

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Technical Preliminaries . . . . .	2
1.2	Our Contributions . . . . .	4
1.3	Closely Related Work . . . . .	5
<b>2</b>	<b>Clipped Stochastic Extragradient</b>	<b>5</b>
<b>3</b>	<b>Clipped Stochastic Gradient Descent-Ascent</b>	<b>7</b>
<b>4</b>	<b>Experiments</b>	<b>7</b>
<b>A</b>	<b>Further Related Work</b>	<b>16</b>
<b>B</b>	<b>Auxiliary Results</b>	<b>17</b>
<b>C</b>	<b>Clipped Stochastic Extragradient: Missing Proofs and Details</b>	<b>18</b>
C.1	Monotone Case . . . . .	18
C.2	Star-Negative Comonotone Case . . . . .	28
C.3	Quasi-Strongly Monotone Case . . . . .	36
<b>D</b>	<b>Clipped Stochastic Gradient Descent-Ascent: Missing Proofs and Details</b>	<b>46</b>
D.1	Monotone Star-Cocoercive Case . . . . .	46
D.2	Star-Cocoercive Case . . . . .	55
D.3	Quasi-Strongly Monotone Star-Cocoercive Case . . . . .	58
<b>E</b>	<b>Extra Experiments</b>	<b>66</b>
E.1	WGAN-GP . . . . .	66
E.2	StyleGAN2 . . . . .	66

## A Further Related Work

**Convergence in expectation.** Convergence in expectation of stochastic methods for solving VIPs is relatively well-studied in the literature. In particular, versions of SEG are studied under bounded variance [Beznosikov et al., 2020, Hsieh et al., 2020], smoothness of stochastic realizations [Mishchenko et al., 2020], and more refined assumptions unifying previously used ones [Gorbunov et al., 2022a]. Recent advances on the in-expectation convergence of SGDA are obtained in [Loizou et al., 2021, Beznosikov et al., 2022].

**Gradient clipping.** In the context of solving minimization problems, gradient clipping [Pascanu et al., 2013] and normalization [Hazan et al., 2015] are known to have a number of favorable properties such as practical robustness to the rapid changes of the loss function [Goodfellow et al., 2016a], provable convergence for structured non-smooth problems with polynomial growth Zhang et al. [2020a], Mai and Johansson [2021] and for the problems with heavy-tailed noise in convex [Nazin et al., 2019, Gorbunov et al., 2020, 2021] and non-convex cases [Zhang et al., 2020b, Cutkosky and Mehta, 2021]. Our work makes a further step towards a better understanding of gradient clipping and is the first to study the theoretical convergence of clipped first-order stochastic methods for VIPs.

**Structured non-monotonicity.** There is a noticeable growing interest of the community in studying the theoretical convergence guarantees of deterministic methods for solving VIP with non-monotone operators  $F(x)$  having a certain structure, e.g., negative comonotonicity [Diakonikolas et al., 2021, Lee and Kim, 2021, Böhm, 2022], quasi-strong monotonicity [Song et al., 2020, Mertikopoulos and Zhou, 2019] and/or star-cocoercivity [Loizou et al., 2021, Gorbunov et al., 2022b,a, Beznosikov et al., 2022]. In the context of stochastic VIPs, SEG (with different extrapolation and update stepsizes) is analyzed under negative comonotonicity by Diakonikolas et al. [2021] and under quasi-strong monotonicity by Gorbunov et al. [2022a], while SGDA is studied under quasi-strong monotonicity and/or star-cocoercivity by [Loizou et al., 2021, Beznosikov et al., 2022]. These results establish in-expectation convergence rates. Our paper continues this line of works and provides the first high-probability analysis of stochastic methods for solving VIPs with structured non-monotonicity.

## B Auxiliary Results

**Useful inequalities.** For all  $a, b \in \mathbb{R}^d$  and  $\alpha > 0$  the following relations hold:

$$2\langle a, b \rangle = \|a\|^2 + \|b\|^2 - \|a - b\|^2, \quad (4)$$

$$\|a + b\|^2 \leq 2\|a\|^2 + 2\|b\|^2, \quad (5)$$

$$-\|a - b\|^2 \leq -\frac{1}{2}\|a\|^2 + \|b\|^2. \quad (6)$$

**Bernstein inequality.** In our proofs, we rely on the following lemma known as *Bernstein inequality for martingale differences* [Bennett, 1962, Dzhaparidze and Van Zanten, 2001, Freedman et al., 1975].

**Lemma B.1.** *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  $b > 0$ ,  $G > 0$  and  $n \geq 1$*

$$\mathbb{P} \left\{ \left| \sum_{i=1}^n X_i \right| > b \text{ and } \sum_{i=1}^n \sigma_i^2 \leq G \right\} \leq 2 \exp \left( -\frac{b^2}{2G + 2cb/3} \right). \quad (7)$$

**Bias and variance of clipped stochastic vector.** We also use the following properties of clipped stochastic estimators from [Gorbunov et al., 2020].

**Lemma B.2** (Simplified version of Lemma F.5 from [Gorbunov et al., 2020]). *Let  $X$  be a random vector in  $\mathbb{R}^d$  and  $\tilde{X} = \text{clip}(X, \lambda)$ . Then,*

$$\left\| \tilde{X} - \mathbb{E}[\tilde{X}] \right\| \leq 2\lambda. \quad (8)$$

Moreover, if for some  $\sigma \geq 0$

$$\mathbb{E}[X] = x \in \mathbb{R}^d, \quad \mathbb{E}[\|X - x\|^2] \leq \sigma^2 \quad (9)$$

and  $x \leq \lambda/2$ , then

$$\left\| \mathbb{E}[\tilde{X}] - x \right\| \leq \frac{4\sigma^2}{\lambda}, \quad (10)$$

$$\mathbb{E} \left[ \left\| \tilde{X} - x \right\|^2 \right] \leq 18\sigma^2, \quad (11)$$

$$\mathbb{E} \left[ \left\| \tilde{X} - \mathbb{E}[\tilde{X}] \right\|^2 \right] \leq 18\sigma^2. \quad (12)$$

*Proof.* The proof of this lemma is identical to the original one, since Gorbunov et al. [2020] rely only on  $\tilde{X} = \text{clip}(X, \lambda)$  to derive (8), and to prove (10)-(12) they use only (9),  $\tilde{X} = \text{clip}(X, \lambda)$  and  $x \leq \lambda/2$   $\square$

## C Clipped Stochastic Extragradient: Missing Proofs and Details

### C.1 Monotone Case

**Lemma C.1.** *Let Assumptions 1.1, 1.2, 1.3 hold for  $Q = B_{4R}(x^*)$ , where  $R \geq R_0 \stackrel{\text{def}}{=} \|x^0 - x^*\|$ , and  $\gamma_1 = \gamma_2 = \gamma$ ,  $0 < \gamma \leq 1/\sqrt{2}L$ . If  $x^k$  and  $\tilde{x}^k$  lie in  $B_{4R}(x^*)$  for all  $k = 0, 1, \dots, K$  for some  $K \geq 0$ , then for all  $u \in B_{4R}(x^*)$  the iterates produced by **clipped-SEG** satisfy*

$$\begin{aligned} \langle F(u), \tilde{x}_{\text{avg}}^K - u \rangle &\leq \frac{\|x^0 - u\|^2 - \|x^{K+1} - u\|^2}{2\gamma(K+1)} + \frac{\gamma}{2(K+1)} \sum_{k=0}^K (\|\theta_k\|^2 + 2\|\omega_k\|^2) \\ &\quad + \frac{1}{K+1} \sum_{k=0}^K \langle x^k - u - \gamma F(\tilde{x}^k), \theta_k \rangle, \end{aligned} \quad (13)$$

$$\tilde{x}_{\text{avg}}^K \stackrel{\text{def}}{=} \frac{1}{K+1} \sum_{k=0}^K \tilde{x}^k, \quad (14)$$

$$\theta_k \stackrel{\text{def}}{=} F(\tilde{x}^k) - \tilde{F}_{\xi_2^k}(\tilde{x}^k), \quad (15)$$

$$\omega_k \stackrel{\text{def}}{=} F(x^k) - \tilde{F}_{\xi_1^k}(x^k). \quad (16)$$

*Proof.* Using the update rule of **clipped-SEG**, for all  $u \in B_{4R}(x^*)$  we obtain

$$\begin{aligned} \|x^{k+1} - u\|^2 &= \|x^k - u\|^2 - 2\gamma \langle x^k - u, \tilde{F}_{\xi_2^k}(\tilde{x}^k) \rangle + \gamma^2 \|\tilde{F}_{\xi_2^k}(\tilde{x}^k)\|^2 \\ &= \|x^k - u\|^2 - 2\gamma \langle x^k - u, F(\tilde{x}^k) \rangle + 2\gamma \langle x^k - u, \theta_k \rangle \\ &\quad + \gamma^2 \|F(\tilde{x}^k)\|^2 - 2\gamma^2 \langle F(\tilde{x}^k), \theta_k \rangle + \gamma^2 \|\theta_k\|^2 \\ &= \|x^k - u\|^2 - 2\gamma \langle \tilde{x}^k - u, F(\tilde{x}^k) \rangle - 2\gamma \langle x^k - \tilde{x}^k, F(\tilde{x}^k) \rangle \\ &\quad + 2\gamma \langle x^k - u - \gamma F(\tilde{x}^k), \theta_k \rangle + \gamma^2 \|F(\tilde{x}^k)\|^2 + \gamma^2 \|\theta_k\|^2 \\ &\stackrel{\text{(Mon)}}{\leq} \|x^k - u\|^2 - 2\gamma \langle \tilde{x}^k - u, F(u) \rangle - 2\gamma^2 \langle \tilde{F}_{\xi_1^k}(x^k), F(\tilde{x}^k) \rangle \\ &\quad + 2\gamma \langle x^k - u - \gamma F(\tilde{x}^k), \theta_k \rangle + \gamma^2 \|F(\tilde{x}^k)\|^2 + \gamma^2 \|\theta_k\|^2 \\ &\stackrel{\text{(4)}}{=} \|x^k - u\|^2 - 2\gamma \langle \tilde{x}^k - u, F(u) \rangle \\ &\quad + \gamma^2 \|\tilde{F}_{\xi_1^k}(x^k) - F(\tilde{x}^k)\|^2 - \gamma^2 \|F(\tilde{x}^k)\|^2 - \gamma^2 \|\tilde{F}_{\xi_1^k}(x^k)\|^2 \\ &\quad + 2\gamma \langle x^k - u - \gamma F(\tilde{x}^k), \theta_k \rangle + \gamma^2 \|F(\tilde{x}^k)\|^2 + \gamma^2 \|\theta_k\|^2 \\ &\stackrel{\text{(5)}}{\leq} \|x^k - u\|^2 - 2\gamma \langle \tilde{x}^k - u, F(u) \rangle \\ &\quad + 2\gamma^2 \|\omega_k\|^2 + 2\gamma^2 \|F(x^k) - F(\tilde{x}^k)\|^2 - \gamma^2 \|\tilde{F}_{\xi_1^k}(x^k)\|^2 \\ &\quad + 2\gamma \langle x^k - u - \gamma F(\tilde{x}^k), \theta_k \rangle + \gamma^2 \|\theta_k\|^2 \\ &\stackrel{\text{(Lip)}}{\leq} \|x^k - u\|^2 - 2\gamma \langle \tilde{x}^k - u, F(u) \rangle - \gamma^2 (1 - 2\gamma^2 L^2) \|\tilde{F}_{\xi_1^k}(x^k)\|^2 \\ &\quad + 2\gamma \langle x^k - u - \gamma F(\tilde{x}^k), \theta_k \rangle + \gamma^2 \|\theta_k\|^2 + 2\gamma^2 \|\omega_k\|^2 \end{aligned}$$

where in the last step we additionally use  $x^k - \tilde{x}^k = \gamma \tilde{F}_{\xi_1^k}(x^k)$  after the application of Lipschitzness of  $F$ . Since  $\gamma \leq 1/\sqrt{2}L$ , we have  $\gamma^2 (1 - 2\gamma^2 L^2) \|\tilde{F}_{\xi_1^k}(x^k)\|^2 \geq 0$ , implying

$$\begin{aligned} 2\gamma \langle F(u), \tilde{x}^k - u \rangle &\leq \|x^k - u\|^2 - \|x^{k+1} - u\|^2 + \gamma^2 (\|\theta_k\|^2 + 2\|\omega_k\|^2) \\ &\quad + 2\gamma \langle x^k - u - \gamma F(\tilde{x}^k), \theta_k \rangle. \end{aligned}$$



Finally, we sum up the above inequality for  $k = 0, 1, \dots, K$  and divide both sides of the result by  $2\gamma(K+1)$ :

$$\begin{aligned}
\langle F(u), \tilde{x}_{\text{avg}}^K - u \rangle &= \frac{1}{K+1} \sum_{k=0}^K \langle F(u), \tilde{x}^k - u \rangle \\
&\leq \frac{1}{2\gamma(K+1)} \sum_{k=0}^K (\|x^k - u\|^2 - \|x^{k+1} - u\|^2) + \frac{\gamma}{2(K+1)} \sum_{k=0}^K \|\theta_k\|^2 \\
&\quad + \frac{1}{K+1} \sum_{k=0}^K \langle x^k - u - \gamma F(\tilde{x}^k), \theta_k \rangle + \frac{\gamma}{K+1} \sum_{k=0}^K \|\omega_k\|^2 \\
&= \frac{\|x^0 - u\|^2 - \|x^{K+1} - u\|^2}{2\gamma(K+1)} + \frac{\gamma}{2(K+1)} \sum_{k=0}^K (\|\theta_k\|^2 + 2\|\omega_k\|^2) \\
&\quad + \frac{1}{K+1} \sum_{k=0}^K \langle x^k - u - \gamma F(\tilde{x}^k), \theta_k \rangle.
\end{aligned}$$

This concludes the proof.  $\square$

**Theorem C.1.** *Let Assumptions 1.1, 1.2, 1.3 hold for  $Q = B_{4R}(x^*)$ , where  $R \geq R_0 \stackrel{\text{def}}{=} \|x^0 - x^*\|$ , and<sup>12</sup>  $\gamma_1 = \gamma_2 = \gamma$ ,*

$$0 < \gamma \leq \frac{1}{160L \ln \frac{6(K+1)}{\beta}}, \quad (17)$$

$$\lambda_{1,k} = \lambda_{2,k} \equiv \lambda = \frac{R}{20\gamma \ln \frac{6(K+1)}{\beta}}, \quad (18)$$

$$m_{1,k} = m_{2,k} \equiv m \geq \max \left\{ 1, \frac{10800(K+1)\gamma^2\sigma^2 \ln \frac{6(K+1)}{\beta}}{R^2} \right\}, \quad (19)$$

for some  $K \geq 0$  and  $\beta \in (0, 1]$  such that  $\ln \frac{6(K+1)}{\beta} \geq 1$ . Then, after  $K$  iterations the iterates produced by [clipped-SEG](#) with probability at least  $1 - \beta$  satisfy

$$\text{Gap}_R(\tilde{x}_{\text{avg}}^K) \leq \frac{9R^2}{2\gamma(K+1)}, \quad (20)$$

where  $\tilde{x}_{\text{avg}}^K$  is defined in (14).

*Proof.* We introduce new notation:  $R_k = \|x^k - x^*\|$  for all  $k \geq 0$ . The proof is based on deriving via induction that  $R_k^2 \leq \tilde{C}R^2$  for some numerical constant  $\tilde{C} > 0$ . In particular, for each  $k = 0, \dots, K+1$  we define probability event  $E_k$  as follows: inequalities

$$\max_{u \in B_R(x^*)} \left\{ \|x^0 - u\|^2 + 2\gamma \sum_{l=0}^{t-1} \langle x^l - u - \gamma F(\tilde{x}^l), \theta_l \rangle + \gamma^2 \sum_{l=0}^{t-1} (\|\theta_l\|^2 + 2\|\omega_l\|^2) \right\} \leq 9R^2, \quad (21)$$

$$\underbrace{\left\| \gamma \sum_{l=0}^{t-1} \theta_l \right\|}_{A_t} \leq R \quad (22)$$

hold for  $t = 0, 1, \dots, k$  simultaneously. Our goal is to prove that  $\mathbb{P}\{E_k\} \geq 1 - k\beta/(K+1)$  for all  $k = 0, 1, \dots, K+1$ . We use the induction to show this statement. For  $k = 0$  the statement

<sup>12</sup>In this and further results, we have relatively large numerical constants in the conditions on step-sizes, batch-sizes, and clipping levels. However, our main goal is deriving results in terms of  $\mathcal{O}(\cdot)$ , where numerical constants are not taken into consideration. Although it is possible to significantly improve the dependence on numerical factors, we do not do it for the sake of proofs' simplicity.

is trivial since  $\|x^0 - u\|^2 \leq 2\|x^0 - x^*\|^2 + 2\|x^* - u\|^2 \leq 4R^2 \leq 9R^2$  and  $\|\gamma \sum_{l=0}^{k-1} \theta_l\| = 0$  for any  $u \in B_R(x^*)$ . Next, assume that the statement holds for  $k = T - 1 \leq K$ , i.e., we have  $\mathbb{P}\{E_{T-1}\} \geq 1 - (T-1)^\beta/(K+1)$ . We need to prove that  $\mathbb{P}\{E_T\} \geq 1 - T^\beta/(K+1)$ . First of all, we show that probability event  $E_{T-1}$  implies  $R_t \leq 3R$  for all  $t = 0, 1, \dots, T$ . For  $t = 0$  we already proved it. Next, assume that we have  $R_t \leq 3R$  for all  $t = 0, 1, \dots, t'$ , where  $t' < T$ . Then, for all  $t = 0, 1, \dots, t'$  we have

$$\begin{aligned} \|\tilde{x}^t - x^*\| &= \|x^t - x^* - \gamma \tilde{F}_{\mathbf{g}_1^t}(x^t)\| \leq \|x^t - x^*\| + \gamma \|\tilde{F}_{\mathbf{g}_1^t}(x^t)\| \\ &\stackrel{(18)}{\leq} \|x^t - x^*\| + \gamma \lambda \leq 3R + \frac{R}{20 \ln \frac{6(K+1)}{\beta}} \leq 4R, \end{aligned} \quad (23)$$

i.e.,  $\tilde{x}^t \in B_{4R}(x^*)$ . This means that the assumptions of Lemma C.1 hold and we have that probability event  $E_{T-1}$  implies

$$\begin{aligned} \max_{u \in B_R(x^*)} \left\{ 2\gamma(t'+1) \langle F(u), \tilde{x}_{\text{avg}}^{t'} - u \rangle + \|x^{t'+1} - u\|^2 \right\} \\ \leq \max_{u \in B_R(x^*)} \left\{ \|x^0 - u\|^2 + 2\gamma \sum_{l=0}^{t'-1} \langle x^l - u - \gamma F(\tilde{x}^l), \theta_l \rangle \right\} \\ \quad + \gamma^2 \sum_{l=0}^{t'-1} (\|\theta_l\|^2 + 2\|\omega_l\|^2) \\ \stackrel{(21)}{\leq} 9R^2, \end{aligned}$$

meaning that

$$\begin{aligned} \|x^{t'+1} - x^*\|^2 &\leq \max_{u \in B_R(x^*)} \left\{ 2\gamma(t'+1) \langle F(u), \tilde{x}_{\text{avg}}^{t'} - u \rangle + \|x^{t'+1} - u\|^2 \right\} \\ &\leq 9R^2, \end{aligned}$$

i.e.,  $R_{t'+1} \leq 3R$ . That is, we proved that probability event  $E_{T-1}$  implies  $R_t \leq 3R$  and

$$\max_{u \in B_R(x^*)} \left\{ 2\gamma(t+1) \langle F(u), \tilde{x}_{\text{avg}}^t - u \rangle + \|x^{t+1} - u\|^2 \right\} \leq 9R^2 \quad (24)$$

for all  $t = 0, 1, \dots, T$ . Moreover, in view of (23)  $E_{T-1}$  also implies that  $\|\tilde{x}^t - x^*\| \leq 4R$  for all  $t = 0, 1, \dots, T$ . Using this, we derive that  $E_{T-1}$  implies

$$\begin{aligned} \|x^t - x^* - \gamma F(\tilde{x}^t)\| &\leq \|x^t - x^*\| + \gamma \|F(\tilde{x}^t)\| \stackrel{(\text{Lip})}{\leq} 3R + \gamma L \|\tilde{x}^t - x^*\| \\ &\stackrel{(23)}{\leq} 3R + 4R\gamma L \stackrel{(17)}{\leq} 5R, \end{aligned} \quad (25)$$

for all  $t = 0, 1, \dots, T$ . Consider random vectors

$$\eta_t = \begin{cases} x^t - x^* - \gamma F(\tilde{x}^t), & \text{if } \|x^t - x^* - \gamma F(\tilde{x}^t)\| \leq 5R, \\ 0, & \text{otherwise,} \end{cases}$$

for all  $t = 0, 1, \dots, T$ . We notice that  $\eta_t$  is bounded with probability 1:

$$\|\eta_t\| \leq 5R \quad (26)$$

for all  $t = 0, 1, \dots, T$ . Moreover, in view of (25), probability event  $E_{T-1}$  implies  $\eta_t = x^t - x^* - \gamma F(\tilde{x}^t)$  for all  $t = 0, 1, \dots, T$ . Therefore,  $E_{T-1}$  implies

