (Complete version of Theorem 3). *Under the conditions of Theorem D.2, let $\gamma = \frac{\lambda}{L}$ 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 than $1 - \delta$, Local Adam with optimal η, ρ yields $\frac{\lambda}{KR} \sum_{r=0}^{R-1} \sum_{k=0}^{K-1} \|\nabla f(\bar{z}_{r,k})\|_{H_r^{-1}}^2 \leq \varepsilon$ if*

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

1836 And equivalently,

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

1845 Further, if $1 - \beta_2 \lesssim \frac{G_\infty^2 + \sigma_\infty^2}{\rho^2 \log \frac{dR}{\delta}}$, where ρ is definded in (D.29), then with probability no less than
1846
1847 $1 - 2\delta$,

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

1856
1857 *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 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
1858 for T in Theorem D.2 will reduce to (D.32). Solve ε and we get the upper bound of
1859

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

1864 Now we turn to bound $\|H_r\|$. Note that $H_{r+1} = \mathbf{diag}(\sqrt{v_{r+1} + \lambda^2})$ and

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

1875 where the last inequality is due to Lemma B.2. Define
1876

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

1880 Further note that
1881

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

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

1890 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

1891
1892
1893 **B.1**, we have $|\sum_{j=0}^{rK-1} [\theta_j]_i| \leq b$ with probability no less than

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

1894 which implies $[H_r]_{i,i} \leq \lambda + 2G_\infty + 2\sigma_\infty$. Therefore, we have

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

1895 And thus

$$\begin{aligned} 1896 \frac{1}{KR} \sum_{r=0}^{R-1} \sum_{k=0}^{K-1} \|\nabla f(\bar{z}_{r,k})\|^2 &\lesssim \left(1 + \frac{G_\infty + \sigma_\infty}{\lambda}\right) \left[\frac{\tau\Delta}{R} + \frac{L\Delta}{KR} + \sqrt{\frac{L\Delta\sigma^2}{MKR}} \log^{\frac{1}{4}} \frac{T}{\delta} + \frac{(L\Delta\sigma)^{\frac{2}{3}}}{K^{\frac{1}{3}}R^{\frac{2}{3}}} \log^{\frac{1}{3}} \frac{MKR}{\delta} \right. \\ 1897 &\quad \left. + \left(\|\boldsymbol{\sigma}\|_{2\alpha} d^{\frac{1}{2} - \frac{1}{2\alpha}}\right)^{\frac{2\alpha}{3\alpha-2}} \left(\frac{L\Delta \log \frac{MKR}{\delta}}{KR}\right)^{\frac{2(\alpha-1)}{3\alpha-2}} \right]. \end{aligned} \quad (\text{D.41})$$

1898
1899 \square

D.2 PRELIMINARIES

1900 We start with theoretical properties of weakly convex function and Moreau envelop, which are re-
1901 peatedly used in our proof.

1902 **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
1903 $L/\lambda \geq \gamma^{-1} \geq 2\tau/\lambda$. Then

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

1904 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

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

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

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

1905 *Proof.* Since y is the minimizer,

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

1906 and note that

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

1907 If $f_\gamma^H(z) - \min f_\gamma^\lambda \leq 2\Delta$, then $f_\gamma^\lambda(z) \leq f_\gamma^H(z)$ and

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

1908 which implies $y, z \in \Omega_0$.

1909 By mean value theorem, there exists a symmetric matrix $-\tau I_d \preceq H_g \preceq L I_d$, such that

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

1944 Hence,

$$1946 \quad \|\nabla f(z) - \nabla f(y)\|_{H^{-1}} \leq \gamma \|H^{-1} \nabla f(y)\|_{H_g H^{-1} H_g} \leq \frac{\gamma L}{\lambda} \|\nabla f_\gamma^H(z)\|_{H^{-1}}. \quad (\text{D.50})$$

$$1948 \quad \|\nabla f(z)\|_{H^{-1}} \leq (1 + \frac{\gamma L}{\lambda}) \|\nabla f_\gamma^H(z)\|_{H^{-1}} \leq \frac{2\gamma L}{\lambda} \|\nabla f_\gamma^H(z)\|_{H^{-1}}. \quad (\text{D.51})$$

1950 Also,

$$1951 \quad H(z - y) - \eta \nabla f(z) = (\gamma I_d - \eta(I_d + \gamma H_g H^{-1})) \nabla f(y) =: \gamma \Lambda \nabla f(y). \quad (\text{D.52})$$

1953 By noticing that

$$1954 \quad -I_d \preceq H^{-1/2} \Lambda H^{1/2} = I_d - \eta \gamma^{-1} - \eta H^{-1/2} H_g H^{-1/2} \preceq I_d, \quad (\text{D.53})$$

1956 we have $\|H(z - y) - \eta \nabla f(z)\|_{H^{-1}} \leq \gamma \|\nabla f(y)\|_{H^{-1}}$.