$$\begin{aligned}
A_T &= \max_{u \in B_R(x^*)} \left\{ \|x^0 - u\|^2 + 2\gamma \sum_{l=0}^{T-1} \langle x^* - u, \theta_l \rangle \right\} \\
&\quad + 2\gamma \sum_{l=0}^{T-1} \langle x^l - x^* - \gamma F(\tilde{x}^l), \theta_l \rangle + \gamma^2 \sum_{l=0}^{T-1} (\|\theta_l\|^2 + 2\|\omega_l\|^2) \\
&\leq 4R^2 + 2\gamma \max_{u \in B_R(x^*)} \left\{ \left\langle x^* - u, \sum_{l=0}^{T-1} \theta_l \right\rangle \right\} \\
&\quad + 2\gamma \sum_{l=0}^{T-1} \langle \eta_l, \theta_l \rangle + \gamma^2 \sum_{l=0}^{T-1} (\|\theta_l\|^2 + 2\|\omega_l\|^2) \\
&= 4R^2 + 2\gamma R \left\| \sum_{l=0}^{T-1} \theta_l \right\| + 2\gamma \sum_{l=0}^{T-1} \langle \eta_l, \theta_l \rangle + \gamma^2 \sum_{l=0}^{T-1} (\|\theta_l\|^2 + 2\|\omega_l\|^2),
\end{aligned}$$

where  $A_T$  is defined in (21). To continue our derivation we introduce new notation:

$$\theta_l^u \stackrel{\text{def}}{=} \mathbb{E}_{\xi_2^l} [\tilde{F}_{\xi_2^l}(\tilde{x}^l)] - \tilde{F}_{\xi_2^l}(\tilde{x}^l), \quad \theta_l^b \stackrel{\text{def}}{=} F(\tilde{x}^l) - \mathbb{E}_{\xi_2^l} [\tilde{F}_{\xi_2^l}(\tilde{x}^l)], \quad (27)$$

$$\omega_l^u \stackrel{\text{def}}{=} \mathbb{E}_{\xi_1^l} [\tilde{F}_{\xi_1^l}(x^l)] - \tilde{F}_{\xi_1^l}(x^l), \quad \omega_l^b \stackrel{\text{def}}{=} F(x^l) - \mathbb{E}_{\xi_1^l} [\tilde{F}_{\xi_1^l}(x^l)], \quad (28)$$

for all  $l = 0, \dots, T-1$ . By definition we have  $\theta_l = \theta_l^u + \theta_l^b$ ,  $\omega_l = \omega_l^u + \omega_l^b$  for all  $l = 0, \dots, T-1$ . Using the introduced notation, we continue our derivation as follows:  $E_{T-1}$  implies

$$\begin{aligned}
A_T &\stackrel{(5)}{\leq} 4R^2 + 2\gamma R \left\| \sum_{l=0}^{T-1} \theta_l \right\| + \underbrace{2\gamma \sum_{l=0}^{T-1} \langle \eta_l, \theta_l^u \rangle}_{\textcircled{1}} + \underbrace{2\gamma \sum_{l=0}^{T-1} \langle \eta_l, \theta_l^b \rangle}_{\textcircled{2}} \\
&\quad + \underbrace{2\gamma^2 \sum_{l=0}^{T-1} \left( \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2] + 2\mathbb{E}_{\xi_1^l} [\|\omega_l^u\|^2] \right)}_{\textcircled{3}} \\
&\quad + \underbrace{2\gamma^2 \sum_{l=0}^{T-1} \left( \|\theta_l^u\|^2 + 2\|\omega_l^u\|^2 - \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2] - 2\mathbb{E}_{\xi_1^l} [\|\omega_l^u\|^2] \right)}_{\textcircled{4}} \\
&\quad + \underbrace{2\gamma^2 \sum_{l=0}^{T-1} (\|\theta_l^b\|^2 + 2\|\omega_l^b\|^2)}_{\textcircled{5}} \quad (29)
\end{aligned}$$

The rest of the proof is based on deriving good enough upper bounds for  $2\gamma R \left\| \sum_{l=0}^{T-1} \theta_l \right\|$ ,  $\textcircled{1}$ ,  $\textcircled{2}$ ,  $\textcircled{3}$ ,  $\textcircled{4}$ ,  $\textcircled{5}$ , i.e., we want to prove that  $2\gamma R \left\| \sum_{l=0}^{T-1} \theta_l \right\| + \textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4} + \textcircled{5} \leq 5R^2$  with high probability.

Before we move on, we need to derive some useful inequalities for operating with  $\theta_l^u, \theta_l^b, \omega_l^u, \omega_l^b$ . First of all, Lemma B.2 implies that

$$\|\theta_l^u\| \leq 2\lambda, \quad \|\omega_l^u\| \leq 2\lambda \quad (30)$$

for all  $l = 0, 1, \dots, T-1$ . Next, since  $\{\xi_1^{i,l}\}_{i=1}^m, \{\xi_2^{i,l}\}_{i=1}^m$  are independently sampled from  $\mathcal{D}$ , we have  $\mathbb{E}_{\xi_1^l}[F_{\xi_1^l}(x^l)] = F(x^l)$ ,  $\mathbb{E}_{\xi_2^l}[F_{\xi_2^l}(\tilde{x}^l)] = F(\tilde{x}^l)$ , and

$$\begin{aligned}\mathbb{E}_{\xi_1^l} \left[ \|F_{\xi_1^l}(x^l) - F(x^l)\|^2 \right] &= \frac{1}{m^2} \sum_{i=1}^m \mathbb{E}_{\xi_1^{i,l}} \left[ \|F_{\xi_1^{i,l}}(x^l) - F(x^l)\|^2 \right] \stackrel{(1)}{\leq} \frac{\sigma^2}{m}, \\ \mathbb{E}_{\xi_2^l} \left[ \|F_{\xi_2^l}(\tilde{x}^l) - F(\tilde{x}^l)\|^2 \right] &= \frac{1}{m^2} \sum_{i=1}^m \mathbb{E}_{\xi_2^{i,l}} \left[ \|F_{\xi_2^{i,l}}(\tilde{x}^l) - F(\tilde{x}^l)\|^2 \right] \stackrel{(1)}{\leq} \frac{\sigma^2}{m},\end{aligned}$$

for all  $l = 0, 1, \dots, T-1$ . Moreover, probability event  $E_{T-1}$  implies

$$\begin{aligned}\|F(x^l)\| &\stackrel{(\text{Lip})}{\leq} L\|x^l - x^*\| \leq 3LR \stackrel{(17)}{\leq} \frac{R}{40\gamma \ln \frac{6(K+1)}{\beta}} \stackrel{(18)}{=} \frac{\lambda}{2}, \\ \|F(\tilde{x}^l)\| &\stackrel{(\text{Lip})}{\leq} L\|\tilde{x}^l - x^*\| \stackrel{(23)}{\leq} 4LR \stackrel{(17)}{\leq} \frac{R}{40\gamma \ln \frac{6(K+1)}{\beta}} \stackrel{(18)}{=} \frac{\lambda}{2}\end{aligned}$$

for all  $l = 0, 1, \dots, T-1$ . Therefore, in view of Lemma B.2,  $E_{T-1}$  implies that

$$\|\theta_l^b\| \leq \frac{4\sigma^2}{m\lambda}, \quad \|\omega_l^b\| \leq \frac{4\sigma^2}{m\lambda}, \quad (31)$$

$$\mathbb{E}_{\xi_2^l} \left[ \|\theta_l\|^2 \right] \leq \frac{18\sigma^2}{m}, \quad \mathbb{E}_{\xi_1^l} \left[ \|\omega_l\|^2 \right] \leq \frac{18\sigma^2}{m}, \quad (32)$$

$$\mathbb{E}_{\xi_2^l} \left[ \|\theta_l^u\|^2 \right] \leq \frac{18\sigma^2}{m}, \quad \mathbb{E}_{\xi_1^l} \left[ \|\omega_l^u\|^2 \right] \leq \frac{18\sigma^2}{m}, \quad (33)$$

for all  $l = 0, 1, \dots, T-1$ .

**Upper bound for ①.** Since  $\mathbb{E}_{\xi_2^l}[\theta_l^u] = 0$ , we have

$$\mathbb{E}_{\xi_2^l} [2\gamma \langle \eta_l, \theta_l^u \rangle] = 0.$$

Next, the summands in ① are bounded with probability 1:

$$|2\gamma \langle \eta_l, \theta_l^u \rangle| \leq 2\gamma \|\eta_l\| \cdot \|\theta_l^u\| \stackrel{(26),(30)}{\leq} 20\gamma R\lambda \stackrel{(18)}{=} \frac{R^2}{\ln \frac{6(K+1)}{\beta}} \stackrel{\text{def}}{=} c. \quad (34)$$

Moreover, these summands have bounded conditional variances  $\sigma_l^2 \stackrel{\text{def}}{=} \mathbb{E}_{\xi_2^l} [4\gamma^2 \langle \eta_l, \theta_l^u \rangle^2]$ :

$$\sigma_l^2 \leq \mathbb{E}_{\xi_2^l} [4\gamma^2 \|\eta_l\|^2 \cdot \|\theta_l^u\|^2] \stackrel{(26)}{\leq} 100\gamma^2 R^2 \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2]. \quad (35)$$

That is, sequence  $\{2\gamma \langle \eta_l, \theta_l^u \rangle\}_{l \geq 0}$  is a bounded martingale difference sequence having bounded conditional variances  $\{\sigma_l^2\}_{l \geq 0}$ . Applying the Bernstein's inequality (Lemma B.1) with  $X_l = 2\gamma \langle \eta_l, \theta_l^u \rangle$ ,  $c$  defined in (34),  $b = R^2$ ,  $G = \frac{R^4}{6 \ln \frac{6(K+1)}{\beta}}$ , we get that

$$\mathbb{P} \left\{ |\textcircled{1}| > R^2 \text{ and } \sum_{l=0}^{T-1} \sigma_l^2 \leq \frac{R^4}{6 \ln \frac{6(K+1)}{\beta}} \right\} \leq 2 \exp \left( -\frac{b^2}{2G + 2cb/3} \right) = \frac{\beta}{3(K+1)}.$$

In other words,  $\mathbb{P}\{E_{\textcircled{1}}\} \geq 1 - \frac{\beta}{3(K+1)}$ , where probability event  $E_{\textcircled{1}}$  is defined as

$$E_{\textcircled{1}} = \left\{ \text{either } \sum_{l=0}^{T-1} \sigma_l^2 > \frac{R^4}{6 \ln \frac{6(K+1)}{\beta}} \text{ or } |\textcircled{1}| \leq R^2 \right\}. \quad (36)$$

Moreover, we notice here that probability event  $E_{T-1}$  implies that

$$\sum_{l=0}^{T-1} \sigma_l^2 \stackrel{(35)}{\leq} 100\gamma^2 R^2 \sum_{l=0}^{T-1} \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2] \stackrel{(33), T \leq K+1}{\leq} \frac{1800(K+1)\gamma^2 R^2 \sigma^2}{m} \stackrel{(19)}{\leq} \frac{R^4}{6 \ln \frac{6(K+1)}{\beta}}. \quad (37)$$

**Upper bound for ②.** Probability event  $E_{T-1}$  implies

$$\begin{aligned} \textcircled{2} &\leq 2\gamma \sum_{l=0}^{T-1} \|\eta_l\| \cdot \|\theta_l^b\| \stackrel{(26),(31), T \leq K+1}{\leq} \frac{40(K+1)\gamma R\sigma^2}{m\lambda} \\ &\stackrel{(18)}{=} \frac{40(K+1)\gamma^2\sigma^2 \ln \frac{6(K+1)}{\beta}}{m} \stackrel{(19)}{\leq} R^2. \end{aligned} \quad (38)$$

**Upper bound for ③.** Probability event  $E_{T-1}$  implies

$$2\gamma^2 \sum_{l=0}^{T-1} \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2] \stackrel{(32), T \leq K+1}{\leq} \frac{36\gamma^2(K+1)\sigma^2}{m} \stackrel{(19)}{\leq} \frac{1}{12}R^2, \quad (39)$$

$$4\gamma^2 \sum_{l=0}^{T-1} \mathbb{E}_{\xi_1^l} [\|\omega_l^u\|^2] \stackrel{(32), T \leq K+1}{\leq} \frac{72\gamma^2(K+1)\sigma^2}{m} \stackrel{(19)}{\leq} \frac{1}{12}R^2, \quad (40)$$

$$\textcircled{3} \stackrel{(39),(40)}{\leq} \frac{1}{6}R^2. \quad (41)$$

**Upper bound for ④.** First of all,

$$2\gamma^2 \mathbb{E}_{\xi_1^l, \xi_2^l} [\|\theta_l^u\|^2 + 2\|\omega_l^u\|^2 - \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2] - 2\mathbb{E}_{\xi_1^l} [\|\omega_l^u\|^2]] = 0.$$

Next, the summands in ④ are bounded with probability 1:

$$\begin{aligned} 2\gamma^2 \left| \|\theta_l^u\|^2 + 2\|\omega_l^u\|^2 - \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2] - 2\mathbb{E}_{\xi_1^l} [\|\omega_l^u\|^2] \right| &\leq 2\gamma^2 \|\theta_l^u\|^2 + 2\gamma^2 \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2] \\ &\quad + 4\gamma^2 \|\omega_l^u\|^2 + 4\gamma^2 \mathbb{E}_{\xi_1^l} [\|\omega_l^u\|^2] \\ &\stackrel{(30)}{\leq} 48\gamma^2 \lambda^2 \\ &\stackrel{(18)}{\leq} \frac{R^2}{6 \ln \frac{6(K+1)}{\beta}} \stackrel{\text{def}}{=} c. \end{aligned} \quad (42)$$

Moreover, these summands have bounded conditional variances  $\tilde{\sigma}_l^2 \stackrel{\text{def}}{=} 4\gamma^4 \mathbb{E}_{\xi_1^l, \xi_2^l} \left[ \left| \|\theta_l^u\|^2 + 2\|\omega_l^u\|^2 - \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2] - 2\mathbb{E}_{\xi_1^l} [\|\omega_l^u\|^2] \right|^2 \right]$ :

$$\begin{aligned} \tilde{\sigma}_l^2 &\stackrel{(42)}{\leq} \frac{\gamma^2 R^2}{3 \ln \frac{6(K+1)}{\beta}} \mathbb{E}_{\xi_1^l, \xi_2^l} \left[ \left| \|\theta_l^u\|^2 + 2\|\omega_l^u\|^2 - \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2] - 2\mathbb{E}_{\xi_1^l} [\|\omega_l^u\|^2] \right|^2 \right] \\ &\leq \frac{2\gamma^2 R^2}{3 \ln \frac{6(K+1)}{\beta}} \mathbb{E}_{\xi_1^l, \xi_2^l} [\|\theta_l^u\|^2 + 2\|\omega_l^u\|^2]. \end{aligned} \quad (43)$$

That is, sequence  $\left\{ 2\gamma^2 \left( \|\theta_l^u\|^2 + 2\|\omega_l^u\|^2 - \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2] - 2\mathbb{E}_{\xi_1^l} [\|\omega_l^u\|^2] \right) \right\}_{l \geq 0}$  is a bounded martingale difference sequence having bounded conditional variances  $\{\tilde{\sigma}_l^2\}_{l \geq 0}$ . Applying the Bernstein's inequality (Lemma B.1) with  $X_l = 2\gamma^2 \left( \|\theta_l^u\|^2 + 2\|\omega_l^u\|^2 - \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2] - 2\mathbb{E}_{\xi_1^l} [\|\omega_l^u\|^2] \right)$ ,  $c$  defined in (42),  $b = \frac{1}{6}R^2$ ,  $G = \frac{R^4}{216 \ln \frac{6(K+1)}{\beta}}$ , we get that

$$\mathbb{P} \left\{ |\textcircled{4}| > \frac{1}{6}R^2 \text{ and } \sum_{l=0}^{T-1} \tilde{\sigma}_l^2 \leq \frac{R^4}{216 \ln \frac{6(K+1)}{\beta}} \right\} \leq 2 \exp \left( -\frac{b^2}{2G + 2cb/3} \right) = \frac{\beta}{3(K+1)}.$$

In other words,  $\mathbb{P}\{E_{\textcircled{4}}\} \geq 1 - \frac{\beta}{3(K+1)}$ , where probability event  $E_{\textcircled{4}}$  is defined as

$$E_{\textcircled{4}} = \left\{ \text{either } \sum_{l=0}^{T-1} \tilde{\sigma}_l^2 > \frac{R^4}{216 \ln \frac{6(K+1)}{\beta}} \text{ or } |\textcircled{4}| \leq \frac{1}{6}R^2 \right\}. \quad (44)$$

Moreover, we notice here that probability event  $E_{T-1}$  implies that

$$\begin{aligned} \sum_{l=0}^{T-1} \tilde{\sigma}_l^2 &\stackrel{(43)}{\leq} \frac{2\gamma^2 R^2}{3 \ln \frac{6(K+1)}{\beta}} \sum_{l=0}^{T-1} \mathbb{E}_{\xi_1^l, \xi_2^l} [\|\theta_l^u\|^2 + 2\|\omega_l^u\|^2] \\ &\stackrel{(33), T \leq K+1}{\leq} \frac{36(K+1)\gamma^2 R^2 \sigma^2}{m} \stackrel{(19)}{\leq} \frac{R^4}{216 \ln \frac{6(K+1)}{\beta}}. \end{aligned} \quad (45)$$

**Upper bound for ⑤.** Probability event  $E_{T-1}$  implies

$$\begin{aligned} \textcircled{5} &= 2\gamma^2 \sum_{l=0}^{T-1} (\|\theta_l^b\|^2 + 2\|\omega_l^b\|^2) \stackrel{(31), T \leq K+1}{\leq} \frac{96\gamma^2 \sigma^4 (K+1)}{m^2 \lambda^2} \\ &\stackrel{(18)}{=} \frac{38400\gamma^4 \sigma^4 (K+1) \ln^2 \frac{6(K+1)}{\beta}}{m^2 R^2} \stackrel{(19)}{\leq} \frac{1}{6} R^2. \end{aligned} \quad (46)$$

**Upper bound for  $2\gamma R \left\| \sum_{l=0}^{T-1} \theta_l \right\|$ .** To handle this term, we introduce new notation:

$$\zeta_l = \begin{cases} \gamma \sum_{r=0}^{l-1} \theta_r, & \text{if } \left\| \gamma \sum_{r=0}^{l-1} \theta_r \right\| \leq R, \\ 0, & \text{otherwise} \end{cases}$$

for  $l = 1, 2, \dots, T-1$ . By definition, we have

$$\|\zeta_l\| \leq R. \quad (47)$$

Therefore, in view of (22), probability event  $E_{T-1}$  implies

$$\begin{aligned} 2\gamma R \left\| \sum_{l=0}^{T-1} \theta_l \right\| &= 2R \sqrt{\gamma^2 \left\| \sum_{l=0}^{T-1} \theta_l \right\|^2} \\ &= 2R \sqrt{\gamma^2 \sum_{l=0}^{T-1} \|\theta_l\|^2 + 2\gamma \sum_{l=0}^{T-1} \left\langle \gamma \sum_{r=0}^{l-1} \theta_r, \theta_l \right\rangle} \\ &= 2R \sqrt{\gamma^2 \sum_{l=0}^{T-1} \|\theta_l\|^2 + 2\gamma \sum_{l=0}^{T-1} \langle \zeta_l, \theta_l \rangle} \\ &\stackrel{(27)}{\leq} 2R \sqrt{\textcircled{3} + \textcircled{4} + \textcircled{5} + \underbrace{2\gamma \sum_{l=0}^{T-1} \langle \zeta_l, \theta_l^u \rangle}_{\textcircled{6}} + \underbrace{2\gamma \sum_{l=0}^{T-1} \langle \zeta_l, \theta_l^b \rangle}_{\textcircled{7}}}. \end{aligned} \quad (48)$$

Following similar steps as before, we bound ⑥ and ⑦.

**Upper bound for ⑥.** Since  $\mathbb{E}_{\xi_2^l} [\theta_l^u] = 0$ , we have

$$\mathbb{E}_{\xi_2^l} [2\gamma \langle \zeta_l, \theta_l^u \rangle] = 0.$$

Next, the summands in ④ are bounded with probability 1:

$$|2\gamma \langle \zeta_l, \theta_l^u \rangle| \leq 2\gamma \|\eta_l\| \cdot \|\theta_l^u\| \stackrel{(47), (30)}{\leq} 4\gamma R \lambda \stackrel{(18)}{\leq} \frac{R^2}{4 \ln \frac{6(K+1)}{\beta}} \stackrel{\text{def}}{=} c. \quad (49)$$

Moreover, these summands have bounded conditional variances  $\hat{\sigma}_l^2 \stackrel{\text{def}}{=} \mathbb{E}_{\xi_2^l} [4\gamma^2 \langle \zeta_l, \theta_l^u \rangle^2]$ :

$$\hat{\sigma}_l^2 \leq \mathbb{E}_{\xi_2^l} [4\gamma^2 \|\zeta_l\|^2 \cdot \|\theta_l^u\|^2] \stackrel{(47)}{\leq} 4\gamma^2 R^2 \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2]. \quad (50)$$

That is, sequence  $\{2\gamma\langle\zeta_l, \theta_l^u\rangle\}_{l \geq 0}$  is a bounded martingale difference sequence having bounded conditional variances  $\{\hat{\sigma}_l^2\}_{l \geq 0}$ . Applying Bernstein's inequality (Lemma B.1) with  $X_l = 2\gamma\langle\zeta_l, \theta_l^u\rangle$ ,  $c$  defined in (34),  $b = \frac{R^2}{4}$ ,  $G = \frac{R^4}{96 \ln \frac{6(K+1)}{\beta}}$ , we get that

$$\mathbb{P}\left\{|\textcircled{5}| > \frac{1}{4}R^2 \text{ and } \sum_{l=0}^{T-1} \hat{\sigma}_l^2 \leq \frac{R^4}{96 \ln \frac{6(K+1)}{\beta}}\right\} \leq 2 \exp\left(-\frac{b^2}{2G + 2cb/3}\right) = \frac{\beta}{3(K+1)}.$$

In other words,  $\mathbb{P}\{E_{\textcircled{6}}\} \geq 1 - \frac{\beta}{3(K+1)}$ , where probability event  $E_{\textcircled{6}}$  is defined as

$$E_{\textcircled{6}} = \left\{ \text{either } \sum_{l=0}^{T-1} \hat{\sigma}_l^2 > \frac{R^4}{96 \ln \frac{6(K+1)}{\beta}} \text{ or } |\textcircled{5}| \leq \frac{1}{4}R^2 \right\}. \quad (51)$$

Moreover, we notice here that probability event  $E_{T-1}$  implies that

$$\sum_{l=0}^{T-1} \hat{\sigma}_l^2 \stackrel{(50)}{\leq} 4\gamma^2 R^2 \sum_{l=0}^{T-1} \mathbb{E}_{\xi_2^{i_l}} [\|\theta_l^u\|^2] \stackrel{(33), T \leq K+1}{\leq} \frac{72(K+1)\gamma^2 R^2 \sigma^2}{m} \stackrel{(19)}{\leq} \frac{R^4}{96 \ln \frac{6(K+1)}{\beta}}. \quad (52)$$

**Upper bound for ⑦.** Probability event  $E_{T-1}$  implies

$$\begin{aligned} \textcircled{7} &\leq 2\gamma \sum_{l=0}^{T-1} \|\zeta_l\| \cdot \|\theta_l^b\| \stackrel{(47), (31), T \leq K+1}{\leq} \frac{8(K+1)\gamma R \sigma^2}{m\lambda} \\ &\stackrel{(18)}{\leq} \frac{160(K+1)\gamma^2 \sigma^2 \ln \frac{6(K+1)}{\beta}}{m} \stackrel{(19)}{\leq} \frac{1}{4}R^2. \end{aligned} \quad (53)$$

**Final derivation.** Putting all bounds together, we get that  $E_{T-1}$  implies

$$\begin{aligned} A_T &\stackrel{(29)}{\leq} 4R^2 + 2\gamma R \left\| \sum_{l=0}^{T-1} \theta_l \right\| + \textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4} + \textcircled{5}, \\ &2\gamma R \left\| \sum_{l=0}^{T-1} \theta_l \right\| \stackrel{(48)}{\leq} 2R\sqrt{\textcircled{3} + \textcircled{4} + \textcircled{5} + \textcircled{6} + \textcircled{7}}, \\ &\textcircled{2} \stackrel{(38)}{\leq} R^2, \quad \textcircled{3} \stackrel{(41)}{\leq} \frac{1}{6}R^2, \quad \textcircled{5} \stackrel{(46)}{\leq} \frac{1}{6}R^2, \quad \textcircled{7} \stackrel{(53)}{\leq} \frac{1}{4}R^2, \\ &\sum_{l=0}^{T-1} \sigma_l^2 \stackrel{(37)}{\leq} \frac{R^4}{6 \ln \frac{6(K+1)}{\beta}}, \quad \sum_{l=0}^{T-1} \tilde{\sigma}_l^2 \stackrel{(45)}{\leq} \frac{R^4}{216 \ln \frac{6(K+1)}{\beta}}, \quad \sum_{l=0}^{T-1} \hat{\sigma}_l^2 \stackrel{(52)}{\leq} \frac{R^4}{96 \ln \frac{6(K+1)}{\beta}}. \end{aligned}$$

Moreover, in view of (36), (44), (51), and our induction assumption, we have

$$\begin{aligned} \mathbb{P}\{E_{T-1}\} &\geq 1 - \frac{(T-1)\beta}{K+1}, \\ \mathbb{P}\{E_{\textcircled{1}}\} &\geq 1 - \frac{\beta}{3(K+1)}, \quad \mathbb{P}\{E_{\textcircled{4}}\} \geq 1 - \frac{\beta}{3(K+1)}, \quad \mathbb{P}\{E_{\textcircled{6}}\} \geq 1 - \frac{\beta}{3(K+1)}, \end{aligned}$$

where probability events  $E_{\textcircled{1}}$ ,  $E_{\textcircled{4}}$ , and  $E_{\textcircled{6}}$  are defined as

$$\begin{aligned} E_{\textcircled{1}} &= \left\{ \text{either } \sum_{l=0}^{T-1} \sigma_l^2 > \frac{R^4}{6 \ln \frac{6(K+1)}{\beta}} \text{ or } |\textcircled{1}| \leq R^2 \right\}, \\ E_{\textcircled{4}} &= \left\{ \text{either } \sum_{l=0}^{T-1} \tilde{\sigma}_l^2 > \frac{R^4}{216 \ln \frac{6(K+1)}{\beta}} \text{ or } |\textcircled{4}| \leq \frac{1}{6}R^2 \right\}, \\ E_{\textcircled{6}} &= \left\{ \text{either } \sum_{l=0}^{T-1} \hat{\sigma}_l^2 > \frac{R^4}{96 \ln \frac{6(K+1)}{\beta}} \text{ or } |\textcircled{6}| \leq \frac{1}{4}R^2 \right\}. \end{aligned}$$

Putting all of these inequalities together, we obtain that probability event  $E_{T-1} \cap E_{\textcircled{1}} \cap E_{\textcircled{4}} \cap E_{\textcircled{6}}$  implies

$$\left\| \gamma \sum_{l=0}^{T-1} \theta_l \right\| \leq \sqrt{\frac{1}{6}R^2 + \frac{1}{6}R^2 + \frac{1}{6}R^2 + \frac{1}{4}R^2 + \frac{1}{4}R^2} = R, \quad (54)$$

$$\begin{aligned} A_T &\leq 4R^2 + 2R\sqrt{\frac{1}{6}R^2 + \frac{1}{6}R^2 + \frac{1}{6}R^2 + \frac{1}{4}R^2 + \frac{1}{4}R^2} \\ &\quad + R^2 + R^2 + \frac{1}{6}R^2 + \frac{1}{6}R^2 + \frac{1}{6}R^2 \\ &\leq 9R^2. \end{aligned} \quad (55)$$

Moreover, union bound for the probability events implies

$$\mathbb{P}\{E_T\} \geq \mathbb{P}\{E_{T-1} \cap E_{\textcircled{1}} \cap E_{\textcircled{4}} \cap E_{\textcircled{6}}\} = 1 - \mathbb{P}\{\bar{E}_{T-1} \cup \bar{E}_{\textcircled{1}} \cup \bar{E}_{\textcircled{4}} \cup \bar{E}_{\textcircled{6}}\} \geq 1 - \frac{T\beta}{K+1}. \quad (56)$$

This is exactly what we wanted to prove (see the paragraph after inequalities (21), (22)). Therefore, for all  $k = 0, 1, \dots, K+1$  we have  $\mathbb{P}\{E_k\} \geq 1 - k\beta/(K+1)$ , i.e., for  $k = K+1$  we have that with probability at least  $1 - \beta$  inequality

$$\begin{aligned} \text{Gap}_R(\tilde{x}_{\text{avg}}^K) &= \max_{u \in B_R(x^*)} \{ \langle F(u), \tilde{x}_{\text{avg}}^K - u \rangle \} \\ &\leq \frac{1}{2\gamma(K+1)} \max_{u \in B_R(x^*)} \{ 2\gamma(K+1) \langle F(u), \tilde{x}_{\text{avg}}^t - u \rangle + \|x^{K+1} - u\|^2 \} \\ &\stackrel{(24)}{\leq} \frac{9R^2}{2\gamma(K+1)} \end{aligned}$$

holds. This concludes the proof.  $\square$

**Corollary C.1.** *Let the assumptions of Theorem C.1 hold. Then, the following statements hold.*

1. **Large stepsize/large batch.** *The choice of stepsize and batchsize*

$$\gamma = \frac{1}{160L \ln \frac{6(K+1)}{\beta}}, \quad m = \max \left\{ 1, \frac{27(K+1)\sigma^2}{64L^2 R^2 \ln \frac{6(K+1)}{\beta}} \right\} \quad (57)$$

satisfies conditions (17) and (19). With such choice of  $\gamma$ ,  $m$ , and the choice of  $\lambda$  as in (18), the iterates produced by **clipped-SEG** after  $K$  iterations with probability at least  $1 - \beta$  satisfy

$$\text{Gap}_R(\tilde{x}_{\text{avg}}^K) \leq \frac{720LR^2 \ln \frac{6(K+1)}{\beta}}{K+1}. \quad (58)$$

In particular, to guarantee  $\text{Gap}_R(\tilde{x}_{\text{avg}}^K) \leq \varepsilon$  with probability at least  $1 - \beta$  for some  $\varepsilon > 0$  **clipped-SEG** requires,

$$\mathcal{O} \left( \frac{LR^2}{\varepsilon} \ln \left( \frac{LR^2}{\varepsilon\beta} \right) \right) \text{ iterations}, \quad (59)$$

$$\mathcal{O} \left( \max \left\{ \frac{LR^2}{\varepsilon}, \frac{\sigma^2 R^2}{\varepsilon^2} \right\} \ln \left( \frac{LR^2}{\varepsilon\beta} \right) \right) \text{ oracle calls}. \quad (60)$$

2. **Small stepsize/small batch.** *The choice of stepsize and batchsize*

$$\gamma = \min \left\{ \frac{1}{160L \ln \frac{6(K+1)}{\beta}}, \frac{R}{60\sigma\sqrt{3(K+1) \ln \frac{6(K+1)}{\beta}}} \right\}, \quad m = 1 \quad (61)$$

satisfies conditions (17) and (19). With such choice of  $\gamma$ ,  $m$ , and the choice of  $\lambda$  as in (18), the iterates produced by **clipped-SEG** after  $K$  iterations with probability at least  $1 - \beta$  satisfy

$$\text{Gap}_R(\tilde{x}_{\text{avg}}^K) \leq \max \left\{ \frac{720LR^2 \ln \frac{6(K+1)}{\beta}}{K+1}, \frac{270\sigma R \sqrt{\ln \frac{6(K+1)}{\beta}}}{\sqrt{K+1}} \right\}. \quad (62)$$



In particular, to guarantee  $\text{Gap}_R(\tilde{x}_{\text{avg}}^K) \leq \varepsilon$  with probability at least  $1 - \beta$  for some  $\varepsilon > 0$ , **clipped-SEG** requires

$$\mathcal{O} \left( \max \left\{ \frac{LR^2}{\varepsilon} \ln \left( \frac{LR^2}{\varepsilon\beta} \right), \frac{\sigma^2 R^2}{\varepsilon^2} \ln \left( \frac{\sigma^2 R^2}{\varepsilon^2 \beta} \right) \right\} \right) \text{ iterations/oracle calls.} \quad (63)$$

*Proof.* 1. **Large stepsize/large batch.** First of all, we verify that the choice of  $\gamma$  and  $m$  from (57) satisfies conditions (17) and (19): (17) trivially holds and (19) holds since

$$m = \max \left\{ 1, \frac{27(K+1)\sigma^2}{64L^2 R^2 \ln \frac{6(K+1)}{\beta}} \right\} = \max \left\{ 1, \frac{10800(K+1)\gamma^2 \sigma^2 \ln \frac{6(K+1)}{\beta}}{R^2} \right\}.$$

Therefore, applying Theorem C.1, we derive that with probability at least  $1 - \beta$

$$\text{Gap}_R(\tilde{x}_{\text{avg}}^K) \leq \frac{9R^2}{2\gamma(K+1)} \stackrel{(57)}{=} \frac{720LR^2 \ln \frac{4(K+1)}{\beta}}{K+1}.$$

To guarantee  $\text{Gap}_R(\tilde{x}_{\text{avg}}^K) \leq \varepsilon$ , we choose  $K$  in such a way that the right-hand side of the above inequality is smaller than  $\varepsilon$  that gives

$$K = \mathcal{O} \left( \frac{LR^2}{\varepsilon} \ln \left( \frac{LR^2}{\varepsilon\beta} \right) \right).$$

The total number of oracle calls equals

$$\begin{aligned} 2m(K+1) &\stackrel{(57)}{=} 2 \max \left\{ K+1, \frac{27(K+1)^2 \sigma^2}{64L^2 R^2 \ln \frac{6(K+1)}{\beta}} \right\} \\ &= \mathcal{O} \left( \max \left\{ \frac{LR^2}{\varepsilon}, \frac{\sigma^2 R^2}{\varepsilon^2} \right\} \ln \left( \frac{LR^2}{\varepsilon\beta} \right) \right). \end{aligned}$$

2. **Small stepsize/small batch.** First of all, we verify that the choice of  $\gamma$  and  $m$  from (57) satisfies conditions (17) and (19):

$$\begin{aligned} \gamma &= \min \left\{ \frac{1}{160L \ln \frac{6(K+1)}{\beta}}, \frac{R}{60\sigma \sqrt{3(K+1) \ln \frac{6(K+1)}{\beta}}} \right\} \leq \frac{1}{160L \ln \frac{6(K+1)}{\beta}}, \\ m &= 1 \stackrel{(61)}{\geq} \frac{10800(K+1)\gamma^2 \sigma^2 \ln \frac{6(K+1)}{\beta}}{R^2}. \end{aligned}$$

Therefore, applying Theorem C.1, we derive that with probability at least  $1 - \beta$

$$\begin{aligned} \text{Gap}_R(\tilde{x}_{\text{avg}}^K) &\leq \frac{9R^2}{2\gamma(K+1)} \\ &\stackrel{(61)}{=} \max \left\{ \frac{720LR^2 \ln \frac{6(K+1)}{\beta}}{K+1}, \frac{270\sigma R \sqrt{\ln \frac{6(K+1)}{\beta}}}{\sqrt{K+1}} \right\}. \end{aligned}$$

To guarantee  $\text{Gap}_R(\tilde{x}_{\text{avg}}^K) \leq \varepsilon$ , we choose  $K$  in such a way that the right-hand side of the above inequality is smaller than  $\varepsilon$  that gives

$$K = \mathcal{O} \left( \max \left\{ \frac{LR^2}{\varepsilon} \ln \left( \frac{LR^2}{\varepsilon\beta} \right), \frac{\sigma^2 R^2}{\varepsilon^2} \ln \left( \frac{\sigma^2 R^2}{\varepsilon^2 \beta} \right) \right\} \right).$$

The total number of oracle calls equals  $2m(K+1) = 2(K+1)$ . □

## C.2 Star-Negative Comonotone Case

**Lemma C.2.** *Let Assumptions 1.2, 1.4 hold for  $Q = B_{3R}(x^*) = \{x \in \mathbb{R}^d \mid \|x - x^*\| \leq 3R\}$ , where  $R \geq R_0 \stackrel{\text{def}}{=} \|x^0 - x^*\|$ , and  $\gamma_2 + 2\rho < \gamma_1 \leq 1/(2L)$ . If  $x^k$  and  $\tilde{x}^k$  lie in  $B_{3R_0}(x^*)$  for all  $k = 0, 1, \dots, K$  for some  $K \geq 0$ , then the iterates produced by **clipped-SEG** satisfy*

$$\begin{aligned} \frac{\gamma_1 \gamma_2}{4(K+1)} \sum_{k=0}^K \|F(x^k)\|^2 &\leq \frac{\|x^0 - x^*\|^2 - \|x^{K+1} - x^*\|^2}{K+1} \\ &\quad + \frac{1}{K+1} \sum_{k=0}^K (\gamma_2^2 \|\omega_k\|^2 + 3\gamma_1 \gamma_2 \|\omega_k\|^2) \\ &\quad + \frac{2\gamma_2}{K+1} \sum_{k=0}^K \langle x^k - x^* - \gamma_2 F(\tilde{x}^k), \theta_k \rangle \end{aligned} \quad (64)$$

where  $\theta_k, \omega_k$  are defined in (15), (16).

*Proof.* Using the update rule of **clipped-SEG**, we obtain

$$\begin{aligned} \|x^{k+1} - x^*\|^2 &= \|x^k - x^*\|^2 - 2\gamma_2 \langle x^k - x^*, \tilde{F}_{\xi_2^k}(\tilde{x}^k) \rangle + \gamma_2^2 \|\tilde{F}_{\xi_2^k}(\tilde{x}^k)\|^2 \\ &= \|x^k - x^*\|^2 - 2\gamma_2 \langle x^k - x^*, F(\tilde{x}^k) \rangle + 2\gamma_2 \langle x^k - x^*, \theta_k \rangle \\ &\quad + \gamma_2^2 \|F(\tilde{x}^k)\|^2 - 2\gamma_2^2 \langle F(\tilde{x}^k), \theta_k \rangle + \gamma_2^2 \|\theta_k\|^2 \\ &= \|x^k - x^*\|^2 - 2\gamma_2 \langle \tilde{x}^k - x^*, F(\tilde{x}^k) \rangle - 2\gamma_2 \langle x^k - \tilde{x}^k, F(\tilde{x}^k) \rangle \\ &\quad + 2\gamma_2 \langle x^k - x^* - \gamma_2 F(\tilde{x}^k), \theta_k \rangle + \gamma_2^2 \|F(\tilde{x}^k)\|^2 + \gamma_2^2 \|\theta_k\|^2 \\ &\stackrel{\text{(SNC)}}{\leq} \|x^k - x^*\|^2 + 2\gamma_2 \rho \|\tilde{F}(\tilde{x}^k)\|^2 - 2\gamma_1 \gamma_2 \langle \tilde{F}_{\xi_1^k}(x^k), F(\tilde{x}^k) \rangle \\ &\quad + 2\gamma_2 \langle x^k - x^* - \gamma_2 F(\tilde{x}^k), \theta_k \rangle + \gamma_2^2 \|F(\tilde{x}^k)\|^2 + \gamma_2^2 \|\theta_k\|^2 \\ &\stackrel{(4)}{=} \|x^k - x^*\|^2 + \gamma_1 \gamma_2 \|\tilde{F}_{\xi_1^k}(x^k) - F(\tilde{x}^k)\|^2 - \gamma_1 \gamma_2 \|F(\tilde{x}^k)\|^2 - \gamma_1 \gamma_2 \|\tilde{F}_{\xi_1^k}(x^k)\|^2 \\ &\quad + 2\gamma_2 \langle x^k - x^* - \gamma_2 F(\tilde{x}^k), \theta_k \rangle + \gamma_2 (2\rho + \gamma_2) \|F(\tilde{x}^k)\|^2 + \gamma_2^2 \|\theta_k\|^2 \\ &\stackrel{(5)}{\leq} \|x^k - x^*\|^2 + 2\gamma_1 \gamma_2 \|\omega_k\|^2 + 2\gamma_1 \gamma_2 \|F(x^k) - F(\tilde{x}^k)\|^2 - \gamma_1 \gamma_2 \|\tilde{F}_{\xi_1^k}(x^k)\|^2 \\ &\quad + 2\gamma_2 \langle x^k - x^* - \gamma_2 F(\tilde{x}^k), \theta_k \rangle + \gamma_2 (2\rho + \gamma_2 - \gamma_1) \|F(\tilde{x}^k)\|^2 + \gamma_2^2 \|\theta_k\|^2 \\ &\stackrel{\text{(Lip)}}{\leq} \|x^k - x^*\|^2 - \gamma_1 \gamma_2 (1 - 2\gamma_1^2 L^2) \|\tilde{F}_{\xi_1^k}(x^k)\|^2 \\ &\quad + 2\gamma_2 \langle x^k - x^* - \gamma_2 F(\tilde{x}^k), \theta_k \rangle + \gamma_2^2 \|\theta_k\|^2 + 2\gamma_1 \gamma_2 \|\omega_k\|^2, \end{aligned}$$

where in the last step we additionally use  $x^k - \tilde{x}^k = \gamma_1 \tilde{F}_{\xi_1^k}(x^k)$  after the application of Lipschitzness of  $F$  and we use our assumption on  $\gamma_1, \gamma_2, \rho$ :  $\gamma_2 + 2\rho \leq \gamma_1$ . Since  $\gamma_1 \leq 1/(2L)$ , we have  $\gamma_1 \gamma_2 (1 - 2\gamma_1^2 L^2) \|\tilde{F}_{\xi_1^k}(x^k)\|^2 \geq 0$  and, using (6) with  $\alpha = 1$ , we derive

$$\begin{aligned} \|x^{k+1} - x^*\|^2 &\leq \|x^k - x^*\|^2 - \frac{\gamma_1 \gamma_2}{2} (1 - 2\gamma_1^2 L^2) \|F(x^k)\|^2 \\ &\quad + 2\gamma_2 \langle x^k - x^* - \gamma_2 F(\tilde{x}^k), \theta_k \rangle + \gamma_2^2 \|\theta_k\|^2 + 2\gamma_1 \gamma_2 \left( \frac{3}{2} - \gamma_1^2 L^2 \right) \|\omega_k\|^2. \end{aligned}$$