1957 Last,

$$1959 \quad \min f_\gamma^\lambda \leq f_\gamma^\lambda(y) \leq f(y) = f_\gamma^H(z) - \frac{1}{2\gamma} \|y - z\|_H^2 = f_\gamma^H(z) - \frac{\gamma}{2} \|\nabla f_\gamma^H(z)\|_{H^{-1}}^2. \quad (\text{D.54})$$

1961 This completes the proof. \square

1963 **Lemma D.5.** If $x, y \in \Omega$, then

$$1964 \quad -\langle x - y, \nabla f(x) - \nabla f(y) \rangle + \frac{1}{L} \|\nabla f(x) - \nabla f(y)\|^2 \leq 2\tau \|x - y\|^2. \quad (\text{D.55})$$

1967 *Proof.* By mean value theorem, there exists a symmetric matrix $-\tau I_d \preceq H \preceq L I_d$, such that

$$1968 \quad \nabla f(x) - \nabla f(y) = H(x - y). \quad (\text{D.56})$$

1970 Therefore,

$$1972 \quad -\langle x - y, \nabla f(x) - \nabla f(y) \rangle + \frac{1}{L} \|\nabla f(x) - \nabla f(y)\|^2 = (x - y)^T (-H + \frac{H^2}{L})(x - y) \\ 1973 \quad \leq (\tau + \frac{\tau^2}{L}) \|x - y\|^2 \\ 1975 \quad \leq 2\tau \|x - y\|^2. \quad (\text{D.57})$$

1977 \square

1979 **Lemma D.6.** If $\gamma = \frac{\lambda}{L}$, then for $z \in \Omega_0$, it holds that $\frac{f(z) - f_*}{2} \leq f_{1/L}(z) - f_* \leq f(z) - f_*$.

1981 *Proof.* By definition of Moreau envelop, the second inequality is trivial. Let $y = \arg \min_x f(x) +$
1983 $\frac{L}{2} \|x - z\|^2$. Note that $x \rightarrow f(x) + \frac{L}{2} \|x - z\|^2$ is $2L$ -smooth. Then we have
1984

$$1986 \quad f(z) \leq f(y) + \frac{L}{2} \|y - z\|^2 + L \|y - z\|^2 = f_{1/L}(z) + L \|y - z\|^2. \quad (\text{D.58})$$

1988 Furthermore, by Lemma D.4

$$1989 \quad \frac{L}{2} \|y - z\|^2 = \frac{1}{2L} \|\nabla f(y)\|^2 \leq f(y) - f_*. \quad (\text{D.59})$$

1992 Therefore, $f(z) - f_* \leq f_{1/L}(z) - f_* + L \|y - z\|^2 \leq 2(f_{1/L}(z) - f_*)$. \square

1994 Next, we show that event E_t implies all the iterates remain in certain area.

1995 **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
1996 $\bar{z}_j \in \Omega_0, x_j^m, \bar{x}_j, z_j^m \in \Omega$. And $\|x_j^m - x_j^n\| \leq \frac{\eta\sigma}{\lambda} \sqrt{KA}$ for all m, n .

1998 *Proof.* Event E_t implies that for all $j \leq t$,

$$2000 f_\gamma^\lambda(\bar{z}_j) - \min f_\gamma^\lambda \leq 2\Delta, \|z_j^m - z_j^n\| \leq \frac{\eta\sigma}{\lambda} \sqrt{KA} \leq \sqrt{\frac{\Delta\gamma}{160\lambda}}. \quad (\text{D.60})$$

2002 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 $\bar{x}_j \in$
 2003 $\mathbf{conv}\{\bar{z}_i\}_{i \leq j} \subset \mathbf{conv}(\Omega_0) \subset \Omega$ and $x_j^m - x_j^n \in \mathbf{conv}\{z_i^m - z_i^n\}_{i \leq j}$. We have
 2004

$$2006 \|x_j^m - x_j^n\| \leq \frac{\eta\sigma}{\lambda} \sqrt{KA}, \|x_j^m - \bar{x}_j\| \leq \frac{\eta\sigma}{\lambda} \sqrt{KA} \leq \sqrt{\frac{\Delta\gamma}{160\lambda}}. \quad (\text{D.61})$$

2008 Therefore by Lemma B.4, $x_j^m \in \mathbf{B}_{R_0}(\mathbf{conv}(\Omega_0)) = \Omega$. \square

2010 The following lemma shows that the second order momentum v_t^m does not change too much from
 2011 $v_{r(t)}$ during local training with high probability, which is also repeatedly used in our proof.

2012 **Lemma D.8.** Let $B := \max \left\{ \frac{6K(G_\infty^2 + \sigma_\infty^2)}{\lambda^2}, \frac{16\rho^2}{\lambda^2} \log \frac{dMT}{\delta}, 2^6 \frac{\sqrt{K}(G_\infty + \sigma_\infty)\sigma_\infty}{\lambda^2} \log^{1/2} \frac{dMT}{\delta} \right\}$.

2015 If $\rho \geq \max\{3\sigma_\infty, 2G_\infty\}$, then the following holds

$$2016 \mathbb{P}(E_{t,1}) \geq \mathbb{P}(E_t) - \frac{\delta}{4T}. \quad (\text{D.62})$$

2019 *Proof.* Let $t = rK + k$. By the update rule of local Adam, we have

$$2021 v_t^m = \beta_2^{k+1} v_r + (1 - \beta_2) \sum_{j=rK}^t \beta_2^{t-j} \widehat{g}_j^m \odot \widehat{g}_j^m \succeq \beta_2^K v_r, \quad (\text{D.63})$$

2024 and hence

$$2025 H_t^m = \mathbf{diag}(\sqrt{v_t^m + \lambda^2}) \succeq \beta_2^{K/2} \mathbf{diag}(\sqrt{v_r + \lambda^2}) = \beta_2^{K/2} H_r. \quad (\text{D.64})$$

2026 For the upper bound, for any index $i \in [d]$, by Lemma B.2,

$$2027 \mathbb{E}_j[\widehat{g}_j^m]^2 \leq \sigma_i^2 + [\mathbb{E}_j[\widehat{g}_j^m]]^2 \leq \sigma_\infty^2 + 3G_\infty^2. \quad (\text{D.65})$$