Rearranging the terms and using  $\frac{3}{2} - \gamma_1^2 L^2 \leq \frac{3}{2}$ ,  $1 - 2\gamma_1^2 L^2 \geq 1/2$ , we derive

$$\begin{aligned} \frac{\gamma_1 \gamma_2}{4} \|F(x^k)\|^2 &\leq \|x^k - x^*\|^2 - \|x^{k+1} - x^*\|^2 + (\gamma_2^2 \|\theta_k\|^2 + 3\gamma_1 \gamma_2 \|\omega_k\|^2) \\ &\quad + 2\gamma_2 \langle x^k - x^* - \gamma_2 F(\tilde{x}^k), \theta_k \rangle. \end{aligned}$$

Finally, we sum up the above inequality for  $k = 0, 1, \dots, K$  and divide both sides of the result by  $(K + 1)$ :

$$\begin{aligned}
\frac{\gamma_1 \gamma_2}{4(K+1)} \sum_{k=0}^K \|F(x^k)\|^2 &\leq \frac{1}{K+1} \sum_{k=0}^K (\|x^k - x^*\|^2 - \|x^{k+1} - x^*\|^2) + \frac{\gamma_2^2}{K+1} \sum_{k=0}^K \|\theta_k\|^2 \\
&\quad + \frac{2\gamma_2}{K+1} \sum_{k=0}^K \langle x^k - x^* - \gamma_2 F(\tilde{x}^k), \theta_k \rangle + \frac{3\gamma_1 \gamma_2}{K+1} \sum_{k=0}^K \|\omega_k\|^2 \\
&= \frac{\|x^0 - x^*\|^2 - \|x^{K+1} - x^*\|^2}{K+1} \\
&\quad + \frac{1}{K+1} \sum_{k=0}^K (\gamma_2^2 \|\theta_k\|^2 + 3\gamma_1 \gamma_2 \|\omega_k\|^2) \\
&\quad + \frac{2\gamma_2}{K+1} \sum_{k=0}^K \langle x^k - x^* - \gamma_2 F(\tilde{x}^k), \theta_k \rangle.
\end{aligned}$$

This finishes the proof.  $\square$

**Theorem C.2.** *Let Assumptions 1.1, 1.2, 1.4 hold for  $Q = B_{3R}(x^*)$ , where  $R \geq R_0 \stackrel{\text{def}}{=} \|x^0 - x^*\|$ , and*

$$\gamma_2 + 2\rho \leq \gamma_1 \leq \frac{1}{160L \ln \frac{6(K+1)}{\beta}}, \quad (65)$$

$$\lambda_{1,k} \equiv \lambda_1 = \frac{R}{20\gamma_1 \ln \frac{6(K+1)}{\beta}}, \quad \lambda_{1,k} \equiv \lambda_2 = \frac{R}{20\gamma_2 \ln \frac{6(K+1)}{\beta}}, \quad (66)$$

$$m_{1,k} \equiv m_1 \geq \max \left\{ 1, \frac{216 \max\{\gamma_1 \gamma_2 (K+1), \sqrt{\gamma_1^3 \gamma_2 (K+1)} \ln \frac{6(K+1)}{\beta}\} \sigma^2}{R^2} \right\}, \quad (67)$$

$$m_{2,k} \equiv m_2 \geq \max \left\{ 1, \frac{3240(K+1)\gamma_2^2 \sigma^2 \ln \frac{6(K+1)}{\beta}}{R^2} \right\}, \quad (68)$$

for some  $K \geq 0$  and  $\beta \in (0, 1]$  such that  $\ln \frac{6(K+1)}{\beta} \geq 1$ . Then, after  $K$  iterations the iterates produced by [clipped-SEG](#) with probability at least  $1 - \beta$  satisfy

$$\frac{1}{K+1} \sum_{k=0}^K \|F(x^k)\|^2 \leq \frac{36R^2}{\gamma_1 \gamma_2 (K+1)}. \quad (69)$$

*Proof.* As in the proof of Theorem C.1, we use the following notation:  $R_k = \|x^k - x^*\|^2$ ,  $k \geq 0$ . We will derive (69) by induction. In particular, for each  $k = 0, \dots, K+1$  we define probability event  $E_k$  as follows: inequalities

$$R_t^2 \leq 4R^2 \quad (70)$$

hold for  $t = 0, 1, \dots, k$  simultaneously. Our goal is to prove that  $\mathbb{P}\{E_k\} \geq 1 - k\beta/(K+1)$  for all  $k = 0, 1, \dots, K+1$ . We use the induction to show this statement. For  $k = 0$  the statement is trivial since  $R_0^2 \leq 4R^2$  by definition. Next, assume that the statement holds for  $k = T-1 \leq K$ , i.e., we have  $\mathbb{P}\{E_{T-1}\} \geq 1 - (T-1)\beta/(K+1)$ . We need to prove that  $\mathbb{P}\{E_T\} \geq 1 - T\beta/(K+1)$ . First of all, since  $R_t^2 \leq 4R^2$ , we have  $x^t \in B_{2R}(x^*)$ . Operator  $F$  is  $L$ -Lipschitz on  $B_{3R}(x^*)$ . Therefore, probability event  $E_{T-1}$  implies

$$\|F(x^t)\| \leq L\|x^t - x^*\| \stackrel{(70)}{\leq} 2LR \stackrel{(65),(66)}{\leq} \frac{\lambda_1}{2}. \quad (71)$$

and

$$\|\omega_t\|^2 \stackrel{(5)}{\leq} 2\|\tilde{F}_{\xi_1}(x^t)\|^2 + 2\|F(x^t)\|^2 \stackrel{(71)}{\leq} \frac{5}{2}\lambda_1^2 \stackrel{(66)}{\leq} \frac{R^2}{4\gamma_1^2} \quad (72)$$

for all  $t = 0, 1, \dots, T-1$ .

Next, we show that probability event  $E_{T-1}$  implies  $\|\tilde{x}^t - x^*\| \leq 3R$  and derive useful inequalities related to  $\theta_t$  for all  $t = 0, 1, \dots, T-1$ . Indeed, due to Lipschitzness of  $F$  probability event  $E_{T-1}$  implies

$$\begin{aligned}
\|\tilde{x}^t - x^*\|^2 &= \|x^t - x^* - \gamma_1 \tilde{F}_{\xi_1}(x^t)\|^2 \stackrel{(5)}{\leq} 2\|x^t - x^*\|^2 + 2\gamma_1^2 \|\tilde{F}_{\xi_1}(x^t)\|^2 \\
&\stackrel{(5)}{\leq} 2R_t^2 + 4\gamma_1^2 \|F(x^t)\|^2 + 4\gamma_1^2 \|\omega_t\|^2 \\
&\stackrel{(\text{Lip})}{\leq} 2(1 + 2\gamma_1^2 L^2)R_t^2 + 4\gamma_1^2 \|\omega_t\|^2 \\
&\stackrel{(65),(72)}{\leq} 7R^2 \leq 9R^2
\end{aligned} \tag{73}$$

and

$$\|F(\tilde{x}^t)\| \leq L\|\tilde{x}^t - x^*\| \stackrel{(65),(66)}{\leq} \sqrt{7}LR \leq \frac{\lambda_2}{2} \tag{74}$$

for all  $t = 0, 1, \dots, T-1$ .

That is,  $E_{T-1}$  implies that  $x^t, \tilde{x}^t \in B_{3R}(x^*)$  for all  $t = 0, 1, \dots, T-1$ . Applying Lemma C.2, we get that probability event  $E_{T-1}$  implies

$$\begin{aligned}
\frac{\gamma_1 \gamma_2}{4T} \sum_{l=0}^{T-1} \|F(x^l)\|^2 &\leq \frac{R^2 - R_T^2}{T} + \frac{2\gamma_2}{T} \sum_{l=0}^{T-1} \langle x^l - x^* - \gamma_2 F(\tilde{x}^l), \theta_l \rangle \\
&\quad + \frac{1}{T} \sum_{l=0}^{T-1} (\gamma_2^2 \|\theta_l\|^2 + 3\gamma_1 \gamma_2 \|\omega_l\|^2) \\
R_T^2 &\leq R^2 + 2\gamma_2 \sum_{l=0}^{T-1} \langle x^l - x^* - \gamma_2 F(\tilde{x}^l), \theta_l \rangle + \sum_{l=0}^{T-1} (\gamma_2^2 \|\theta_l\|^2 + 3\gamma_1 \gamma_2 \|\omega_l\|^2).
\end{aligned} \tag{75}$$

To estimate the sums in the right-hand side, we introduce new vectors:

$$\eta_t = \begin{cases} x^t - x^* - \gamma_2 F(\tilde{x}^t), & \text{if } \|x^t - x^* - \gamma_2 F(\tilde{x}^t)\| \leq \sqrt{7}(1 + \gamma_2 L)R, \\ 0, & \text{otherwise,} \end{cases} \tag{76}$$

for  $t = 0, 1, \dots, T-1$ . First of all, we point out that vectors  $\zeta_t$  and  $\eta_t$  are bounded with probability 1, i.e., with probability 1

$$\|\eta_t\| \leq \sqrt{7}(1 + \gamma_2 L)R \tag{77}$$

for all  $t = 0, 1, \dots, T-1$ . Next, we notice that  $E_{T-1}$  implies

$$\begin{aligned}
\|x^t - x^* - \gamma_2 F(\tilde{x}^t)\| &\leq \|x^t - x^*\| + \gamma_2 \|F(\tilde{x}^t)\| \\
&\stackrel{(73),(74)}{\leq} \sqrt{7}(1 + \gamma_2 L)R
\end{aligned}$$

for  $t = 0, 1, \dots, T-1$ , i.e., probability event  $E_{T-1}$  implies  $\eta_t = x^t - x^* - \gamma_2 F(\tilde{x}^t)$  for all  $t = 0, 1, \dots, T-1$ . Therefore,  $E_{T-1}$  implies

$$R_T^2 \leq R^2 + 2\gamma_2 \sum_{l=0}^{T-1} \langle \eta_l, \theta_l \rangle + \sum_{l=0}^{T-1} (\gamma_2^2 \|\theta_l\|^2 + 3\gamma_1 \gamma_2 \|\omega_l\|^2).$$

As in the monotone case, to continue the derivation, we introduce vectors  $\theta_l^u, \theta_l^b, \omega_l^u, \omega_l^b$  defined as

$$\theta_l^u \stackrel{\text{def}}{=} \mathbb{E}_{\xi_2^l} [\tilde{F}_{\xi_2^l}(\tilde{x}^l)] - \tilde{F}_{\xi_2^l}(\tilde{x}^l), \quad \theta_l^b \stackrel{\text{def}}{=} F(\tilde{x}^l) - \mathbb{E}_{\xi_2^l} [\tilde{F}_{\xi_2^l}(\tilde{x}^l)], \tag{78}$$

$$\omega_l^u \stackrel{\text{def}}{=} \mathbb{E}_{\xi_1^l} [\tilde{F}_{\xi_1^l}(x^l)] - \tilde{F}_{\xi_1^l}(x^l), \quad \theta_l^b \stackrel{\text{def}}{=} F(x^l) - \mathbb{E}_{\xi_1^l} [\tilde{F}_{\xi_1^l}(x^l)], \tag{79}$$

for all  $l = 0, \dots, T-1$ . By definition we have  $\theta_l = \theta_l^u + \theta_l^b$ ,  $\omega_l = \omega_l^u + \omega_l^b$  for all  $l = 0, \dots, T-1$ . Using the introduced notation, we continue our derivation as follows:  $E_{T-1}$  implies

$$\begin{aligned}
R_T^2 &\stackrel{(5)}{\leq} R^2 + \underbrace{2\gamma_2 \sum_{l=0}^{T-1} \langle \eta_l, \theta_l^u \rangle}_{\textcircled{1}} + \underbrace{2\gamma_2 \sum_{l=0}^{T-1} \langle \eta_l, \theta_l^b \rangle}_{\textcircled{2}} + \underbrace{2\gamma_2^2 \sum_{l=0}^{T-1} \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2]}_{\textcircled{3}} \\
&\quad + \underbrace{2\gamma_2^2 \sum_{l=0}^{T-1} \left( \|\theta_l^u\|^2 - \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2] \right)}_{\textcircled{4}} + \underbrace{2\gamma_2^2 \sum_{l=0}^{T-1} \|\theta_l^b\|^2}_{\textcircled{5}} + \underbrace{6\gamma_1\gamma_2 \sum_{l=0}^{T-1} \mathbb{E}_{\xi_1^l} [\|\omega_l^u\|^2]}_{\textcircled{6}} \\
&\quad + \underbrace{6\gamma_1\gamma_2 \sum_{l=0}^{T-1} \left( \|\omega_l^u\|^2 - \mathbb{E}_{\xi_1^l} [\|\omega_l^u\|^2] \right)}_{\textcircled{7}} + \underbrace{6\gamma_1\gamma_2 \sum_{l=0}^{T-1} \|\omega_l^b\|^2}_{\textcircled{8}}. \tag{80}
\end{aligned}$$

The rest of the proof is based on deriving good enough upper bounds for  $\textcircled{1}, \textcircled{2}, \textcircled{3}, \textcircled{4}, \textcircled{5}, \textcircled{6}, \textcircled{7}, \textcircled{8}$ , i.e., we want to prove that  $\textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4} + \textcircled{5} + \textcircled{6} + \textcircled{7} + \textcircled{8} \leq 8R^2$  with high probability.

Before we move on, we need to derive some useful inequalities for operating with  $\theta_l^u, \theta_l^b, \omega_l^u, \omega_l^b$ . First of all, Lemma B.2 implies that

$$\|\theta_l^u\| \leq 2\lambda_2, \quad \|\omega_l^u\| \leq 2\lambda_1 \tag{81}$$

for all  $l = 0, 1, \dots, T-1$ . Next, since  $\{\xi_1^{i,l}\}_{i=1}^{m_1}, \{\xi_2^{i,l}\}_{i=1}^{m_2}$  are independently sampled from  $\mathcal{D}$ , we have  $\mathbb{E}_{\xi_1^l} [F_{\xi_1^l}(x^l)] = F(x^l)$ ,  $\mathbb{E}_{\xi_2^l} [F_{\xi_2^l}(\tilde{x}^l)] = F(\tilde{x}^l)$ , and

$$\begin{aligned}
\mathbb{E}_{\xi_1^l} \left[ \|F_{\xi_1^l}(x^l) - F(x^l)\|^2 \right] &= \frac{1}{m_1^2} \sum_{i=1}^{m_1} \mathbb{E}_{\xi_1^{i,l}} \left[ \|F_{\xi_1^{i,l}}(x^l) - F(x^l)\|^2 \right] \stackrel{(1)}{\leq} \frac{\sigma^2}{m_1}, \\
\mathbb{E}_{\xi_2^l} \left[ \|F_{\xi_2^l}(\tilde{x}^l) - F(\tilde{x}^l)\|^2 \right] &= \frac{1}{m_2^2} \sum_{i=1}^{m_2} \mathbb{E}_{\xi_2^{i,l}} \left[ \|F_{\xi_2^{i,l}}(\tilde{x}^l) - F(\tilde{x}^l)\|^2 \right] \stackrel{(1)}{\leq} \frac{\sigma^2}{m_2},
\end{aligned}$$

for all  $l = 0, 1, \dots, T-1$ . Moreover, as we already derived, probability event  $E_{T-1}$  implies that  $\|F(x^l)\| \leq \lambda_1/2$  and  $\|F(\tilde{x}^l)\| \leq \lambda_1/2$  for all  $l = 0, 1, \dots, T-1$  (see (71) and (74)). Therefore, in view of Lemma B.2,  $E_{T-1}$  implies that

$$\|\theta_l^b\| \leq \frac{4\sigma^2}{m_2\lambda_2}, \quad \|\omega_l^b\| \leq \frac{4\sigma^2}{m_1\lambda_1}, \tag{82}$$

$$\mathbb{E}_{\xi_2^l} \left[ \|\theta_l\|^2 \right] \leq \frac{18\sigma^2}{m_2}, \quad \mathbb{E}_{\xi_1^l} \left[ \|\omega_l\|^2 \right] \leq \frac{18\sigma^2}{m_1}, \tag{83}$$

$$\mathbb{E}_{\xi_2^l} \left[ \|\theta_l^u\|^2 \right] \leq \frac{18\sigma^2}{m_2}, \quad \mathbb{E}_{\xi_1^l} \left[ \|\omega_l^u\|^2 \right] \leq \frac{18\sigma^2}{m_1}, \tag{84}$$

for all  $l = 0, 1, \dots, T-1$ .

**Upper bound for  $\textcircled{1}$ .** Since  $\mathbb{E}_{\xi_2^l} [\theta_l^u] = 0$ , we have

$$\mathbb{E}_{\xi_2^l} [2\gamma_2 \langle \eta_l, \theta_l^u \rangle] = 0.$$

Next, the summands in  $\textcircled{1}$  are bounded with probability 1:

$$|2\gamma_2 \langle \eta_l, \theta_l^u \rangle| \leq 2\gamma_2 \|\eta_l\| \cdot \|\theta_l^u\| \stackrel{(77),(81)}{\leq} 4\sqrt{7}\gamma_2(1 + \gamma_2 L)R\lambda_l \stackrel{(65),(66)}{\leq} \frac{R^2}{\ln \frac{6(K+1)}{\beta}} \stackrel{\text{def}}{=} c. \tag{85}$$

Moreover, these summands have bounded conditional variances  $\sigma_l^2 \stackrel{\text{def}}{=} \mathbb{E}_{\xi_2^l} [4\gamma_2^2 \langle \eta_l, \theta_l^u \rangle^2]$ :

$$\sigma_l^2 \leq \mathbb{E}_{\xi_2^l} [4\gamma_2^2 \|\eta_l\|^2 \cdot \|\theta_l^u\|^2] \stackrel{(77)}{\leq} 28\gamma_2^2(1 + \gamma_2 L)^2 R^2 \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2] \stackrel{(65)}{\leq} 30\gamma_2^2 R^2 \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2]. \tag{86}$$

That is, sequence  $\{2\gamma_2\langle\eta_l, \theta_l^u\rangle\}_{l \geq 0}$  is a bounded martingale difference sequence having bounded conditional variances  $\{\sigma_l^2\}_{l \geq 0}$ . Applying Bernstein's inequality (Lemma B.1) with  $X_l = 2\gamma_2\langle\eta_l, \theta_l^u\rangle$ ,  $c$  defined in (85),  $b = R^2$ ,  $G = \frac{R^4}{6 \ln \frac{6(K+1)}{\beta}}$ , we get that

$$\mathbb{P}\left\{|\textcircled{1}| > R^2 \text{ and } \sum_{l=0}^{T-1} \sigma_l^2 \leq \frac{R^4}{6 \ln \frac{6(K+1)}{\beta}}\right\} \leq 2 \exp\left(-\frac{b^2}{2G + 2cb/3}\right) = \frac{\beta}{3(K+1)}.$$

In other words,  $\mathbb{P}\{E_{\textcircled{1}}\} \geq 1 - \frac{\beta}{3(K+1)}$ , where probability event  $E_{\textcircled{1}}$  is defined as

$$E_{\textcircled{1}} = \left\{ \text{either } \sum_{l=0}^{T-1} \sigma_l^2 > \frac{R^4}{6 \ln \frac{6(K+1)}{\beta}} \text{ or } |\textcircled{1}| \leq R^2 \right\}. \quad (87)$$

Moreover, we notice here that probability event  $E_{T-1}$  implies that

$$\sum_{l=0}^{T-1} \sigma_l^2 \stackrel{(86)}{\leq} 30\gamma_2^2 R^2 \sum_{l=0}^{T-1} \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2] \stackrel{(84), T \leq K+1}{\leq} \frac{540\gamma_2^2 R^2 \sigma^2 (K+1)}{m_2} \stackrel{(68)}{\leq} \frac{R^4}{6 \ln \frac{6(K+1)}{\beta}}. \quad (88)$$

**Upper bound for ②.** Probability event  $E_{T-1}$  implies

$$\begin{aligned} \textcircled{2} &\leq 2\gamma_2 \sum_{l=0}^{T-1} \|\eta_l\| \cdot \|\theta_l^b\| \stackrel{(77), (82), T \leq K+1}{\leq} \frac{8\sqrt{7}\gamma_2(1 + \gamma_2 L)\sigma^2 R(K+1)}{m_2 \lambda_2} \\ &\stackrel{(65), (66)}{=} \frac{161\sqrt{7}\gamma_2^2 \sigma^2 (K+1) \ln \frac{6(K+1)}{\beta}}{m_2} \stackrel{(68)}{\leq} R^2. \end{aligned} \quad (89)$$

**Upper bound for ③.** Probability event  $E_{T-1}$  implies

$$\textcircled{3} = 2\gamma_2^2 \sum_{l=0}^{T-1} \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2] \stackrel{(84), T \leq K+1}{\leq} \frac{36\gamma_2^2 \sigma^2 (K+1)}{m_2} \stackrel{(68)}{\leq} R^2. \quad (90)$$

**Upper bound for ④.** We have

$$2\gamma_2^2 \mathbb{E}_{\xi_2^l} \left[ \|\theta_l^u\|^2 - \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2] \right] = 0.$$

Next, the summands in ④ are bounded with probability 1:

$$\begin{aligned} 2\gamma_2^2 \left| \|\theta_l^u\|^2 - \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2] \right| &\leq 2\gamma_2^2 \left( \|\theta_l^u\|^2 + \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2] \right) \stackrel{(81)}{\leq} 16\gamma_2^2 \lambda_2^2 \\ &\stackrel{(66)}{\leq} \frac{R^2}{\ln \frac{6(K+1)}{\beta}} \stackrel{\text{def}}{=} c. \end{aligned} \quad (91)$$

Moreover, these summands have bounded conditional variances  $\tilde{\sigma}_l^2 \stackrel{\text{def}}{=} 4\gamma_2^4 \mathbb{E}_{\xi_2^l} \left[ \left( \|\theta_l^u\|^2 - \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2] \right)^2 \right]$ :

$$\tilde{\sigma}_l^2 \stackrel{(91)}{\leq} \frac{2\gamma_2^2 R^2}{\ln \frac{6(K+1)}{\beta}} \mathbb{E}_{\xi_2^l} \left[ \left| \|\theta_l^u\|^2 - \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2] \right| \right] \leq \frac{4\gamma_2^2 R^2}{\ln \frac{6(K+1)}{\beta}} \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2] \quad (92)$$

That is, sequence  $\{\|\theta_l^u\|^2 - \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2]\}_{l \geq 0}$  is a bounded martingale difference sequence having bounded conditional variances  $\{\tilde{\sigma}_l^2\}_{l \geq 0}$ . Applying Bernstein's inequality (Lemma B.1) with  $X_l = \|\theta_l^u\|^2 - \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2]$ ,  $c$  defined in (91),  $b = R^2$ ,  $G = \frac{R^4}{6 \ln \frac{6(K+1)}{\beta}}$ , we get that

$$\mathbb{P}\left\{|\textcircled{4}| > R^2 \text{ and } \sum_{l=0}^{T-1} \tilde{\sigma}_l^2 \leq \frac{R^4}{6 \ln \frac{6(K+1)}{\beta}}\right\} \leq 2 \exp\left(-\frac{b^2}{2G + 2cb/3}\right) = \frac{\beta}{3(K+1)}.$$

In other words,  $\mathbb{P}\{E_{\textcircled{4}}\} \geq 1 - \frac{\beta}{3(K+1)}$ , where probability event  $E_{\textcircled{4}}$  is defined as

$$E_{\textcircled{4}} = \left\{ \text{either } \sum_{l=0}^{T-1} \tilde{\sigma}_l^2 > \frac{R^4}{6 \ln \frac{6(K+1)}{\beta}} \quad \text{or } |\textcircled{4}| \leq R^2 \right\}. \quad (93)$$

Moreover, we notice here that probability event  $E_{T-1}$  implies that

$$\begin{aligned} \sum_{l=0}^{T-1} \tilde{\sigma}_l^2 &\stackrel{(92)}{\leq} \frac{4\gamma_2^2 R^2}{\ln \frac{6(K+1)}{\beta}} \sum_{l=0}^{T-1} \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2] \stackrel{(84), T \leq K+1}{\leq} \frac{72\gamma_2^2 R^2 \sigma^2 (K+1)}{m_2 \ln \frac{6(K+1)}{\beta}} \\ &\stackrel{(68)}{\leq} \frac{R^4}{6 \ln \frac{6(K+1)}{\beta}}. \end{aligned} \quad (94)$$

**Upper bound for ⑤.** Probability event  $E_{T-1}$  implies

$$\textcircled{5} = 2\gamma_2^2 \sum_{l=0}^{T-1} \|\theta_l^b\|^2 \stackrel{(82), T \leq K+1}{\leq} \frac{32\gamma_2^2 \sigma^4 (K+1)}{m_2^2 \lambda_2^2} \stackrel{(66)}{\leq} \frac{12800\gamma_2^4 \sigma^4 (K+1) \ln^2 \frac{6(K+1)}{\beta}}{m_2^2 R^2} \stackrel{(68)}{\leq} R^2. \quad (95)$$

**Upper bound for ⑥.** Probability event  $E_{T-1}$  implies

$$\textcircled{6} = 6\gamma_1\gamma_2 \sum_{l=0}^{T-1} \mathbb{E}_{\xi_1^l} [\|\omega_l^u\|^2] \stackrel{(84), T \leq K+1}{\leq} \frac{108\gamma_1\gamma_2 \sigma^2 (K+1)}{m_1} \stackrel{(67)}{\leq} R^2. \quad (96)$$

**Upper bound for ⑦.** We have

$$6\gamma_1\gamma_2 \mathbb{E}_{\xi_1^l} [\|\omega_l^u\|^2 - \mathbb{E}_{\xi_1^l} [\|\omega_l^u\|^2]] = 0.$$

Next, the summands in ⑦ are bounded with probability 1:

$$\begin{aligned} 6\gamma_1\gamma_2 \left| \|\omega_l^u\|^2 - \mathbb{E}_{\xi_1^l} [\|\omega_l^u\|^2] \right| &\leq 6\gamma_1\gamma_2 \left( \|\omega_l^u\|^2 + \mathbb{E}_{\xi_1^l} [\|\omega_l^u\|^2] \right) \stackrel{(81)}{\leq} 48\gamma_1\gamma_2 \lambda_2^2 \\ &\stackrel{(66)}{\leq} \frac{\gamma_2 R^2}{\gamma_1 \ln \frac{6(K+1)}{\beta}} \stackrel{\gamma_2 \leq \gamma_1}{\leq} \frac{R^2}{\ln \frac{6(K+1)}{\beta}} \stackrel{\text{def}}{=} c. \end{aligned} \quad (97)$$

Moreover, these summands have bounded conditional variances  $\hat{\sigma}_l^2 \stackrel{\text{def}}{=} 36\gamma_1^2\gamma_2^2 \mathbb{E}_{\xi_1^l} \left[ \left( \|\omega_l^u\|^2 - \mathbb{E}_{\xi_1^l} [\|\omega_l^u\|^2] \right)^2 \right]$ :

$$\hat{\sigma}_l^2 \stackrel{(97)}{\leq} \frac{6\gamma_2^2 R^2}{\ln \frac{6(K+1)}{\beta}} \mathbb{E}_{\xi_1^l} \left[ \left| \|\omega_l^u\|^2 - \mathbb{E}_{\xi_1^l} [\|\omega_l^u\|^2] \right| \right] \leq \frac{12\gamma_2^2 R^2}{\ln \frac{6(K+1)}{\beta}} \mathbb{E}_{\xi_1^l} [\|\omega_l^u\|^2] \quad (98)$$

That is, sequence  $\{\|\omega_l^u\|^2 - \mathbb{E}_{\xi_1^l} [\|\omega_l^u\|^2]\}_{l \geq 0}$  is a bounded martingale difference sequence having bounded conditional variances  $\{\hat{\sigma}_l^2\}_{l \geq 0}$ . Applying Bernstein's inequality (Lemma B.1) with  $X_l = \|\omega_l^u\|^2 - \mathbb{E}_{\xi_1^l} [\|\omega_l^u\|^2]$ ,  $c$  defined in (97),  $b = R^2$ ,  $G = \frac{R^4}{6 \ln \frac{6(K+1)}{\beta}}$ , we get that

$$\mathbb{P} \left\{ |\textcircled{7}| > R^2 \text{ and } \sum_{l=0}^{T-1} \hat{\sigma}_l^2 \leq \frac{R^4}{6 \ln \frac{6(K+1)}{\beta}} \right\} \leq 2 \exp \left( -\frac{b^2}{2G + 2cb/3} \right) = \frac{\beta}{3(K+1)}.$$

In other words,  $\mathbb{P}\{E_{\textcircled{7}}\} \geq 1 - \frac{\beta}{3(K+1)}$ , where probability event  $E_{\textcircled{7}}$  is defined as

$$E_{\textcircled{7}} = \left\{ \text{either } \sum_{l=0}^{T-1} \hat{\sigma}_l^2 > \frac{R^4}{6 \ln \frac{6(K+1)}{\beta}} \quad \text{or } |\textcircled{7}| \leq R^2 \right\}. \quad (99)$$

Moreover, we notice here that probability event  $E_{T-1}$  implies that

$$\begin{aligned} \sum_{l=0}^{T-1} \hat{\sigma}_l^2 &\stackrel{(98)}{\leq} \frac{12\gamma_2^2 R^2}{\ln \frac{6(K+1)}{\beta}} \sum_{l=0}^{T-1} \mathbb{E}_{\xi_1^l} [\|\omega_l^u\|^2] \stackrel{(84), T \leq K+1}{\leq} \frac{216\gamma_2^2 R^2 \sigma^2 (K+1)}{m_1 \ln \frac{6(K+1)}{\beta}} \\ &\stackrel{(67)}{\leq} \frac{R^4}{6 \ln \frac{6(K+1)}{\beta}}. \end{aligned} \quad (100)$$

**Upper bound for ⑧.** Probability event  $E_{T-1}$  implies

$$\begin{aligned} \textcircled{8} &= 6\gamma_1\gamma_2 \sum_{l=0}^{T-1} \|\omega_l^b\|^2 \stackrel{(82), T \leq K+1}{\leq} \frac{96\gamma_1\gamma_2\sigma^4(K+1)}{m_1^2\lambda_1^2} \\ &\stackrel{(66)}{\leq} \frac{38400\gamma_1^3\gamma_2\sigma^4(K+1)\ln^2 \frac{6(K+1)}{\beta}}{m_1^2 R^2} \stackrel{(67)}{\leq} R^2. \end{aligned} \quad (101)$$

**Final derivation.** Putting all bounds together, we get that  $E_{T-1}$  implies

$$\begin{aligned} R_T^2 &\stackrel{(80)}{\leq} R^2 + \textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4} + \textcircled{5} + \textcircled{6} + \textcircled{7} + \textcircled{8}, \\ \textcircled{2} &\stackrel{(89)}{\leq} R^2, \quad \textcircled{3} \stackrel{(90)}{\leq} R^2, \quad \textcircled{5} \stackrel{(95)}{\leq} R^2, \quad \textcircled{6} \stackrel{(96)}{\leq} R^2, \quad \textcircled{8} \stackrel{(101)}{\leq} R^2, \\ \sum_{l=0}^{T-1} \sigma_l^2 &\stackrel{(88)}{\leq} \frac{R^4}{6 \ln \frac{6(K+1)}{\beta}}, \quad \sum_{l=0}^{T-1} \tilde{\sigma}_l^2 \stackrel{(94)}{\leq} \frac{R^4}{6 \ln \frac{6(K+1)}{\beta}}, \quad \sum_{l=0}^{T-1} \hat{\sigma}_l^2 \stackrel{(100)}{\leq} \frac{R^4}{6 \ln \frac{6(K+1)}{\beta}}. \end{aligned}$$

Moreover, in view of (87), (93), (99), and our induction assumption, we have

$$\begin{aligned} \mathbb{P}\{E_{T-1}\} &\geq 1 - \frac{(T-1)\beta}{K+1}, \\ \mathbb{P}\{E_{\textcircled{1}}\} &\geq 1 - \frac{\beta}{3(K+1)}, \quad \mathbb{P}\{E_{\textcircled{4}}\} \geq 1 - \frac{\beta}{3(K+1)}, \quad \mathbb{P}\{E_{\textcircled{7}}\} \geq 1 - \frac{\beta}{3(K+1)}, \end{aligned}$$

where probability events  $E_{\textcircled{1}}$ ,  $E_{\textcircled{4}}$ , and  $E_{\textcircled{7}}$  are defined as

$$\begin{aligned} E_{\textcircled{1}} &= \left\{ \text{either } \sum_{l=0}^{T-1} \sigma_l^2 > \frac{R^4}{6 \ln \frac{6(K+1)}{\beta}} \text{ or } |\textcircled{1}| \leq R^2 \right\}, \\ E_{\textcircled{4}} &= \left\{ \text{either } \sum_{l=0}^{T-1} \tilde{\sigma}_l^2 > \frac{R^4}{6 \ln \frac{6(K+1)}{\beta}} \text{ or } |\textcircled{4}| \leq R^2 \right\}, \\ E_{\textcircled{7}} &= \left\{ \text{either } \sum_{l=0}^{T-1} \hat{\sigma}_l^2 > \frac{R^4}{6 \ln \frac{6(K+1)}{\beta}} \text{ or } |\textcircled{7}| \leq R^2 \right\}. \end{aligned}$$

Putting all of these inequalities together, we obtain that probability event  $E_{T-1} \cap E_{\textcircled{1}} \cap E_{\textcircled{4}} \cap E_{\textcircled{7}}$  implies

$$\begin{aligned} R_T^2 &\stackrel{(80)}{\leq} R^2 + \textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4} + \textcircled{5} + \textcircled{6} + \textcircled{7} + \textcircled{8} \\ &\leq 9R^2. \end{aligned}$$

Moreover, union bound for the probability events implies

$$\mathbb{P}\{E_T\} \geq \mathbb{P}\{E_{T-1} \cap E_{\textcircled{1}} \cap E_{\textcircled{4}} \cap E_{\textcircled{7}}\} = 1 - \mathbb{P}\{\bar{E}_{T-1} \cup \bar{E}_{\textcircled{1}} \cup \bar{E}_{\textcircled{4}} \cup \bar{E}_{\textcircled{7}}\} \geq 1 - \frac{T\beta}{K+1}. \quad (102)$$

This is exactly what we wanted to prove (see the paragraph after inequality (70)). In particular,  $E_{K+1}$  implies

$$\begin{aligned} \frac{1}{K+1} \sum_{k=0}^{K+1} \|F(x^k)\|^2 &\stackrel{(75)}{\leq} \frac{4(R^2 - R_{K+1}^2)}{\gamma_1\gamma_2(K+1)} + \frac{4(\textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4} + \textcircled{5} + \textcircled{6} + \textcircled{7} + \textcircled{8})}{\gamma_1\gamma_2(K+1)} \\ &\leq \frac{36R^2}{\gamma_1\gamma_2(K+1)}. \end{aligned}$$

This finishes the proof.  $\square$

**Corollary C.2.** *Let the assumptions of Theorem C.2 hold and*

$$\rho \leq \frac{1}{640L \ln \frac{6(K+1)}{\beta}}. \quad (103)$$



Then, the choice of step-sizes and batch-sizes

$$2\gamma_2 = \gamma_1 = \frac{1}{160L \ln \frac{6(K+1)}{\beta}}, \quad m_1 = m_2 = \max \left\{ 1, \frac{81(K+1)\sigma^2}{640L^2R^2 \ln \frac{6(K+1)}{\beta}} \right\} \quad (104)$$

satisfies conditions (65), (67), (68). With such choice of  $\gamma, m_1, m_2$ , and the choice of  $\lambda_1, \lambda_2$  as in (66), the iterates produced by clipped-SEG after  $K$  iterations with probability at least  $1 - \beta$  satisfy

$$\frac{1}{K+1} \sum_{k=0}^K \|F(x^k)\|^2 \leq \frac{1843200L^2R^2 \ln^2 \frac{6(K+1)}{\beta}}{K+1}. \quad (105)$$

In particular, to guarantee  $\frac{1}{K+1} \sum_{k=0}^K \|F(x^k)\|^2 \leq \varepsilon$  with probability at least  $1 - \beta$  for some  $\varepsilon > 0$  clipped-SEG requires,

$$\mathcal{O} \left( \frac{L^2R^2}{\varepsilon} \ln^2 \left( \frac{L^2R^2}{\varepsilon\beta} \right) \right) \text{ iterations}, \quad (106)$$

$$\mathcal{O} \left( \max \left\{ \frac{L^2R^2}{\varepsilon} \ln^2 \left( \frac{L^2R^2}{\varepsilon\beta} \right), \frac{L^2\sigma^2R^2}{\varepsilon^2} \ln^3 \left( \frac{LR^2}{\varepsilon\beta} \right) \right\} \right) \text{ oracle calls}. \quad (107)$$

*Proof.* First of all, we verify that the choice of  $\gamma_1, \gamma_2$  and  $m_1, m_2$  from (104) satisfies conditions (65), (67), (68). Inequality (65) holds since

$$\gamma_2 + 2\rho \stackrel{(104)}{=} \frac{1}{320L \ln \frac{6(K+1)}{\beta}} + 2\rho \stackrel{(103)}{\leq} \frac{1}{320L \ln \frac{6(K+1)}{\beta}} + \frac{1}{320L \ln \frac{6(K+1)}{\beta}} \stackrel{(104)}{=} \gamma_1$$

and (67), (68) are satisfied since

$$\begin{aligned} m_1 &= \max \left\{ 1, \frac{81(K+1)\sigma^2}{640L^2R^2 \ln \frac{6(K+1)}{\beta}} \right\} \\ &\geq \max \left\{ 1, \frac{216 \max\{\gamma_1\gamma_2(K+1), \sqrt{\gamma_1^3\gamma_2(K+1) \ln \frac{6(K+1)}{\beta}}\} \sigma^2}{R^2} \right\}, \\ m_2 &= \max \left\{ 1, \frac{81(K+1)\sigma^2}{640L^2R^2 \ln \frac{6(K+1)}{\beta}} \right\} \geq \max \left\{ 1, \frac{3240(K+1)\gamma_2^2\sigma^2 \ln \frac{6(K+1)}{\beta}}{R^2} \right\}. \end{aligned}$$

Therefore, applying Theorem C.2, we derive that with probability at least  $1 - \beta$

$$\frac{1}{K+1} \sum_{k=0}^K \|F(x^k)\|^2 \leq \frac{36R^2}{\gamma_1\gamma_2(K+1)} \stackrel{(104)}{=} \frac{1843200L^2R^2 \ln^2 \frac{6(K+1)}{\beta}}{K+1}.$$

To guarantee  $\frac{1}{K+1} \sum_{k=0}^K \|F(x^k)\|^2 \leq \varepsilon$ , we choose  $K$  in such a way that the right-hand side of the above inequality is smaller than  $\varepsilon$  that gives

$$K = \mathcal{O} \left( \frac{L^2R^2}{\varepsilon} \ln^2 \left( \frac{L^2R^2}{\varepsilon\beta} \right) \right).$$

The total number of oracle calls equals

$$\begin{aligned} 2m(K+1) &\stackrel{(104)}{=} 2 \max \left\{ K+1, \frac{81(K+1)^2\sigma^2}{640L^2R^2 \ln \frac{6(K+1)}{\beta}} \right\} \\ &= \mathcal{O} \left( \max \left\{ \frac{L^2R^2}{\varepsilon} \ln^2 \left( \frac{L^2R^2}{\varepsilon\beta} \right), \frac{L^2\sigma^2R^2}{\varepsilon^2} \ln^3 \left( \frac{L^2R^2}{\varepsilon\beta} \right) \right\} \right). \end{aligned}$$

□

### C.3 Quasi-Strongly Monotone Case

**Lemma C.3.** *Let Assumptions 1.2, 1.5 hold for  $Q = B_{3R}(x^*) = \{x \in \mathbb{R}^d \mid \|x - x^*\| \leq 3R\}$ , where  $R \geq R_0 \stackrel{\text{def}}{=} \|x^0 - x^*\|$ , and  $\gamma_1 = \gamma_2 = \gamma$ ,  $0 < \gamma \leq 1/2(L+2\mu)$ . If  $x^k$  and  $\tilde{x}^k \stackrel{\text{def}}{=} x^k - \gamma F(x^k)$  lie in  $B_{3R}(x^*)$  for all  $k = 0, 1, \dots, K$  for some  $K \geq 0$ , then the iterates produced by **clipped-SEG** satisfy*

$$\begin{aligned} \|x^{K+1} - x^*\|^2 &\leq (1 - \gamma\mu)^{K+1} \|x^0 - x^*\|^2 - 4\gamma^3\mu \sum_{k=0}^K (1 - \gamma\mu)^{K-k} \langle F(x^k), \omega_k \rangle \\ &\quad + 2\gamma \sum_{k=0}^K (1 - \gamma\mu)^{K-k} \langle x^k - x^* - \gamma F(\tilde{x}^k), \theta_k \rangle \\ &\quad + \gamma^2 \sum_{k=0}^K (1 - \gamma\mu)^{K-k} (\|\theta_k\|^2 + 4\|\omega_k\|^2), \end{aligned} \quad (108)$$

where  $\theta_k, \omega_k$  are defined in (15), (16).