2029 Therefore,

$$2031 [v_t^m]_i \leq [v_r]_i + (1 - \beta_2) K(\sigma_\infty^2 + 3G_\infty^2) + (1 - \beta_2) \sum_{j=rK}^t \left[[\widehat{g}_j^m]^2 - \mathbb{E}_j[\widehat{g}_j^m]^2 \right]. \quad (\text{D.66})$$

2033 Define

$$2034 [\theta_j^m]_i = \begin{cases} [\widehat{g}_j^m]^2 - \mathbb{E}_j[\widehat{g}_j^m]^2, & \text{if event } E_j \text{ holds,} \\ 0, & \text{otherwise.} \end{cases} \quad (\text{D.67})$$

2037 Event E_t implies $[\theta_j^m]_i = [\widehat{g}_j^m]^2 - \mathbb{E}_j[\widehat{g}_j^m]^2$. Further note that $|[\theta_j^m]_i| \leq \rho^2 \stackrel{\text{def}}{=} c$,

$$\begin{aligned} 2039 \text{Var}_j([\theta_j^m]_i) &\leq \mathbb{E}_j \left[[\widehat{g}_j^m]^2 - [\nabla f(x_j^m)]_i^2 \right]^2 \\ 2040 &= \mathbb{E}_j \left[[\widehat{g}_j^m]_i - [\nabla f(x_j^m)]_i \right]^2 \left[[\widehat{g}_j^m]_i - [\nabla f(x_j^m)]_i + 2[\nabla f(x_j^m)]_i \right]^2 \\ 2041 &\stackrel{\text{AM-GM}}{\leq} 2\mathbb{E}_j \left[[\widehat{g}_j^m]_i - [\nabla f(x_j^m)]_i \right]^4 + 8\mathbb{E}_j \left[[\widehat{g}_j^m]_i - [\nabla f(x_j^m)]_i \right]^2 [\nabla f(x_j^m)]_i^2 \\ 2043 &\stackrel{\text{Lemma B.2}}{\leq} 2\sigma_\infty^4 + 8\sigma_\infty^2 G_\infty^2. \end{aligned} \quad (\text{D.68})$$

2047 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
 2048 probability no less than
 2049

$$2051 1 - 2 \exp \left(-\frac{b^2}{2V + 2cb/3} \right) \geq 1 - \frac{\delta}{4dMT}, \quad (\text{D.69})$$

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

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

and thus

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

□

D.3 PROOF OF CONTRACTION

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

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

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

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:

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

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

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

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

Define

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

Then we have

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

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

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 $\left| \sum_{j=rK}^t [\Xi_j^{m,n}]_i \right| \leq b$ with probability no less than

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

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

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

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

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

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

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

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

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

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

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

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

then the following holds:

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

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 that $t \bmod K \not\equiv -1$. The update rules implies

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

2160 Note that the first order term is
 2161

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

2169 And for the first term above,
 2170

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

2180 By definition of $\{z_t^m\}$ and event $E_{t,2}$,
 2181

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

2190 Besides,
 2191

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

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

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

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

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

Here we repeatedly apply $\|H_r(H_t^n)^{-1} - I_d\| \leq (1 - \beta_2)B$ and $\|H_r((H_t^m)^{-1} - (H_t^n)^{-1})\| \leq (1 - \beta_2)B_1$ 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{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{(D.84)}{=} 2\eta \langle x_t^m - x_t^n, \nabla f(x_t^m) - \nabla f(x_t^n) \rangle \\
& \quad + 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 \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{(D.87)}{=} 2\eta K \cdot (*) + 2\eta K \cdot (***) \right] \\
& \quad + 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 (1 + \frac{1}{2K}) \|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 + (\frac{\eta}{2L} + \frac{4\eta^2}{\lambda}) \|\nabla f(x_t^m) - \nabla f(x_t^n)\|^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}_{(\sharp)} \\
& \quad + 8\eta^2 K ((*) + (**)) \\
& \leq (1 + \frac{1}{2K}) \|z_t^m - z_t^n\|_{H_r}^2 \underbrace{- 2\eta \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)\|^2}_{(\sharp\sharp)} \\
& \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 (\sharp) \\
& \quad + \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] \\
& \quad + 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{(D.90)}{=} 16\eta^2 K (1 - \beta_2)^2 B_1^2 \|\widehat{g}_t^m\|_{H_r^{-1}}^2 \\
& \leq (1 + \frac{1}{2K}) \|z_t^m - z_t^n\|_{H_r}^2 + (\sharp\sharp) + 8\eta^2 \cdot (\sharp) \\
& \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 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 \\
& \quad \underbrace{}_{(\sharp\sharp\sharp)} \\
& \leq (1 + \frac{1}{K}) \|z_t^m - z_t^n\|_{H_r}^2 + (\sharp\sharp) + 8\eta^2 \cdot (\sharp) - 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 + (\sharp\sharp\sharp) \stackrel{\text{Lemma B.2}}{=} \frac{8\eta^2 K}{\lambda} \cdot \frac{\|2\sigma\|_{2\alpha}^{2\alpha}}{\rho^{2(\alpha-1)}}. \tag{D.91}
\end{aligned}$$

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 notice that by definition of $\{u_t^m\}$,

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

which implies

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

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

And thus

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

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

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

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,

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

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

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

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

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

Define

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

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

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

by Lemma B.2,

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

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

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

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

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

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

Also note that

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

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

Let $b = \frac{\eta^2\sigma^2}{4\lambda} KA$, $V = \frac{72\eta^4\sigma^4}{\lambda^2} K^2 A$. Then by Lemma B.1, $|\sum_{j=rK}^t \zeta_j^{m,n}| \leq b$ with probability no less than $1 - 2\exp\left(-\frac{b^2}{2V + 2cb/3}\right)$.