*Proof.* Using the update rule of **clipped-SEG**, we obtain

$$\begin{aligned} \|x^{k+1} - x^*\|^2 &= \|x^k - x^*\|^2 - 2\gamma \langle x^k - x^*, \tilde{F}_{\xi_2^k}(\tilde{x}^k) \rangle + \gamma^2 \|\tilde{F}_{\xi_2^k}(\tilde{x}^k)\|^2 \\ &= \|x^k - x^*\|^2 - 2\gamma \langle x^k - x^*, F(\tilde{x}^k) \rangle + 2\gamma \langle x^k - x^*, \theta_k \rangle \\ &\quad + \gamma^2 \|F(\tilde{x}^k)\|^2 - 2\gamma^2 \langle F(\tilde{x}^k), \theta_k \rangle + \gamma^2 \|\theta_k\|^2 \\ &= \|x^k - x^*\|^2 - 2\gamma \langle \tilde{x}^k - x^*, F(\tilde{x}^k) \rangle - 2\gamma \langle x^k - \tilde{x}^k, F(\tilde{x}^k) \rangle \\ &\quad + 2\gamma \langle x^k - x^* - \gamma F(\tilde{x}^k), \theta_k \rangle + \gamma^2 \|F(\tilde{x}^k)\|^2 + \gamma^2 \|\theta_k\|^2. \end{aligned}$$

Since  $F$  is  $\mu$ -quasi strongly monotone, we have

$$\begin{aligned} -2\gamma \langle \tilde{x}^k - x^*, F(\tilde{x}^k) \rangle &\leq -2\gamma\mu \|\tilde{x}^k - x^*\|^2 \stackrel{(6)}{\leq} -\gamma\mu \|x^k - x^*\|^2 + 2\gamma\mu \|\tilde{x}^k - x^k\|^2 \\ &= -\gamma\mu \|x^k - x^*\|^2 + 2\gamma^3\mu \|\tilde{F}_{\xi_1}(x^k)\|^2 \\ &= -\gamma\mu \|x^k - x^*\|^2 + 2\gamma^3\mu \|F(x^k)\|^2 - 4\gamma^3\mu \langle F(x^k), \omega_k \rangle + 2\gamma^3\mu \|\omega_k\|^2. \end{aligned}$$

Moreover,  $-2\gamma \langle x^k - \tilde{x}^k, F(\tilde{x}^k) \rangle$  can be rewritten as

$$\begin{aligned} -2\gamma \langle x^k - \tilde{x}^k, F(\tilde{x}^k) \rangle &= -2\gamma^2 \langle \tilde{F}_{\xi_1}(x^k), F(\tilde{x}^k) \rangle \\ &= \gamma^2 \|\tilde{F}_{\xi_1}(x^k) - F(x^k)\|^2 - \gamma^2 \|\tilde{F}_{\xi_1}(x^k)\|^2 - \gamma^2 \|F(\tilde{x}^k)\|^2. \end{aligned}$$

Putting all together, we get

$$\begin{aligned}
\|x^{k+1} - x^*\|^2 &\leq (1 - \gamma\mu)\|x^k - x^*\|^2 + 2\gamma^3\mu\|F(x^k)\|^2 - 4\gamma^3\mu\langle F(x^k), \omega_k \rangle + 2\gamma^3\mu\|\omega_k\|^2 \\
&\quad + \gamma^2\|\tilde{F}_{\xi_1}(x^k) - F(x^k)\|^2 - \gamma^2\|\tilde{F}_{\xi_1}(x^k)\|^2 - \gamma^2\|F(\tilde{x}^k)\|^2 \\
&\quad + 2\gamma\langle x^k - x^* - \gamma F(\tilde{x}^k), \theta_k \rangle + \gamma^2\|F(\tilde{x}^k)\|^2 + \gamma^2\|\theta_k\|^2 \\
&\stackrel{(5)}{\leq} (1 - \gamma\mu)\|x^k - x^*\|^2 + 2\gamma^3\mu\|F(x^k)\|^2 - 4\gamma^3\mu\langle F(x^k), \omega_k \rangle + 2\gamma^3\mu\|\omega_k\|^2 \\
&\quad + 2\gamma^2\|\omega_k\|^2 + 2\gamma^2\|F(x^k) - F(\tilde{x}^k)\|^2 - \gamma^2\|\tilde{F}_{\xi_1}(x^k)\|^2 \\
&\quad + 2\gamma\langle x^k - x^* - \gamma F(\tilde{x}^k), \theta_k \rangle + \gamma^2\|\theta_k\|^2 \\
&\stackrel{(\text{Lip})}{\leq} (1 - \gamma\mu)\|x^k - x^*\|^2 + 2\gamma^3\mu\|F(x^k)\|^2 - 4\gamma^3\mu\langle F(x^k), \omega_k \rangle \\
&\quad + 2\gamma^2(1 + \gamma\mu)\|\omega_k\|^2 - \gamma^2(1 - 2\gamma^2L^2)\|\tilde{F}_{\xi_1}(x^k)\|^2 \\
&\quad + 2\gamma\langle x^k - x^* - \gamma F(\tilde{x}^k), \theta_k \rangle + \gamma^2\|\theta_k\|^2 \\
&\stackrel{(6)}{\leq} (1 - \gamma\mu)\|x^k - x^*\|^2 - \gamma^2\left(\frac{1}{2} - \gamma^2L^2 - 2\gamma\mu\right)\|F(x^k)\|^2 \\
&\quad - 4\gamma^3\mu\langle F(x^k), \omega_k \rangle + \gamma^2(3 - 2\gamma^2L^2 + 2\gamma\mu)\|\omega_k\|^2 \\
&\quad + 2\gamma\langle x^k - x^* - \gamma F(\tilde{x}^k), \theta_k \rangle + \gamma^2\|\theta_k\|^2 \\
&\leq (1 - \gamma\mu)\|x^k - x^*\|^2 - 4\gamma^3\mu\langle F(x^k), \omega_k \rangle \\
&\quad + 2\gamma\langle x^k - x^* - \gamma F(\tilde{x}^k), \theta_k \rangle + \gamma^2(\|\theta_k\|^2 + 4\|\omega_k\|^2),
\end{aligned}$$

where in the last step we apply  $0 < \gamma \leq 1/2(L+2\mu)$ . Unrolling the recurrence, we obtain (108).  $\square$

**Theorem C.3.** *Let Assumptions 1.1, 1.2, 1.5, hold for  $Q = B_{3R}(x^*) = \{x \in \mathbb{R}^d \mid \|x - x^*\| \leq 3R\}$ , where  $R \geq R_0 \stackrel{\text{def}}{=} \|x^0 - x^*\|$ , and  $\gamma_1 = \gamma_2 = \gamma$ ,*

$$0 < \gamma \leq \frac{1}{650L \ln \frac{6(K+1)}{\beta}}, \quad (109)$$

$$\lambda_{1,k} = \lambda_{2,k} = \lambda_k = \frac{\exp(-\gamma\mu(1 + k/2))R}{120\gamma \ln \frac{6(K+1)}{\beta}}, \quad (110)$$

$$m_{1,k} = m_{2,k} = m_k \geq \max\left\{1, \frac{264600\gamma^2(K+1)\sigma^2 \ln \frac{6(K+1)}{\beta}}{\exp(-\gamma\mu k)R^2}\right\}, \quad (111)$$

for some  $K \geq 0$  and  $\beta \in (0, 1]$  such that  $\ln \frac{6(K+1)}{\beta} \geq 1$ . Then, after  $K$  iterations the iterates produced by clipped-SEG with probability at least  $1 - \beta$  satisfy

$$\|x^{K+1} - x^*\|^2 \leq 2 \exp(-\gamma\mu(K+1))R^2. \quad (112)$$

*Proof.* As in the proof of Theorem C.1, we use the following notation:  $R_k = \|x^k - x^*\|^2$ ,  $k \geq 0$ . We will derive (112) by induction. In particular, for each  $k = 0, \dots, K+1$  we define probability event  $E_k$  as follows: inequalities

$$R_t^2 \leq 2 \exp(-\gamma\mu t)R^2 \quad (113)$$

hold for  $t = 0, 1, \dots, k$  simultaneously. Our goal is to prove that  $\mathbb{P}\{E_k\} \geq 1 - k\beta/(K+1)$  for all  $k = 0, 1, \dots, K+1$ . We use the induction to show this statement. For  $k = 0$  the statement is trivial since  $R_0^2 \leq 2R^2$  by definition. Next, assume that the statement holds for  $k = T-1 \leq K$ , i.e., we have  $\mathbb{P}\{E_{T-1}\} \geq 1 - (T-1)\beta/(K+1)$ . We need to prove that  $\mathbb{P}\{E_T\} \geq 1 - T\beta/(K+1)$ . First of all, since  $R_t^2 \leq 2 \exp(-\gamma\mu t)R^2 \leq 9R^2$ , we have  $x^t \in B_{3R}(x^*)$ . Operator  $F$  is  $L$ -Lipschitz on  $B_{3R}(x^*)$ . Therefore, probability event  $E_{T-1}$  implies

$$\|F(x^t)\| \leq L\|x^t - x^*\| \stackrel{(113)}{\leq} \sqrt{2}L \exp(-\gamma\mu t/2)R \stackrel{(109),(110)}{\leq} \frac{\lambda_t}{2}. \quad (114)$$

and

$$\|\omega_t\|^2 \stackrel{(5)}{\leq} 2\|\tilde{F}_{\xi_1}(x^t)\|^2 + 2\|F(x^t)\|^2 \stackrel{(114)}{\leq} \frac{5}{2}\lambda_t^2 \stackrel{(110)}{\leq} \frac{\exp(-\gamma\mu t)R^2}{4\gamma^2} \quad (115)$$

for all  $t = 0, 1, \dots, T-1$ .

Next, we show that probability event  $E_{T-1}$  implies  $\|\tilde{x}^t - x^*\| \leq 3R$  and derive useful inequalities related to  $\theta_t$  for all  $t = 0, 1, \dots, T-1$ . Indeed, due to Lipschitzness of  $F$  probability event  $E_{T-1}$  implies

$$\begin{aligned} \|\tilde{x}^t - x^*\|^2 &= \|x^t - x^* - \gamma \tilde{F}_{\xi_1}(x^t)\|^2 \stackrel{(5)}{\leq} 2\|x^t - x^*\|^2 + 2\gamma^2 \|\tilde{F}_{\xi_1}(x^t)\|^2 \\ &\stackrel{(5)}{\leq} 2R_t^2 + 4\gamma^2 \|F(x^t)\|^2 + 4\gamma^2 \|\omega_t\|^2 \\ &\stackrel{(\text{Lip})}{\leq} 2(1 + 2\gamma^2 L^2)R_t^2 + 4\gamma^2 \|\omega_t\|^2 \\ &\stackrel{(109),(115)}{\leq} 7 \exp(-\gamma\mu t) R^2 \leq 9R^2 \end{aligned} \quad (116)$$

and

$$\|F(\tilde{x}^t)\| \leq L\|\tilde{x}^t - x^*\| \leq \sqrt{7}L \exp(-\gamma\mu t/2) R \stackrel{(109),(110)}{\leq} \frac{\lambda_t}{2} \quad (117)$$

for all  $t = 0, 1, \dots, T-1$ .

That is,  $E_{T-1}$  implies that  $x^t, \tilde{x}^t \in B_{3R}(x^*)$  for all  $t = 0, 1, \dots, T-1$ . Applying Lemma C.3 and  $(1 - \gamma\mu)^T \leq \exp(-\gamma\mu T)$ , we get that probability event  $E_{T-1}$  implies

$$\begin{aligned} R_T^2 &\leq \exp(-\gamma\mu T) R^2 - 4\gamma^3 \mu \sum_{l=0}^{T-1} (1 - \gamma\mu)^{T-1-l} \langle F(x^l), \omega_l \rangle \\ &\quad + 2\gamma \sum_{l=0}^{T-1} (1 - \gamma\mu)^{T-1-l} \langle x^l - x^* - \gamma F(\tilde{x}^l), \theta_l \rangle \\ &\quad + \gamma^2 \sum_{l=0}^{T-1} (1 - \gamma\mu)^{T-1-l} (\|\theta_l\|^2 + 4\|\omega_l\|^2). \end{aligned}$$

To estimate the sums in the right-hand side, we introduce new vectors:

$$\zeta_t = \begin{cases} F(x^t), & \text{if } \|F(x^t)\| \leq \sqrt{2}L \exp(-\gamma\mu t/2) R, \\ 0, & \text{otherwise,} \end{cases} \quad (118)$$

$$\eta_t = \begin{cases} x^t - x^* - \gamma F(\tilde{x}^t), & \text{if } \|x^t - x^* - \gamma F(\tilde{x}^t)\| \leq \sqrt{7}(1 + \gamma L) \exp(-\gamma\mu t/2) R, \\ 0, & \text{otherwise,} \end{cases} \quad (119)$$

for  $t = 0, 1, \dots, T-1$ . First of all, we point out that vectors  $\zeta_t$  and  $\eta_t$  are bounded with probability 1, i.e., with probability 1

$$\|\zeta_t\| \leq \sqrt{2}L \exp(-\gamma\mu t/2) R, \quad \|\eta_t\| \leq \sqrt{7}(1 + \gamma L) \exp(-\gamma\mu t/2) R \quad (120)$$

for all  $t = 0, 1, \dots, T-1$ . Next, we notice that  $E_{T-1}$  implies  $\|F(x^t)\| \leq \sqrt{2}L \exp(-\gamma\mu t/2) R$  (due to (114)) and

$$\begin{aligned} \|x^t - x^* - \gamma F(\tilde{x}^t)\| &\leq \|x^t - x^*\| + \gamma \|F(\tilde{x}^t)\| \\ &\stackrel{(116),(117)}{\leq} \sqrt{7}(1 + \gamma L) \exp(-\gamma\mu t/2) R \end{aligned}$$

for  $t = 0, 1, \dots, T-1$ , i.e., probability event  $E_{T-1}$  implies  $\zeta_t = F(x^t)$  and  $\eta_t = x^t - x^* - \gamma F(\tilde{x}^t)$  for all  $t = 0, 1, \dots, T-1$ . Therefore,  $E_{T-1}$  implies

$$\begin{aligned} R_T^2 &\leq \exp(-\gamma\mu T) R^2 - 4\gamma^3 \mu \sum_{l=0}^{T-1} (1 - \gamma\mu)^{T-1-l} \langle \zeta_l, \omega_l \rangle \\ &\quad + 2\gamma \sum_{l=0}^{T-1} (1 - \gamma\mu)^{T-1-l} \langle \eta_l, \theta_l \rangle + \gamma^2 \sum_{l=0}^{T-1} (1 - \gamma\mu)^{T-1-l} (\|\theta_l\|^2 + 4\|\omega_l\|^2). \end{aligned}$$

As in the monotone case, to continue the derivation, we introduce vectors  $\theta_l^u, \theta_l^b, \omega_l^u, \omega_l^b$  defined as

$$\theta_l^u \stackrel{\text{def}}{=} \mathbb{E}_{\xi_2^l} \left[ \tilde{F}_{\xi_2^l}(\tilde{x}^l) \right] - \tilde{F}_{\xi_2^l}(\tilde{x}^l), \quad \theta_l^b \stackrel{\text{def}}{=} F(\tilde{x}^l) - \mathbb{E}_{\xi_2^l} \left[ \tilde{F}_{\xi_2^l}(\tilde{x}^l) \right], \quad (121)$$

$$\omega_l^u \stackrel{\text{def}}{=} \mathbb{E}_{\xi_1^l} \left[ \tilde{F}_{\xi_1^l}(x^l) \right] - \tilde{F}_{\xi_1^l}(x^l), \quad \omega_l^b \stackrel{\text{def}}{=} F(x^l) - \mathbb{E}_{\xi_1^l} \left[ \tilde{F}_{\xi_1^l}(x^l) \right], \quad (122)$$

for all  $l = 0, \dots, T-1$ . By definition we have  $\theta_l = \theta_l^u + \theta_l^b, \omega_l = \omega_l^u + \omega_l^b$  for all  $l = 0, \dots, T-1$ . Using the introduced notation, we continue our derivation as follows:  $E_{T-1}$  implies

$$\begin{aligned} R_T^2 &\stackrel{(5)}{\leq} \underbrace{\exp(-\gamma\mu T)R^2 - 4\gamma^3\mu \sum_{l=0}^{T-1} (1-\gamma\mu)^{T-1-l} \langle \zeta_l, \omega_l^u \rangle - 4\gamma^3\mu \sum_{l=0}^{T-1} (1-\gamma\mu)^{T-1-l} \langle \zeta_l, \omega_l^b \rangle}_{\textcircled{1}} \\ &\quad + \underbrace{2\gamma \sum_{l=0}^{T-1} (1-\gamma\mu)^{T-1-l} \langle \eta_l, \theta_l^u \rangle + 2\gamma \sum_{l=0}^{T-1} (1-\gamma\mu)^{T-1-l} \langle \eta_l, \theta_l^b \rangle}_{\textcircled{2}} \\ &\quad + \underbrace{2\gamma^2 \sum_{l=0}^{T-1} (1-\gamma\mu)^{T-1-l} \left( \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2] + 4\mathbb{E}_{\xi_1^l} [\|\omega_l^u\|^2] \right)}_{\textcircled{3}} \\ &\quad + \underbrace{2\gamma^2 \sum_{l=0}^{T-1} (1-\gamma\mu)^{T-1-l} \left( \|\theta_l^u\|^2 + 4\|\omega_l^u\|^2 - \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2] - 4\mathbb{E}_{\xi_1^l} [\|\omega_l^u\|^2] \right)}_{\textcircled{4}} \\ &\quad + \underbrace{2\gamma^2 \sum_{l=0}^{T-1} (1-\gamma\mu)^{T-1-l} \left( \|\theta_l^b\|^2 + 4\|\omega_l^b\|^2 \right)}_{\textcircled{5}}. \end{aligned} \quad (123)$$

The rest of the proof is based on deriving good enough upper bounds for  $\textcircled{1}, \textcircled{2}, \textcircled{3}, \textcircled{4}, \textcircled{5}, \textcircled{6}, \textcircled{7}$ , i.e., we want to prove that  $\textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4} + \textcircled{5} + \textcircled{6} + \textcircled{7} \leq \exp(-\gamma\mu T)R^2$  with high probability.

Before we move on, we need to derive some useful inequalities for operating with  $\theta_l^u, \theta_l^b, \omega_l^u, \omega_l^b$ . First of all, Lemma B.2 implies that

$$\|\theta_l^u\| \leq 2\lambda_l, \quad \|\omega_l^u\| \leq 2\lambda_l \quad (124)$$

for all  $l = 0, 1, \dots, T-1$ . Next, since  $\{\xi_1^{i,l}\}_{i=1}^{m_l}, \{\xi_2^{i,l}\}_{i=1}^{m_l}$  are independently sampled from  $\mathcal{D}$ , we have  $\mathbb{E}_{\xi_1^l} [F_{\xi_1^l}(x^l)] = F(x^l), \mathbb{E}_{\xi_2^l} [F_{\xi_2^l}(\tilde{x}^l)] = F(\tilde{x}^l)$ , and

$$\begin{aligned} \mathbb{E}_{\xi_1^l} \left[ \|F_{\xi_1^l}(x^l) - F(x^l)\|^2 \right] &= \frac{1}{m_l^2} \sum_{i=1}^{m_l} \mathbb{E}_{\xi_1^{i,l}} \left[ \|F_{\xi_1^{i,l}}(x^l) - F(x^l)\|^2 \right] \stackrel{(1)}{\leq} \frac{\sigma^2}{m_l}, \\ \mathbb{E}_{\xi_2^l} \left[ \|F_{\xi_2^l}(\tilde{x}^l) - F(\tilde{x}^l)\|^2 \right] &= \frac{1}{m_l^2} \sum_{i=1}^{m_l} \mathbb{E}_{\xi_2^{i,l}} \left[ \|F_{\xi_2^{i,l}}(\tilde{x}^l) - F(\tilde{x}^l)\|^2 \right] \stackrel{(1)}{\leq} \frac{\sigma^2}{m_l}, \end{aligned}$$

for all  $l = 0, 1, \dots, T-1$ . Moreover, as we already derived, probability event  $E_{T-1}$  implies that  $\|F(x^l)\| \leq \lambda_l/2$  and  $\|F(\tilde{x}^l)\| \leq \lambda_l/2$  for all  $l = 0, 1, \dots, T-1$  (see (114) and (117)). Therefore, in view of Lemma B.2,  $E_{T-1}$  implies that

$$\|\theta_l^b\| \leq \frac{4\sigma^2}{m_l\lambda_l}, \quad \|\omega_l^b\| \leq \frac{4\sigma^2}{m_l\lambda_l}, \quad (125)$$

$$\mathbb{E}_{\xi_2^l} \left[ \|\theta_l\|^2 \right] \leq \frac{18\sigma^2}{m_l}, \quad \mathbb{E}_{\xi_1^l} \left[ \|\omega_l\|^2 \right] \leq \frac{18\sigma^2}{m_l}, \quad (126)$$

$$\mathbb{E}_{\xi_2^l} \left[ \|\theta_l^u\|^2 \right] \leq \frac{18\sigma^2}{m_l}, \quad \mathbb{E}_{\xi_1^l} \left[ \|\omega_l^u\|^2 \right] \leq \frac{18\sigma^2}{m_l}, \quad (127)$$

for all  $l = 0, 1, \dots, T-1$ .