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

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

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

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

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

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),

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

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

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

□

D.4 PROOF OF DESCENT LEMMA

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

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

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

and

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

2430 Then the following holds:
 2431

$$2432 \quad \mathbb{P}(E_{t+1}) \geq \mathbb{P}(E_{t,3}) - \frac{\delta}{4T}. \quad (\text{D.115})$$

2436
 2437 *Proof.* For any $x \in \mathbb{R}^d$, since $\nabla^2 f(\cdot) \succeq -\tau I_d$ and $H_r \succeq \lambda I_d$, $y \mapsto f(y) + \frac{1}{2\gamma} \|x - y\|_{H_r}^2$ is $(\frac{1}{\gamma} - \frac{\tau}{\lambda})$ -
 2438 convex with respect to $\|\cdot\|_{H_r}$. Note that under event E_t , $\bar{z}_t \in \Omega_0$. Let $y_t := \arg \min_y f(y) + \frac{1}{2\gamma} \|\bar{z}_t - y\|_{H_r}^2$, and by Lemma D.4, $y_t \in \Omega_0$. Then
 2440
 2441

$$2442 \quad f(y_t) + \frac{1}{2\gamma} \|y_t - \bar{z}_t\|_{H_r}^2 \leq f(\bar{z}_{t+1}) + \frac{1}{2\gamma} \|\bar{z}_{t+1} - \bar{z}_t\|_{H_r}^2 - \frac{1}{2} \left(\frac{1}{\gamma} - \frac{\tau}{\lambda} \right) \|\bar{z}_{t+1} - y_t\|_{H_r}^2. \quad (\text{D.116})$$

2445 Recall that the definition of $\{z_t^m\}$ implies
 2446

$$\begin{aligned} 2447 \quad z_{t+1}^m - z_t^m &= -\frac{\eta(H_t^m)^{-1} u_t^m}{1 - \beta_1} + \frac{\eta\beta_1(H_{t-1}^m)^{-1} u_{t-1}^m}{1 - \beta_1} \\ 2448 \quad &= -\frac{\eta\beta_1}{1 - \beta_1} [(H_t^m)^{-1} - (H_{t-1}^m)^{-1}] u_{t-1}^m - \eta(H_t^m)^{-1} \widehat{g}_t^m \\ 2449 \quad &= -\eta(H_t^m)^{-1} (\widehat{g}_t^m + e_t^m). \end{aligned} \quad (\text{D.117})$$

2454 Here $e_t^m = \frac{\beta_1}{1 - \beta_1} (I_d - H_t^m (H_{t-1}^m)^{-1}) u_{t-1}^m$.

2455 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
 2456

$$\begin{aligned} 2457 \quad f(\bar{z}_{t+1}) - f(y_t) &\leq f(\bar{z}_t) + \langle \nabla f(\bar{z}_t), \bar{z}_{t+1} - \bar{z}_t \rangle + \frac{L}{2} \|\bar{z}_{t+1} - \bar{z}_t\|^2 - f(y_t) \\ 2458 \quad &\leq \langle \nabla f(\bar{z}_t), \bar{z}_{t+1} - y_t \rangle + \frac{\tau}{2} \|\bar{z}_t - y_t\|^2 + \frac{L}{2} \|\bar{z}_{t+1} - \bar{z}_t\|^2 \\ 2459 \quad &\leq \langle \nabla f(\bar{z}_t), \bar{z}_{t+1} - y_t \rangle + \frac{\tau}{2\lambda} \|\bar{z}_t - y_t\|_{H_r}^2 + \frac{L}{2\lambda} \|\bar{z}_{t+1} - \bar{z}_t\|_{H_r}^2. \end{aligned} \quad (\text{D.118})$$

2460 Combine this with (D.116),
 2461

$$\begin{aligned} 2462 \quad &\frac{\frac{1}{\eta} + \frac{1}{\gamma} - \frac{\tau}{\lambda}}{2} \|\bar{z}_{t+1} - y_t\|_{H_r}^2 - \frac{\frac{1}{\eta} - \frac{1}{\gamma} + \frac{\tau}{\lambda}}{2} \|\bar{z}_t - y_t\|_{H_r}^2 + \frac{\frac{1}{\eta} + \frac{1}{\gamma} - \frac{L}{\lambda}}{2} \|\bar{z}_{t+1} - \bar{z}_t\|_{H_r}^2 \\ 2463 \quad &\leq \left\langle \bar{z}_{t+1} - y_t, \nabla f(\bar{z}_t) + \frac{H_r(\bar{z}_{t+1} - \bar{z}_t)}{\eta} \right\rangle \\ 2464 \quad &= \left\langle \bar{z}_t - \eta \mathbb{E}_m[(H_t^m)^{-1}(\widehat{g}_t^m + e_t^m)] - y_t, \nabla f(\bar{z}_t) - H_r \mathbb{E}_m[(H_t^m)^{-1}(\widehat{g}_t^m + e_t^m)] \right\rangle \\ 2465 \quad &= \left\langle \bar{z}_t - \eta H_r^{-1} \nabla f(\bar{z}_t) - y_t, \nabla f(\bar{z}_t) - H_r \mathbb{E}_m[(H_t^m)^{-1}(\widehat{g}_t^m + e_t^m)] \right\rangle \\ 2466 \quad &\quad + \eta \|\nabla f(\bar{z}_t) - H_r \mathbb{E}_m[(H_t^m)^{-1}(\widehat{g}_t^m + e_t^m)]\|_{H_r^{-1}}^2 \\ 2467 \quad &\leq \left\langle \bar{z}_t - \eta H_r^{-1} \nabla f(\bar{z}_t) - y_t, \nabla f(\bar{z}_t) - H_r \mathbb{E}_m[(H_t^m)^{-1}(\widehat{g}_t^m + e_t^m)] \right\rangle \\ 2468 \quad &\quad + 4\eta \|\nabla f(\bar{z}_t) - \mathbb{E}_m[\nabla f(x_t^m)]\|_{H_r^{-1}}^2 + 4\eta \|\mathbb{E}_m[\nabla f(x_t^m) - \widehat{g}_t^m]\|_{H_r^{-1}}^2 \\ 2469 \quad &\quad + 4\eta \left\| \mathbb{E}_m[(H_r(H_t^m)^{-1} - I_d)\widehat{g}_t^m] \right\|_{H_r^{-1}}^2 + 4\eta \|\mathbb{E}_m[(H_t^m)^{-1}e_t^m]\|_{H_r}^2. \end{aligned} \quad (\text{D.119})$$

2484 By Lemma D.4, we have
 2485

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

2499 Also,
 2500

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

2505 Further noticing that $\eta \leq \frac{\gamma}{4}$ and by AM-GM inequality, we conclude that
 2506

2508 LHS of (D.119)

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

2516 If $t \bmod K \equiv -1$, then $r(t+1) = r(t) + 1 = r + 1$ and event $E_{t,1}$ implies
 2517

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

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

2534 On the other hand, if $t \bmod K \not\equiv -1$, then $r(t+1) = r(t) = r$,
 2535

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

2538 Hence the following always holds:
 2539

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

2561 Sum over t and we get
 2562

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

2574 By AM-GM inequality and notice that $\bar{x}_t, \bar{z}_t \in \Omega$,

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

2579 Under event $E_{t,3}$,

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

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

2586 By the definition of u_{t-1}^m , we have

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

2592 Plug these inequalities above in (D.127),
 2593

$$\begin{aligned}
 2594 \quad & f_{\gamma}^{H_r(t+1)}(\bar{z}_{t+1}) \leq f_{\gamma}^{\lambda}(x_0) - \frac{\eta}{8} \sum_{j=0}^t \|\nabla f(y_j)\|_{H_{r(j)}^{-1}}^2 + \frac{5\eta^2}{\lambda\gamma} \sum_{j=0}^t \|\mathbb{E}_m[\nabla f(x_j^m) - \widehat{g}_j^m]\|^2 \\
 2595 \quad & + \frac{20\eta}{\lambda} \sum_{j=0}^t [L^2 \|\bar{z}_j - \bar{x}_j\|^2 + \|\nabla f(\bar{x}_j) - \mathbb{E}_m[\nabla f(x_j^m)]\|^2] \\
 2596 \quad & + \stackrel{(D.128)}{\eta} \left(\frac{48\beta_1^2}{(1-\beta_1)(\sqrt{\beta_2}-\beta_1)} + 10 \right) (1-\beta_2)^2 B^2 \sum_{j=0}^t \mathbb{E}_m \left[\|\widehat{g}_j^m\|_{H_{r(j)}^{-1}}^2 \right] + (*). \\
 2597 \quad & \\
 2598 \quad & \\
 2599 \quad & \\
 2600 \quad & \\
 2601 \quad & \\
 2602 \quad & \\
 2603 \quad & \\
 2604 \quad & \\
 2605 \quad & \\
 2606 \quad & \\
 2607 \quad & \\
 2608 \quad & \\
 2609 \quad & \\
 2610 \quad & \\
 2611 \quad & \\
 2612 \quad & \\
 2613 \quad & \\
 2614 \quad & \\
 2615 \quad & \\
 2616 \quad & \\
 2617 \quad & \\
 2618 \quad & \\
 2619 \quad & \\
 2620 \quad & \\
 2621 \quad & \\
 2622 \quad & \\
 2623 \quad & \\
 2624 \quad & \\
 2625 \quad & \\
 2626 \quad & \\
 2627 \quad & \\
 2628 \quad & \\
 2629 \quad & \\
 2630 \quad & \\
 2631 \quad & \\
 2632 \quad & \\
 2633 \quad & \\
 2634 \quad & \\
 2635 \quad & \\
 2636 \quad & \\
 2637 \quad & \\
 2638 \quad & \\
 2639 \quad & \\
 2640 \quad & \\
 2641 \quad & \\
 2642 \quad & \\
 2643 \quad & \\
 2644 \quad & \\
 2645 \quad &
 \end{aligned} \tag{D.132}$$