**Upper bound for ①.** Since  $\mathbb{E}_{\xi_l^i}[\omega_l^u] = 0$ , we have

$$\mathbb{E}_{\xi_l^i}[-4\gamma^3\mu(1-\gamma\mu)^{T-1-l}\langle\zeta_l, \omega_l^u\rangle] = 0.$$

Next, the summands in ① are bounded with probability 1:

$$\begin{aligned} |-4\gamma^3\mu(1-\gamma\mu)^{T-1-l}\langle\zeta_l, \omega_l^u\rangle| &\leq 4\gamma^3\mu \exp(-\gamma\mu(T-1-l))\|\zeta_l\| \cdot \|\omega_l^u\| \\ &\stackrel{(120),(124)}{\leq} 8\sqrt{2}\gamma^3\mu L \exp(-\gamma\mu(T-1-l/2))R\lambda_l \\ &\stackrel{(109),(110)}{\leq} \frac{\exp(-\gamma\mu T)R^2}{7 \ln \frac{6(K+1)}{\beta}} \stackrel{\text{def}}{=} c. \end{aligned} \quad (128)$$

Moreover, these summands have bounded conditional variances  $\sigma_l^2 \stackrel{\text{def}}{=} \mathbb{E}_{\xi_l^i}[16\gamma^6\mu^2(1-\gamma\mu)^{2T-2-2l}\langle\zeta_l, \omega_l^u\rangle^2]$ :

$$\begin{aligned} \sigma_l^2 &\leq \mathbb{E}_{\xi_l^i}[16\gamma^6\mu^2 \exp(-\gamma\mu(2T-2-2l))\|\zeta_l\|^2 \cdot \|\omega_l^u\|^2] \\ &\stackrel{(120)}{\leq} 36\gamma^6\mu^2 L^2 \exp(-\gamma\mu(2T-2-2l))R^2 \mathbb{E}_{\xi_l^i}[\|\omega_l^u\|^2] \\ &\stackrel{(109)}{\leq} \frac{4\gamma^2 \exp(-\gamma\mu(2T-l))R^2}{2809 \ln \frac{6(K+1)}{\beta}} \mathbb{E}_{\xi_l^i}[\|\omega_l^u\|^2]. \end{aligned} \quad (129)$$

That is, sequence  $\{-4\gamma^3\mu(1-\gamma\mu)^{T-1-l}\langle\zeta_l, \omega_l^u\rangle\}_{l \geq 0}$  is a bounded martingale difference sequence having bounded conditional variances  $\{\sigma_l^2\}_{l \geq 0}$ . Applying Bernstein's inequality (Lemma B.1) with  $X_l = -4\gamma^3\mu(1-\gamma\mu)^{T-1-l}\langle\zeta_l, \omega_l^u\rangle$ ,  $c$  defined in (128),  $b = \frac{1}{7} \exp(-\gamma\mu T)R^2$ ,  $G = \frac{\exp(-2\gamma\mu T)R^4}{294 \ln \frac{6(K+1)}{\beta}}$ , we get that

$$\mathbb{P}\left\{|\textcircled{1}| > \frac{1}{7} \exp(-\gamma\mu T)R^2 \text{ and } \sum_{l=0}^{T-1} \sigma_l^2 \leq \frac{\exp(-2\gamma\mu T)R^4}{294 \ln \frac{6(K+1)}{\beta}}\right\} \leq 2 \exp\left(-\frac{b^2}{2G + 2cb/3}\right) = \frac{\beta}{3(K+1)}.$$

In other words,  $\mathbb{P}\{E_{\textcircled{1}}\} \geq 1 - \frac{\beta}{3(K+1)}$ , where probability event  $E_{\textcircled{1}}$  is defined as

$$E_{\textcircled{1}} = \left\{ \text{either } \sum_{l=0}^{T-1} \sigma_l^2 > \frac{\exp(-2\gamma\mu T)R^4}{294 \ln \frac{6(K+1)}{\beta}} \text{ or } |\textcircled{1}| \leq \frac{1}{7} \exp(-\gamma\mu T)R^2 \right\}. \quad (130)$$

Moreover, we notice here that probability event  $E_{T-1}$  implies that

$$\begin{aligned} \sum_{l=0}^{T-1} \sigma_l^2 &\stackrel{(129)}{\leq} \frac{4\gamma^2 \exp(-2\gamma\mu T)R^2}{2809 \ln \frac{6(K+1)}{\beta}} \sum_{l=0}^{T-1} \frac{\mathbb{E}_{\xi_l^i}[\|\omega_l^u\|^2]}{\exp(-\gamma\mu l)} \\ &\stackrel{(127), T \leq K+1}{\leq} \frac{72\gamma^2 \exp(-2\gamma\mu T)R^2 \sigma^2}{2809 \ln \frac{6(K+1)}{\beta}} \sum_{l=0}^K \frac{1}{m_l \exp(-\gamma\mu l)} \\ &\stackrel{(111)}{\leq} \frac{\exp(-2\gamma\mu T)R^4}{294 \ln \frac{6(K+1)}{\beta}}. \end{aligned} \quad (131)$$

**Upper bound for ②.** Probability event  $E_{T-1}$  implies

$$\begin{aligned} \textcircled{2} &\leq 4\gamma^3\mu \sum_{l=0}^{T-1} \exp(-\gamma\mu(T-1-l))\|\zeta_l\| \cdot \|\omega_l^b\| \\ &\stackrel{(120),(125)}{\leq} 16\sqrt{2} \exp(-\gamma\mu(T-1))\gamma^3\mu LR \sum_{l=0}^{T-1} \frac{\sigma^2}{m_l \lambda_l \exp(-\gamma\mu l/2)} \\ &\stackrel{(110)}{=} 1920\sqrt{2} \exp(-\gamma\mu(T-2))\gamma^4\mu L \sum_{l=0}^{T-1} \frac{\sigma^2 \ln \frac{6(K+1)}{\beta}}{m_l \exp(-\gamma\mu l)} \\ &\stackrel{(109),(111), T \leq K+1}{\leq} \frac{1}{7} \exp(-\gamma\mu T)R^2. \end{aligned} \quad (132)$$

**Upper bound for ③.** Since  $\mathbb{E}_{\xi_2^l}[\theta_l^u] = 0$ , we have

$$\mathbb{E}_{\xi_2^l} [2\gamma(1 - \gamma\mu)^{T-1-l} \langle \eta_l, \theta_l^u \rangle] = 0.$$

Next, the summands in ③ are bounded with probability 1:

$$\begin{aligned} |2\gamma(1 - \gamma\mu)^{T-1-l} \langle \eta_l, \theta_l^u \rangle| &\leq 2\gamma \exp(-\gamma\mu(T-1-l)) \|\eta_l\| \cdot \|\theta_l^u\| \\ &\stackrel{(120),(124)}{\leq} 4\sqrt{7}\gamma(1 + \gamma L) \exp(-\gamma\mu(T-1-l/2)) R \lambda_l \\ &\stackrel{(109),(110)}{\leq} \frac{\exp(-\gamma\mu T) R^2}{7 \ln \frac{6(K+1)}{\beta}} \stackrel{\text{def}}{=} c. \end{aligned} \quad (133)$$

Moreover, these summands have bounded conditional variances  $\tilde{\sigma}_l^2 \stackrel{\text{def}}{=} \mathbb{E}_{\xi_2^l} [4\gamma^2(1 - \gamma\mu)^{2T-2-2l} \langle \eta_l, \theta_l^u \rangle^2]$ :

$$\begin{aligned} \tilde{\sigma}_l^2 &\leq \mathbb{E}_{\xi_2^l} [4\gamma^2 \exp(-\gamma\mu(2T-2-2l)) \|\eta_l\|^2 \cdot \|\theta_l^u\|^2] \\ &\stackrel{(120)}{\leq} 49\gamma^2(1 + \gamma L)^2 \exp(-\gamma\mu(2T-2-l)) R^2 \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2] \\ &\stackrel{(109)}{\leq} 50\gamma^2 \exp(-\gamma\mu(2T-l)) R^2 \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2]. \end{aligned} \quad (134)$$

That is, sequence  $\{2\gamma(1 - \gamma\mu)^{T-1-l} \langle \eta_l, \theta_l^u \rangle\}_{l \geq 0}$  is a bounded martingale difference sequence having bounded conditional variances  $\{\tilde{\sigma}_l^2\}_{l \geq 0}$ . Applying Bernstein's inequality (Lemma B.1) with  $X_l = 2\gamma(1 - \gamma\mu)^{T-1-l} \langle \eta_l, \theta_l^u \rangle$ ,  $c$  defined in (133),  $b = \frac{1}{7} \exp(-\gamma\mu T) R^2$ ,  $G = \frac{\exp(-2\gamma\mu T) R^4}{294 \ln \frac{6(K+1)}{\beta}}$ , we get that

$$\mathbb{P} \left\{ |\textcircled{3}| > \frac{1}{7} \exp(-\gamma\mu T) R^2 \text{ and } \sum_{l=0}^{T-1} \tilde{\sigma}_l^2 \leq \frac{\exp(-2\gamma\mu T) R^4}{294 \ln \frac{6(K+1)}{\beta}} \right\} \leq 2 \exp \left( -\frac{b^2}{2G + 2cb/3} \right) = \frac{\beta}{3(K+1)}.$$

In other words,  $\mathbb{P}\{E_{\textcircled{3}}\} \geq 1 - \frac{\beta}{3(K+1)}$ , where probability event  $E_{\textcircled{3}}$  is defined as

$$E_{\textcircled{3}} = \left\{ \text{either } \sum_{l=0}^{T-1} \tilde{\sigma}_l^2 > \frac{\exp(-2\gamma\mu T) R^4}{294 \ln \frac{6(K+1)}{\beta}} \text{ or } |\textcircled{3}| \leq \frac{1}{7} \exp(-\gamma\mu T) R^2 \right\}. \quad (135)$$

Moreover, we notice here that probability event  $E_{T-1}$  implies that

$$\begin{aligned} \sum_{l=0}^{T-1} \tilde{\sigma}_l^2 &\stackrel{(134)}{\leq} 50\gamma^2 \exp(-2\gamma\mu T) R^2 \sum_{l=0}^{T-1} \frac{\mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2]}{\exp(-\gamma\mu l)} \\ &\stackrel{(127), T \leq K+1}{\leq} 900\gamma^2 \exp(-2\gamma\mu T) R^2 \sigma^2 \sum_{l=0}^K \frac{1}{m_l \exp(-\gamma\mu l)} \\ &\stackrel{(111)}{\leq} \frac{\exp(-2\gamma\mu T) R^4}{294 \ln \frac{6(K+1)}{\beta}}. \end{aligned} \quad (136)$$

**Upper bound for ④.** Probability event  $E_{T-1}$  implies

$$\begin{aligned} \textcircled{4} &\leq 2\gamma \exp(-\gamma\mu(T-1)) \sum_{l=0}^{T-1} \frac{\|\eta_l\| \cdot \|\theta_l^b\|}{\exp(-\gamma\mu l)} \\ &\stackrel{(120),(125)}{\leq} 8\sqrt{7}\gamma(1 + \gamma L) \exp(-\gamma\mu(T-1)) R \sum_{l=0}^{T-1} \frac{\sigma^2}{m_l \lambda_l \exp(-\gamma\mu l/2)} \\ &\stackrel{(110)}{\leq} 960\sqrt{7}\gamma^2(1 + \gamma L) \exp(-\gamma\mu(T-2)) \sum_{l=0}^{T-1} \frac{\sigma^2 \ln \frac{6(K+1)}{\beta}}{m_l \exp(-\gamma\mu l)} \\ &\stackrel{(111), T \leq K+1}{\leq} \frac{1}{7} \exp(-\gamma\mu T) R^2. \end{aligned} \quad (137)$$

**Upper bound for ⑤.** Probability event  $E_{T-1}$  implies

$$\begin{aligned}
\textcircled{5} &= 2\gamma^2 \exp(-\gamma\mu(T-1)) \sum_{l=0}^{T-1} \frac{\mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2] + 4\mathbb{E}_{\xi_1^l} [\|\omega_l^u\|^2]}{\exp(-\gamma\mu l)} \\
&\stackrel{(127)}{\leq} 180\gamma^2 \exp(-\gamma\mu(T-1)) \sum_{l=0}^{T-1} \frac{\sigma^2}{m_l \exp(-\gamma\mu l)} \\
&\stackrel{(111), T \leq K+1}{\leq} \frac{1}{7} \exp(-\gamma\mu T) R^2. \tag{138}
\end{aligned}$$

**Upper bound for ⑥.** First of all, we have

$$2\gamma^2(1-\gamma\mu)^{T-1-l} \mathbb{E}_{\xi_1^l, \xi_2^l} \left[ \|\theta_l^u\|^2 + 4\|\omega_l^u\|^2 - \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2] - 4\mathbb{E}_{\xi_1^l} [\|\omega_l^u\|^2] \right] = 0.$$

Next, the summands in ⑥ are bounded with probability 1:

$$\begin{aligned}
2\gamma^2(1-\gamma\mu)^{T-1-l} \left| \|\theta_l^u\|^2 + 4\|\omega_l^u\|^2 - \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2] - 4\mathbb{E}_{\xi_1^l} [\|\omega_l^u\|^2] \right| &\stackrel{(124)}{\leq} \frac{80\gamma^2 \exp(-\gamma\mu T) \lambda_l^2}{\exp(-\gamma\mu(1+l))} \\
&\stackrel{(110)}{\leq} \frac{\exp(-\gamma\mu T) R^2}{7 \ln \frac{6(K+1)}{\beta}} \\
&\stackrel{\text{def}}{=} c. \tag{139}
\end{aligned}$$

Moreover, these summands have bounded conditional variances  $\hat{\sigma}_l^2 \stackrel{\text{def}}{=} \mathbb{E}_{\xi_1^l, \xi_2^l} \left[ 4\gamma^4(1-\gamma\mu)^{2T-2-2l} \left| \|\theta_l^u\|^2 + 4\|\omega_l^u\|^2 - \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2] - 4\mathbb{E}_{\xi_1^l} [\|\omega_l^u\|^2] \right|^2 \right]$ :

$$\begin{aligned}
\hat{\sigma}_l^2 &\stackrel{(139)}{\leq} \frac{2\gamma^2 \exp(-2\gamma\mu T) R^2}{7 \exp(-\gamma\mu(1+l)) \ln \frac{6(K+1)}{\beta}} \mathbb{E}_{\xi_1^l, \xi_2^l} \left[ \left| \|\theta_l^u\|^2 + 4\|\omega_l^u\|^2 - \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2] - 4\mathbb{E}_{\xi_1^l} [\|\omega_l^u\|^2] \right|^2 \right] \\
&\leq \frac{4\gamma^2 \exp(-2\gamma\mu T) R^2}{7 \exp(-\gamma\mu(1+l)) \ln \frac{6(K+1)}{\beta}} \mathbb{E}_{\xi_1^l, \xi_2^l} [\|\theta_l^u\|^2 + 4\|\omega_l^u\|^2]. \tag{140}
\end{aligned}$$

That is, sequence  $\left\{ 2\gamma^2(1-\gamma\mu)^{T-1-l} \left( \|\theta_l^u\|^2 + 4\|\omega_l^u\|^2 - \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2] - 4\mathbb{E}_{\xi_1^l} [\|\omega_l^u\|^2] \right) \right\}_{l \geq 0}$  is a bounded martingale difference sequence having bounded conditional variances  $\{\hat{\sigma}_l^2\}_{l \geq 0}$ . Applying Bernstein's inequality (Lemma B.1) with  $X_l = 2\gamma^2(1-\gamma\mu)^{T-1-l} \left( \|\theta_l^u\|^2 + 4\|\omega_l^u\|^2 - \mathbb{E}_{\xi_2^l} [\|\theta_l^u\|^2] - 4\mathbb{E}_{\xi_1^l} [\|\omega_l^u\|^2] \right)$ ,  $c$  defined in (139),  $b = \frac{1}{7} \exp(-\gamma\mu T) R^2$ ,  $G = \frac{\exp(-2\gamma\mu T) R^4}{294 \ln \frac{6(K+1)}{\beta}}$ , we get that

$$\mathbb{P} \left\{ |\textcircled{6}| > \frac{1}{7} \exp(-\gamma\mu T) R^2 \text{ and } \sum_{l=0}^{T-1} \hat{\sigma}_l^2 \leq \frac{\exp(-2\gamma\mu T) R^4}{294 \ln \frac{6(K+1)}{\beta}} \right\} \leq 2 \exp \left( -\frac{b^2}{2G + 2cb/3} \right) = \frac{\beta}{3(K+1)}.$$

In other words,  $\mathbb{P}\{E_{\textcircled{6}}\} \geq 1 - \frac{\beta}{3(K+1)}$ , where probability event  $E_{\textcircled{6}}$  is defined as

$$E_{\textcircled{6}} = \left\{ \text{either } \sum_{l=0}^{T-1} \hat{\sigma}_l^2 > \frac{\exp(-2\gamma\mu T) R^4}{294 \ln \frac{6(K+1)}{\beta}} \text{ or } |\textcircled{6}| \leq \frac{1}{7} \exp(-\gamma\mu T) R^2 \right\}. \tag{141}$$

Moreover, we notice here that probability event  $E_{T-1}$  implies that

$$\begin{aligned}
\sum_{l=0}^{T-1} \hat{\sigma}_l^2 &\stackrel{(140)}{\leq} \frac{4\gamma^2 \exp(-\gamma\mu(2T-1)) R^2}{7 \ln \frac{6(K+1)}{\beta}} \sum_{l=0}^{T-1} \frac{\mathbb{E}_{\xi_1^l, \xi_2^l} [\|\theta_l^u\|^2 + 4\|\omega_l^u\|^2]}{\exp(-\gamma\mu l)} \\
&\stackrel{(127), T \leq K+1}{\leq} \frac{360\gamma^2 \exp(-\gamma\mu(2T-1)) R^2 \sigma^2}{7 \ln \frac{6(K+1)}{\beta}} \sum_{l=0}^K \frac{1}{m_l \exp(-\gamma\mu l)} \\
&\stackrel{(111)}{\leq} \frac{\exp(-2\gamma\mu T) R^4}{294 \ln \frac{6(K+1)}{\beta}}. \tag{142}
\end{aligned}$$



**Upper bound for ⑦.** Probability event  $E_{T-1}$  implies

$$\begin{aligned}
\textcircled{7} &= 2\gamma^2 \sum_{l=0}^{T-1} \exp(-\gamma\mu(T-1-l)) (\|\theta_l^b\|^2 + 4\|\omega_l^b\|^2) \\
&\stackrel{(125)}{\leq} 160\gamma^2 \exp(-\gamma\mu(T-1)) \sum_{l=0}^{T-1} \frac{\sigma^4}{m_l^2 \lambda_l^2 \exp(-\gamma\mu l)} \\
&\stackrel{(110)}{=} 2304000\gamma^4 \exp(-\gamma\mu(T-3)) \sum_{l=0}^{T-1} \frac{\sigma^4 \ln^2 \frac{6(K+1)}{\beta}}{m_l^2 R^2 \exp(-2\gamma\mu l)} \\
&\stackrel{(111), T \leq K+1}{\leq} \frac{1}{7} \exp(-\gamma\mu T) R^2. \tag{143}
\end{aligned}$$

**Final derivation.** Putting all bounds together, we get that  $E_{T-1}$  implies

$$\begin{aligned}
R_T^2 &\stackrel{(123)}{\leq} \exp(-\gamma\mu T) R^2 + \textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4} + \textcircled{5} + \textcircled{6} + \textcircled{7}, \\
\textcircled{2} &\stackrel{(132)}{\leq} \frac{1}{7} \exp(-\gamma\mu T) R^2, \quad \textcircled{4} \stackrel{(137)}{\leq} \frac{1}{7} \exp(-\gamma\mu T) R^2, \\
\textcircled{5} &\stackrel{(138)}{\leq} \frac{1}{7} \exp(-\gamma\mu T) R^2, \quad \textcircled{7} \stackrel{(143)}{\leq} \frac{1}{7} \exp(-\gamma\mu T) R^2, \\
\sum_{l=0}^{T-1} \sigma_l^2 &\stackrel{(131)}{\leq} \frac{\exp(-2\gamma\mu T) R^4}{294 \ln \frac{6(K+1)}{\beta}}, \quad \sum_{l=0}^{T-1} \tilde{\sigma}_l^2 \stackrel{(136)}{\leq} \frac{\exp(-2\gamma\mu T) R^4}{294 \ln \frac{6(K+1)}{\beta}}, \quad \sum_{l=0}^{T-1} \hat{\sigma}_l^2 \stackrel{(142)}{\leq} \frac{\exp(-2\gamma\mu T) R^4}{294 \ln \frac{6(K+1)}{\beta}}.
\end{aligned}$$

Moreover, in view of (130), (135), (141), and our induction assumption, we have

$$\begin{aligned}
\mathbb{P}\{E_{T-1}\} &\geq 1 - \frac{(T-1)\beta}{K+1}, \\
\mathbb{P}\{E_{\textcircled{1}}\} &\geq 1 - \frac{\beta}{3(K+1)}, \quad \mathbb{P}\{E_{\textcircled{3}}\} \geq 1 - \frac{\beta}{3(K+1)}, \quad \mathbb{P}\{E_{\textcircled{6}}\} \geq 1 - \frac{\beta}{3(K+1)},
\end{aligned}$$

where probability events  $E_{\textcircled{1}}$ ,  $E_{\textcircled{3}}$ , and  $E_{\textcircled{6}}$  are defined as

$$\begin{aligned}
E_{\textcircled{1}} &= \left\{ \text{either } \sum_{l=0}^{T-1} \sigma_l^2 > \frac{\exp(-2\gamma\mu T) R^4}{294 \ln \frac{6(K+1)}{\beta}} \text{ or } |\textcircled{1}| \leq \frac{1}{7} \exp(-\gamma\mu T) R^2 \right\}, \\
E_{\textcircled{3}} &= \left\{ \text{either } \sum_{l=0}^{T-1} \tilde{\sigma}_l^2 > \frac{\exp(-2\gamma\mu T) R^4}{294 \ln \frac{6(K+1)}{\beta}} \text{ or } |\textcircled{3}| \leq \frac{1}{7} \exp(-\gamma\mu T) R^2 \right\}, \\
E_{\textcircled{6}} &= \left\{ \text{either } \sum_{l=0}^{T-1} \hat{\sigma}_l^2 > \frac{\exp(-2\gamma\mu T) R^4}{294 \ln \frac{6(K+1)}{\beta}} \text{ or } |\textcircled{6}| \leq \frac{1}{7} \exp(-\gamma\mu T) R^2 \right\}.
\end{aligned}$$