By AM-GM inequality and Lemma D.4,

$$\begin{aligned}
 2605 \quad & \mathbb{E}_m \left[\|\widehat{g}_t^m\|_{H_r^{-1}}^2 \right] \leq 4\mathbb{E}_m \left[\|\widehat{g}_t^m - \nabla f(x_t^m)\|_{H_r^{-1}}^2 + \|\nabla f(x_t^m) - \nabla f(\bar{x}_t)\|_{H_r^{-1}}^2 \right. \\
 2606 \quad & \quad \left. + \|\nabla f(\bar{x}_t) - \nabla f(\bar{z}_t)\|_{H_r^{-1}}^2 + \|\nabla f(\bar{z}_t)\|_{H_r^{-1}}^2 \right] \\
 2607 \quad & \leq \frac{4}{\lambda} \left[\mathbb{E}_m \|\widehat{g}_t^m - \nabla f(x_t^m)\|^2 + L^2 \mathbb{E}_m [\|x_t^m - \bar{x}_t\|^2] + L^2 \|\bar{z}_t - \bar{x}_t\|^2 \right] + \frac{16(\gamma L)^2}{\lambda^2} \|\nabla f_{\gamma}^{H_r}(\bar{z}_t)\|_{H_r^{-1}}^2. \\
 2608 \quad & \\
 2609 \quad & \\
 2610 \quad & \\
 2611 \quad & \\
 2612 \quad & \\
 2613 \quad & \\
 2614 \quad & \\
 2615 \quad & \\
 2616 \quad & \\
 2617 \quad & \\
 2618 \quad & \\
 2619 \quad & \\
 2620 \quad & \\
 2621 \quad & \\
 2622 \quad & \\
 2623 \quad & \\
 2624 \quad & \\
 2625 \quad & \\
 2626 \quad & \\
 2627 \quad & \\
 2628 \quad & \\
 2629 \quad & \\
 2630 \quad & \\
 2631 \quad & \\
 2632 \quad & \\
 2633 \quad & \\
 2634 \quad & \\
 2635 \quad & \\
 2636 \quad & \\
 2637 \quad & \\
 2638 \quad & \\
 2639 \quad & \\
 2640 \quad & \\
 2641 \quad & \\
 2642 \quad & \\
 2643 \quad & \\
 2644 \quad & \\
 2645 \quad &
 \end{aligned} \tag{D.133}$$

Therefore, we achieve that

$$\begin{aligned}
 2613 \quad & f_{\gamma}^{H_r(t+1)}(\bar{z}_{t+1}) \leq f_{\gamma}^{H_0}(x_0) - \frac{\eta}{9} \sum_{j=0}^t \|\nabla f(y_j)\|_{H_{r(j)}^{-1}}^2 + \frac{5\eta^2}{\lambda\gamma} \sum_{j=0}^t \|\mathbb{E}_m[\nabla f(x_j^m) - \widehat{g}_j^m]\|^2 \\
 2614 \quad & + \frac{40\eta}{\lambda} \sum_{j=0}^t [L^2 \|\bar{z}_j - \bar{x}_j\|^2 + \|\nabla f(\bar{x}_j) - \mathbb{E}_m[\nabla f(x_j^m)]\|^2] \\
 2615 \quad & + \frac{160\eta(1-\beta_2)^2 B^2}{\lambda(1-\beta_1)(\sqrt{\beta_2}-\beta_1)} \sum_{j=0}^t \left[\mathbb{E}_m \|\widehat{g}_j^m - \nabla f(x_j^m)\|^2 + L^2 \mathbb{E}_m [\|x_j^m - \bar{x}_j\|^2] \right] + (*). \\
 2616 \quad & \\
 2617 \quad & \\
 2618 \quad & \\
 2619 \quad & \\
 2620 \quad & \\
 2621 \quad & \\
 2622 \quad & \\
 2623 \quad & \\
 2624 \quad & \\
 2625 \quad & \\
 2626 \quad & \\
 2627 \quad & \\
 2628 \quad & \\
 2629 \quad & \\
 2630 \quad & \\
 2631 \quad & \\
 2632 \quad & \\
 2633 \quad & \\
 2634 \quad & \\
 2635 \quad & \\
 2636 \quad & \\
 2637 \quad & \\
 2638 \quad & \\
 2639 \quad & \\
 2640 \quad & \\
 2641 \quad & \\
 2642 \quad & \\
 2643 \quad & \\
 2644 \quad & \\
 2645 \quad &
 \end{aligned} \tag{D.134}$$

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

$$\begin{aligned}
 2624 \quad & \|\bar{z}_j - \bar{x}_j\|^2 \leq \left(\frac{\beta_1}{1-\beta_1} \right)^2 \left[64\eta^2 \left(\|\nabla f(\bar{z}_j)\|_{H_{r(j)}^{-2}}^2 + \frac{L^2}{\lambda^2} \Lambda_{j-1} \right) \right. \\
 2625 \quad & \quad \left. + \frac{36\eta^2}{\lambda^2} (1-\beta_1) \sum_{i=r(j)K}^{j-1} \beta_1^{j-i-1} \left[\frac{\eta^2 L^2 \sigma^2}{\lambda^2} K A + \mathbb{E}_m \|\widehat{g}_i^m - \nabla f(x_i^m)\|^2 \right] \right]. \\
 2626 \quad & \\
 2627 \quad & \\
 2628 \quad & \\
 2629 \quad & \\
 2630 \quad & \\
 2631 \quad & \\
 2632 \quad & \\
 2633 \quad & \\
 2634 \quad & \\
 2635 \quad & \\
 2636 \quad & \\
 2637 \quad & \\
 2638 \quad & \\
 2639 \quad & \\
 2640 \quad & \\
 2641 \quad & \\
 2642 \quad & \\
 2643 \quad & \\
 2644 \quad & \\
 2645 \quad &
 \end{aligned} \tag{D.135}$$

Hence

$$\begin{aligned}
 2632 \quad & \sum_{j=0}^t \|\bar{z}_j - \bar{x}_j\|^2 \leq \left(\frac{\beta_1}{1-\beta_1} \right)^2 \left[64\eta^2 \sum_{j=0}^t \left(\|\nabla f(\bar{z}_j)\|_{H_{r(j)}^{-2}}^2 + \frac{L^2}{\lambda^2} \Lambda_{j-1} \right) \right. \\
 2633 \quad & \quad \left. + \frac{36\eta^2}{\lambda^2} \sum_{j=0}^{t-1} \left[\frac{\eta^2 L^2 \sigma^2}{\lambda^2} K A + \mathbb{E}_m \|\widehat{g}_j^m - \nabla f(x_j^m)\|^2 \right] \right]. \\
 2634 \quad & \\
 2635 \quad & \\
 2636 \quad & \\
 2637 \quad & \\
 2638 \quad & \\
 2639 \quad & \\
 2640 \quad & \\
 2641 \quad & \\
 2642 \quad & \\
 2643 \quad & \\
 2644 \quad & \\
 2645 \quad &
 \end{aligned} \tag{D.136}$$

Additionally by Lemma D.12,

$$\begin{aligned}
 2640 \quad & \Lambda_t + \frac{(1-\beta_1)^2}{2} \sum_{j=0}^{t-1} \Lambda_j \leq \frac{64\eta^2}{1-\beta_1} \sum_{j=0}^t \|\nabla f(\bar{z}_j)\|_{H_{r(j)}^{-2}}^2 \\
 2641 \quad & + \frac{36\eta^2}{\lambda^2} (1-\beta_1) \sum_{j=0}^{t-1} \left[\frac{\eta^2 L^2 \sigma^2}{\lambda^2} K A + \mathbb{E}_m \|\widehat{g}_j^m - \nabla f(x_j^m)\|^2 \right]. \\
 2642 \quad & \\
 2643 \quad & \\
 2644 \quad & \\
 2645 \quad &
 \end{aligned} \tag{D.137}$$

2646 Therefore, by noticing that $\Lambda_t \geq 0$ and $\frac{\eta L}{\lambda} \leq \frac{(1 - \beta_1)^2}{16}$,

2647

2648

2649

2650

2651

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

2655

2656 For the third term of RHS of (D.130),

2657

2658

2659

2660

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

2674 For the (*) term of RHS of (D.130),

2675

2676

2677

2678

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

2689 $\stackrel{\text{AM-GM}}{\leq} \frac{2\eta}{\gamma} \sum_{j=0}^t \left[\frac{1}{120\gamma} \|H_{r(j)}(\bar{z}_j - y_j) - \eta \nabla f(\bar{z}_j)\|_{H_{r(j)}^{-1}}^2 + 30\gamma \frac{\|2\sigma\|_{2\alpha}^{2\alpha}}{\lambda\rho^{2(\alpha-1)}} \right] + \textcircled{2}$

2690

2691

2692 $\stackrel{(\text{D.44})}{\leq} \frac{\eta}{60} \sum_{j=0}^t \|\nabla f(y_j)\|_{H_{r(j)}^{-1}}^2 + \frac{60\eta T}{\lambda} \frac{\|2\sigma\|_{2\alpha}^{2\alpha}}{\rho^{2(\alpha-1)}} + \textcircled{2}$

2693

2694

2695

2696

2697

2698

2699 Here we remark that $\textcircled{2}$ is a martingale because $H_{r(j)}$ only depends on stochastic gradients drawn strictly before round $r(j)$ and thus independent of $\widehat{g_j^m}$, which is drawn during round $r(j)$.

2700 Plug (D.138), (D.139), (D.140) in (D.130),
 2701

$$\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(j)}^{-1}}^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} KA + \mathbb{E}_m \|\widehat{g}_j^m - \nabla f(x_j^m)\|^2 \right] + \frac{\eta^2 L^2 \sigma^2}{\lambda^2} KA \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} KA \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(j)}^{-1}}^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} KA \\
 &\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(j)}^{-1}}^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{20(\frac{\eta L}{\lambda})^2}{(1-\beta_1)(\sqrt{\beta_2} - \beta_1)} \sum_{j=0}^t \underbrace{\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(KA + \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(j)}^{-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)}} \\
 &\quad + \textcircled{1} + \textcircled{2} + \textcircled{3}. \tag{D.141}
 \end{aligned}$$

2739 where in the third inequality, we apply $(1-\beta_2)B \leq \frac{\eta L}{\lambda}$.
 2740

2741 For $\textcircled{1}$, define
 2742

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

2746 Then event E_t implies $\textcircled{1} = \sum_{j=0}^t \theta_j$ and notice that
 2747

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

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

2754 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
 2755 than
 2756

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

2760 For ③, define

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

2765 Note that

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

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

2775 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
 2776 probability no less than
 2777

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

2781 For ②, define

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

2786 Then event E_t implies ② = $\sum_{j=0}^t \zeta_j$ and notice that by Lemma D.4,
 2787

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

2797 Therefore,

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

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

2804 Let $b = \Delta/4$, $V = \frac{100\eta\sigma^2\Delta}{\lambda M}$. Then by Lemma B.1,
 2805

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

Note that by Lemma D.4 and event E_t ,

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

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

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

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

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

In the last inequality, we apply

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

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

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 the same conditions in Lemma D.11, event $E_{t,3}$ implies

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

Proof. By the update rule, it always holds that

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

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

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

2862 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
 2863

2864

2865

2866

2867

2868

2869

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

2881

2882

2883

2884

2885

2886

2887

2888

2889

2890

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

2913

2914

2915

(D.163)

2916 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
 2917
 2918

2919 $\Lambda_t = \sum_{j=0}^{t-1} a_{t,j} \|\bar{x}_j - \bar{x}_{j+1}\|^2$, we conclude that
 2920

$$\begin{aligned} 2922 \|\bar{x}_t - \bar{x}_{t-1}\|^2 &\leq 64\eta^2 \left[\|\nabla f(\bar{z}_t)\|_{H_r^{-2}}^2 + \frac{L^2}{\lambda^2} \Lambda_{t-1} \right] + \frac{4\eta^2}{\lambda^2} (1-\beta_1) \sum_{j=rK}^{t-1} \beta_1^{t-j-1} \left[\frac{\eta^2 L^2 \sigma^2}{\lambda^2} KA + \mathbb{E}_m \|\widehat{g}_j^m - \nabla f(x_j^m)\|^2 \right] \\ 2923 &\quad + \frac{32\eta^2(1-\beta_1)}{\lambda^2} \sum_{j=0}^{t-1} \beta_1^{t-j-1} \left[\|\mathbb{E}_m [\widehat{g}_j^m - \nabla f(x_j^m)]\|^2 + \|\mathbb{E}_m [\nabla f(x_j^m) - \nabla f(\bar{x}_j)]\|^2 \right] \\ 2924 &\leq 64\eta^2 \left[\|\nabla f(\bar{z}_t)\|_{H_r^{-2}}^2 + \frac{L^2}{\lambda^2} \Lambda_{t-1} \right] \\ 2925 &\quad + \frac{36\eta^2}{\lambda^2} (1-\beta_1) \sum_{j=rK}^{t-1} \beta_1^{t-j-1} \left[\frac{\eta^2 L^2 \sigma^2}{\lambda^2} KA + \mathbb{E}_m \|\widehat{g}_j^m - \nabla f(x_j^m)\|^2 \right], \\ 2926 & \end{aligned} \tag{D.164}$$

2927 and

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

2928 This completes the proof. \square
 2929

2930 D.5 FURTHER DISCUSSION

2940 **Compared to other results under centralized weakly convex setting.** Theorem D.2 can reduce
 2941 to Minibatch Adam (by substituting M, K with 1 and σ with $\frac{\sigma}{\sqrt{MK}}$ in (D.27) (Petrov, 1992)), and
 2942 the convergence guarantee is
 2943

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

2944 Therefore, in centralized setting with iteration number R and batch size 1, our guarantee for squared
 2945 norm of gradient of Moreau envelope is
 2946

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

2947 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
 2948 term can be dominated by the first term. Then we obtain
 2949

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

2950 In the previous literature of weakly convex function (Davis & Drusvyatskiy, 2019; Alacaoglu et al.,
 2951 2020; Mai & Johansson, 2021), f is typically non-smooth and stochastic gradient is assumed to have
 2952 bounded second order moment. This is weaker than the smoothness assumption but stronger than
 2953 that of noise with bounded moment. There are a few existing results for smooth objective (Davis
 2954 & Drusvyatskiy, 2019; Mai & Johansson, 2020; Deng & Gao, 2021), but they set $\tau = L$. Overall,
 2955 our result is the first convergence guarantee for smooth weakly convex function with $\tau \ll L$ and
 2956 bounded-moment noise.
 2957

2958 **Dependence on β_2 .** The default setting of β_2 in the Adam optimizer of PyTorch is 0.999, which
 2959 is a constant close to 1. Adam with small β_2 has been shown to diverge in some examples (Reddi
 2960 et al., 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 stagnant to provide adaptivity. Therefore, to derive a proper range for β_2 is crucial in the theoretical analysis of Adam.

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

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

Dependence on λ . λ in the denominator of Adam is aimed to avoid numerical instability, and usually a small constant in practice. Note $H_r = \text{diag}(\sqrt{V_r + \lambda^2})$ and v_r is the EMA of squared past gradients. Informally, v_r vanishes as r grows and thus H_r would gradually reduce to λI_d . In the worst case, H_r can be bounded by a constant. In conclusion, the LHS in (4.9) is roughly the averaged squared gradient norm if λ is not too small. It is worth noting that λ can be arbitrarily small or even 0 in (Défossez et al., 2020; Wang et al., 2022; 2024). However, their results all depend on $\text{poly}(d)$. It is still an interesting question to get dimension-free result with small λ .

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

E FAILURE OF STANDARD SGD WITH HEAVY-TAILED NOISE

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

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

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

Proof. We follow the construction of the counter example in Sadiev et al. (2023). To prove the above theorem, we consider a simple 1D-problem $f(x) = Lx^2/2$. It is easy to see that the considered problem is L -strongly convex, L -smooth, and has optimum at $x_* = 0$. We construct the noise in an adversarial way with respect to the parameters of the SGD. Concretely, the noise depends on the

number of iterates t , learning rate η , target precision ε , the starting point x_0 , and the moment bound σ such that

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

where

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

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,

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

which implies that Assumption 3 holds.

We are interested in the situation when

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

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 contradiction and assume that

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

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

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

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

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

Furthermore,

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

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

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

3078
 3079 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, $A = \frac{2\sqrt{\frac{2\varepsilon}{L}}}{\eta\sigma}$
 3080 and we have
 3081

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

3082
 3083 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
 3084

$$3085 \quad T = \Omega\left(\frac{\sigma}{\delta^{1/\alpha}} \sqrt{\frac{L}{\varepsilon}} \log \frac{Lx_0^2}{2\varepsilon}\right). \quad (\text{E.11})$$

3086 This concludes the proof. □
 3087

3088
 3089
 3090
 3091
 3092
 3093
 3094
 3095
 3096
 3097
 3098
 3099
 3100
 3101
 3102
 3103
 3104
 3105
 3106
 3107
 3108
 3109
 3110
 3111
 3112
 3113
 3114
 3115
 3116
 3117
 3118
 3119
 3120
 3121
 3122
 3123
 3124
 3125
 3126
 3127
 3128
 3129
 3130
 3131