Putting all of these inequalities together, we obtain that probability event  $E_{T-1} \cap E_{\textcircled{1}} \cap E_{\textcircled{3}} \cap E_{\textcircled{6}}$  implies

$$\begin{aligned}
R_T^2 &\stackrel{(123)}{\leq} \exp(-\gamma\mu T) R^2 + \textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4} + \textcircled{5} + \textcircled{6} + \textcircled{7} \\
&\leq 2 \exp(-\gamma\mu T) R^2.
\end{aligned}$$

Moreover, union bound for the probability events implies

$$\mathbb{P}\{E_T\} \geq \mathbb{P}\{E_{T-1} \cap E_{\textcircled{1}} \cap E_{\textcircled{3}} \cap E_{\textcircled{6}}\} = 1 - \mathbb{P}\{\bar{E}_{T-1} \cup \bar{E}_{\textcircled{1}} \cup \bar{E}_{\textcircled{3}} \cup \bar{E}_{\textcircled{6}}\} \geq 1 - \frac{T\beta}{K+1}. \tag{144}$$

This is exactly what we wanted to prove (see the paragraph after inequality (113)). In particular, with probability at least  $1 - \beta$  satisfy we have

$$\|x^{K+1} - x^*\|^2 \leq 2 \exp(-\gamma\mu(K+1)) R^2,$$

which finishes the proof.  $\square$

**Corollary C.3.** *Let the assumptions of Theorem C.3 hold. Then, the following statements hold.*

1. **Large stepsize/large batch.** *The choice of stepsize and batchsize*

$$\gamma = \frac{1}{650L \ln \frac{6(K+1)}{\beta}}, \quad m_k = \max \left\{ 1, \frac{264600\gamma^2(K+1)\sigma^2 \ln \frac{6(K+1)}{\beta}}{\exp(-\gamma\mu k)R^2} \right\} \quad (145)$$

satisfies conditions (109) and (111). With such choice of  $\gamma$ ,  $m_k$ , and the choice of  $\lambda_k$  as in (110), the iterates produced by **clipped-SEG** after  $K$  iterations with probability at least  $1 - \beta$  satisfy

$$\|x^{K+1} - x^*\|^2 \leq 2 \exp \left( -\frac{\mu(K+1)}{650L \ln \frac{6(K+1)}{\beta}} \right) R^2. \quad (146)$$

In particular, to guarantee  $\|x^{K+1} - x^*\|^2 \leq \varepsilon$  with probability at least  $1 - \beta$  for some  $\varepsilon > 0$  **clipped-SEG** requires

$$\mathcal{O} \left( \frac{L}{\mu} \ln \left( \frac{R^2}{\varepsilon} \right) \ln \left( \frac{L}{\mu\beta} \ln \left( \frac{R^2}{\varepsilon} \right) \right) \right) \text{ iterations}, \quad (147)$$

$$\mathcal{O} \left( \max \left\{ \frac{L}{\mu}, \frac{\sigma^2}{\mu^2\varepsilon} \right\} \ln \left( \frac{R^2}{\varepsilon} \right) \ln \left( \frac{L}{\mu\beta} \ln \left( \frac{R^2}{\varepsilon} \right) \right) \right) \text{ oracle calls}. \quad (148)$$

2. **Small stepsize/small batch.** *The choice of stepsize and batchsize*

$$\gamma = \min \left\{ \frac{1}{650L \ln \frac{6(K+1)}{\beta}}, \frac{\ln(B_K)}{\mu(K+1)} \right\}, \quad m_k \equiv 1 \quad (149)$$

satisfies conditions (109) and (111), where  $B_K = \max \left\{ 2, \frac{(K+1)\mu^2 R^2}{264600\sigma^2 \ln \left( \frac{6(K+1)}{\beta} \right) \ln^2(B_K)} \right\} = \mathcal{O} \left( \max \left\{ 2, \frac{(K+1)\mu^2 R^2}{264600\sigma^2 \ln \left( \frac{6(K+1)}{\beta} \right) \ln^2 \left( \max \left\{ 2, \frac{(K+1)\mu^2 R^2}{264600\sigma^2 \ln \left( \frac{6(K+1)}{\beta} \right)} \right\} \right)} \right\} \right)$ . With such choice

of  $\gamma$ ,  $m_k$ , and the choice of  $\lambda_k$  as in (110), the iterates produced by **clipped-SEG** after  $K$  iterations with probability at least  $1 - \beta$  satisfy

$$\|x^{K+1} - x^*\|^2 \leq \max \left\{ 2 \exp \left( -\frac{\mu(K+1)}{650L \ln \frac{6(K+1)}{\beta}} \right) R^2, \frac{529200\sigma^2 \ln \left( \frac{6(K+1)}{\beta} \right) \ln^2(B_K)}{\mu^2(K+1)} \right\}. \quad (150)$$

In particular, to guarantee  $\|x^{K+1} - x^*\|^2 \leq \varepsilon$  with probability at least  $1 - \beta$  for some  $\varepsilon > 0$  **clipped-SEG** requires

$$\mathcal{O} \left( \max \left\{ \frac{L}{\mu} \ln \left( \frac{R^2}{\varepsilon} \right) \ln \left( \frac{L}{\mu\beta} \ln \left( \frac{R^2}{\varepsilon} \right) \right), \frac{\sigma^2}{\mu^2\varepsilon} \ln \left( \frac{\sigma^2}{\mu^2\varepsilon\beta} \right) \ln^2(B_\varepsilon) \right\} \right) \quad (151)$$

iterations/oracle calls, where

$$B_\varepsilon = \max \left\{ 2, \frac{R^2}{\varepsilon \ln \left( \frac{\sigma^2}{\mu^2\varepsilon\beta} \right) \ln^2 \left( \max \left\{ 2, \frac{R^2}{\varepsilon \ln \left( \frac{\sigma^2}{\mu^2\varepsilon\beta} \right)} \right\} \right)} \right\}.$$

*Proof.* 1. **Large stepsize/large batch.** First of all, it is easy to see that the choice of  $\gamma$  and  $m_k$  from (145) satisfies conditions (109) and (111). Therefore, applying Theorem C.3, we derive that with probability at least  $1 - \beta$

$$\|x^{K+1} - x^*\|^2 \leq 2 \exp(-\gamma\mu(K+1))R^2 \stackrel{(145)}{=} 2 \exp \left( -\frac{\mu(K+1)}{650L \ln \frac{6(K+1)}{\beta}} \right) R^2.$$

To guarantee  $\|x^{K+1} - x^*\|^2 \leq \varepsilon$ , we choose  $K$  in such a way that the right-hand side of the above inequality is smaller than  $\varepsilon$  that gives

$$K = \mathcal{O}\left(\frac{L}{\mu} \ln\left(\frac{R^2}{\varepsilon}\right) \ln\left(\frac{L}{\mu\beta} \ln\left(\frac{R^2}{\varepsilon}\right)\right)\right).$$

The total number of oracle calls equals

$$\begin{aligned} \sum_{k=0}^K 2m_k &\stackrel{(145)}{=} 2 \sum_{k=0}^K \max\left\{1, \frac{264600\gamma^2(K+1)\sigma^2 \ln \frac{6(K+1)}{\beta}}{\exp(-\gamma\mu k)R^2}\right\} \\ &= \mathcal{O}\left(\max\left\{K, \frac{\gamma(K+1) \exp(\gamma\mu(K+1))\sigma^2 \ln \frac{6(K+1)}{\beta}}{\mu R^2}\right\}\right) \\ &= \mathcal{O}\left(\max\left\{\frac{L}{\mu}, \frac{\sigma^2}{\mu^2\varepsilon}\right\} \ln\left(\frac{R^2}{\varepsilon}\right) \ln\left(\frac{L}{\mu\beta} \ln\left(\frac{R^2}{\varepsilon}\right)\right)\right). \end{aligned}$$

2. **Small stepsize/small batch.** First of all, we verify that the choice of  $\gamma$  and  $m_k$  from (149) satisfies conditions (109) and (111): (109) trivially holds and (111) holds since for all  $k = 0, \dots, K$

$$\begin{aligned} \frac{264600\gamma^2(K+1)\sigma^2 \ln \frac{6(K+1)}{\beta}}{\exp(-\gamma\mu k)R^2} &\leq \frac{264600\gamma^2(K+1)\sigma^2 \ln \frac{6(K+1)}{\beta}}{\exp(-\gamma\mu(K+1))R^2} \\ &\stackrel{(149)}{\leq} \frac{264600 \ln^2(B_K) \exp(\gamma\mu(K+1))\sigma^2 \ln \frac{6(K+1)}{\beta}}{\mu^2(K+1)R^2} \\ &\stackrel{(149)}{\leq} 1. \end{aligned}$$

Therefore, applying Theorem C.3, we derive that with probability at least  $1 - \beta$

$$\begin{aligned} \|x^{K+1} - x^*\|^2 &\leq 2 \exp(-\gamma\mu(K+1))R^2 \\ &\stackrel{(149)}{=} \max\left\{2 \exp\left(-\frac{\mu(K+1)}{650L \ln \frac{6(K+1)}{\beta}}\right) R^2, \frac{2R^2}{B_K}\right\} \\ &= \max\left\{2 \exp\left(-\frac{\mu(K+1)}{650L \ln \frac{6(K+1)}{\beta}}\right) R^2, \frac{529200\sigma^2 \ln\left(\frac{6(K+1)}{\beta}\right) \ln^2(B_K)}{\mu^2(K+1)}\right\}. \end{aligned}$$

To guarantee  $\|x^{K+1} - x^*\|^2 \leq \varepsilon$ , we choose  $K$  in such a way that the right-hand side of the above inequality is smaller than  $\varepsilon$  that gives  $K$  of the order

$$\mathcal{O}\left(\max\left\{\frac{L}{\mu} \ln\left(\frac{R^2}{\varepsilon}\right) \ln\left(\frac{L}{\mu\beta} \ln\left(\frac{R^2}{\varepsilon}\right)\right), \frac{\sigma^2}{\mu^2\varepsilon} \ln\left(\frac{\sigma^2}{\mu^2\varepsilon\beta}\right) \ln^2(B_\varepsilon)\right\}\right),$$

where

$$B_\varepsilon = \max\left\{2, \frac{R^2}{\varepsilon \ln\left(\frac{\sigma^2}{\mu^2\varepsilon\beta}\right) \ln^2\left(\max\left\{2, \frac{R^2}{\varepsilon \ln\left(\frac{\sigma^2}{\mu^2\varepsilon\beta}\right)}\right\}\right)}\right\}.$$

The total number of oracle calls equals  $\sum_{k=0}^K 2m_k = 2(K+1)$ .

□

## D Clipped Stochastic Gradient Descent-Ascent: Missing Proofs and Details

### D.1 Monotone Star-Cocoercive Case

**Lemma D.1.** *Let Assumption 1.3 hold for  $Q = B_{2R}(x^*)$ , where  $R \geq R_0 \stackrel{\text{def}}{=} \|x^0 - x^*\|$  and  $0 < \gamma \leq 2/\ell$ . If  $x^k$  lies in  $B_{2R}(x^*)$  for all  $k = 0, 1, \dots, K$  for some  $K \geq 0$ , then for all  $u \in B_{3R}(x^*)$  the iterates produced by **clipped-SGDA** satisfy*

$$\begin{aligned} 2\gamma \langle F(u), x_{\text{avg}}^K - u \rangle &\leq \frac{\|x^0 - u\|^2 - \|x^{K+1} - u\|^2}{K+1} \\ &\quad + \frac{2\gamma}{K+1} \sum_{k=0}^K \langle x^k - u - \gamma F(x^k), \omega_k \rangle \\ &\quad + \frac{\gamma^2}{K+1} \sum_{k=0}^K (\|F(x^k)\|^2 + \|\omega_k\|^2), \end{aligned} \quad (152)$$

$$x_{\text{avg}}^K \stackrel{\text{def}}{=} \frac{1}{K+1} \sum_{k=0}^K x^k, \quad (153)$$

$$\omega_k \stackrel{\text{def}}{=} F(x^k) - \tilde{F}_{\xi^k}(x^k). \quad (154)$$

*Proof.* Using the update rule of **clipped-SGDA**, we obtain

$$\begin{aligned} \|x^{k+1} - u\|^2 &= \|x^k - u\|^2 - 2\gamma \langle x^k - u, \tilde{F}_{\xi^k}(x^k) \rangle + \gamma^2 \|\tilde{F}_{\xi^k}(x^k)\|^2 \\ &= \|x^k - u\|^2 - 2\gamma \langle x^k - u, F(x^k) \rangle + 2\gamma \langle x^k - u, \omega_k \rangle \\ &\quad + \gamma^2 \|F(x^k)\|^2 - 2\gamma^2 \langle F(x^k), \omega_k \rangle + \gamma^2 \|\omega_k\|^2 \\ &\stackrel{\text{(Mon)}}{\leq} \|x^k - u\|^2 - 2\gamma \langle x^k - u, F(u) \rangle + 2\gamma \langle x^k - u - \gamma F(x^k), \omega_k \rangle \\ &\quad + \gamma^2 (\|F(x^k)\|^2 + \|\omega_k\|^2). \end{aligned}$$

Rearranging the terms, we derive

$$\begin{aligned} 2\gamma \langle F(u), x^k - u \rangle &\leq \|x^k - u\|^2 - \|x^{k+1} - u\|^2 + 2\gamma \langle x^k - u - \gamma F(x^k), \omega_k \rangle \\ &\quad + \gamma^2 (\|F(x^k)\|^2 + \|\omega_k\|^2). \end{aligned}$$

Finally, we sum up the above inequality for  $k = 0, 1, \dots, K$  and divide both sides of the result by  $(K+1)$ :

$$\begin{aligned} 2\gamma \langle F(u), x_{\text{avg}}^K - u \rangle &\leq \frac{1}{K+1} \sum_{k=0}^K (\|x^k - u\|^2 - \|x^{k+1} - u\|^2) \\ &\quad + \frac{2\gamma}{K+1} \sum_{k=0}^K \langle x^k - u - \gamma F(x^k), \omega_k \rangle \\ &\quad + \frac{\gamma^2}{K+1} \sum_{k=0}^K (\|F(x^k)\|^2 + \|\omega_k\|^2) \\ &= \frac{\|x^0 - u\|^2 - \|x^{K+1} - u\|^2}{K+1} \\ &\quad + \frac{2\gamma}{K+1} \sum_{k=0}^K \langle x^k - u - \gamma F(x^k), \omega_k \rangle \\ &\quad + \frac{\gamma^2}{K+1} \sum_{k=0}^K (\|F(x^k)\|^2 + \|\omega_k\|^2). \end{aligned}$$

This finishes the proof.  $\square$

We also derive the following lemma, which we use in the analysis of the star-cocoercive case as well.

**Lemma D.2.** *Let Assumption 1.6 hold for  $Q = B_{2R}(x^*)$ , where  $R \geq R_0 \stackrel{\text{def}}{=} \|x^0 - x^*\|$  and  $0 < \gamma \leq 2/\ell$ . If  $x^k$  lies in  $B_{2R}(x^*)$  for all  $k = 0, 1, \dots, K$  for some  $K \geq 0$ , then the iterates produced by clipped-SGDA satisfy*

$$\begin{aligned} \frac{\gamma}{K+1} \left( \frac{2}{\ell} - \gamma \right) \sum_{k=0}^K \|F(x^k)\|^2 &\leq \frac{\|x^0 - x^*\|^2 - \|x^{K+1} - x^*\|^2}{K+1} \\ &\quad + \frac{2\gamma}{K+1} \sum_{k=0}^K \langle x^k - x^* - \gamma F(x^k), \omega_k \rangle \\ &\quad + \frac{\gamma^2}{K+1} \sum_{k=0}^K \|\omega_k\|^2, \end{aligned} \tag{155}$$

where  $\omega_k$  is defined in (154).

*Proof.* Using the update rule of clipped-SGDA, we obtain

$$\begin{aligned} \|x^{k+1} - x^*\|^2 &= \|x^k - x^*\|^2 - 2\gamma \langle x^k - x^*, \tilde{F}_{\xi^k}(x^k) \rangle + \gamma^2 \|\tilde{F}_{\xi^k}(x^k)\|^2 \\ &= \|x^k - x^*\|^2 - 2\gamma \langle x^k - x^*, F(x^k) \rangle + 2\gamma \langle x^k - x^*, \omega_k \rangle \\ &\quad + \gamma^2 \|F(x^k)\|^2 - 2\gamma^2 \langle F(x^k), \omega_k \rangle + \gamma^2 \|\omega_k\|^2 \\ &\stackrel{\text{(SC)}}{\leq} \|x^k - x^*\|^2 + 2\gamma \langle x^k - x^*, \omega_k \rangle - 2\gamma^2 \langle F(x^k), \omega_k \rangle \\ &\quad + \gamma \left( \gamma - \frac{2}{\ell} \right) \|F(x^k)\|^2 + \gamma^2 \|\omega_k\|^2. \end{aligned}$$

Since  $0 < \gamma \leq 2/\ell$ , we have  $\gamma(2/\ell - \gamma) \|F(x^k)\|^2 \geq 0$  and, rearranging the terms, we derive

$$\begin{aligned} \gamma \left( \frac{2}{\ell} - \gamma \right) \|F(x^k)\|^2 &\leq \|x^k - x^*\|^2 - \|x^{k+1} - x^*\|^2 + 2\gamma \langle x^k - x^*, \omega_k \rangle \\ &\quad - 2\gamma^2 \langle F(x^k), \omega_k \rangle + \gamma^2 \|\omega_k\|^2. \end{aligned}$$

Finally, we sum up the above inequality for  $k = 0, 1, \dots, K$  and divide both sides of the result by  $(K+1)$ :

$$\begin{aligned} \frac{\gamma}{K+1} \left( \frac{2}{\ell} - \gamma \right) \sum_{k=0}^K \|F(x^k)\|^2 &\leq \frac{1}{K+1} \sum_{k=0}^K (\|x^k - x^*\|^2 - \|x^{k+1} - x^*\|^2) + \frac{\gamma^2}{K+1} \sum_{k=0}^K \|\omega_k\|^2 \\ &\quad + \frac{2\gamma}{K+1} \sum_{k=0}^K \langle x^k - x^*, \omega_k \rangle - \frac{2\gamma^2}{K+1} \sum_{k=0}^K \langle F(x^k), \omega_k \rangle \\ &= \frac{\|x^0 - x^*\|^2 - \|x^{K+1} - x^*\|^2}{K+1} + \frac{\gamma^2}{K+1} \sum_{k=0}^K \|\omega_k\|^2 \\ &\quad + \frac{2\gamma}{K+1} \sum_{k=0}^K \langle x^k - x^*, \omega_k \rangle - \frac{2\gamma^2}{K+1} \sum_{k=0}^K \langle F(x^k), \omega_k \rangle. \end{aligned}$$

This finishes the proof.  $\square$

**Theorem D.1.** *Let Assumptions 1.1, 1.3, 1.6, hold for  $Q = B_{2R}(x^*)$ , where  $R \geq R_0 \stackrel{\text{def}}{=} \|x^0 - x^*\|$ , and*

$$\gamma \leq \frac{1}{170\ell \ln \frac{6(K+1)}{\beta}}, \tag{156}$$

$$\lambda = \frac{R}{60\gamma \ln \frac{6(K+1)}{\beta}}, \tag{157}$$

$$m \geq \max \left\{ 1, \frac{97200(K+1)\gamma^2\sigma^2 \ln \frac{6(K+1)}{\beta}}{R^2} \right\}, \tag{158}$$

for some  $K \geq 0$  and  $\beta \in (0, 1]$  such that  $\ln \frac{6(K+1)}{\beta} \geq 1$ . Then, after  $K$  iterations the iterates produced by **clipped-SGDA** with probability at least  $1 - \beta$  satisfy

$$\text{Gap}_R(x_{\text{avg}}^K) \leq \frac{9R^2}{2\gamma(K+1)}. \quad (159)$$

*Proof.* We introduce new notation:  $R_k = \|x^k - x^*\|$  for all  $k \geq 0$ . The proof is based on the induction. In particular, for each  $k = 0, \dots, K+1$  we define the probability event  $E_k$  as follows: inequalities

$$\|x^t - x^*\|^2 \leq 2R^2 \quad \text{and} \quad \gamma \left\| \sum_{l=0}^{t-1} \omega_l \right\| \leq R \quad (160)$$

hold for  $t = 0, 1, \dots, k$  simultaneously. Our goal is to prove that  $\mathbb{P}\{E_k\} \geq 1 - k\beta/(K+1)$  for all  $k = 0, 1, \dots, K+1$ . We use the induction to show this statement. For  $k = 0$  the statement is trivial since  $R_0^2 \leq 2R^2$  by definition and  $\sum_{l=0}^{-1} \omega_l = 0$ . Next, assume that the statement holds for  $k = T \leq K$ , i.e., we have  $\mathbb{P}\{E_T\} \geq 1 - T\beta/(K+1)$ . We need to prove that  $\mathbb{P}\{E_{T+1}\} \geq 1 - (T+1)\beta/(K+1)$ . Let us notice that probability event  $E_T$  implies  $x^t \in B_{2R}(x^*)$  for all  $t = 0, 1, \dots, T$ . This means that the assumptions of Lemma D.2 hold and we have that probability event  $E_T$  implies  $(\gamma < 1/\ell)$

$$\begin{aligned} \frac{\gamma}{\ell(T+1)} \sum_{t=0}^T \|F(x^t)\|^2 &\leq \frac{\|x^0 - x^*\|^2 - \|x^{T+1} - x^*\|^2}{T+1} \\ &\quad + \frac{2\gamma}{T+1} \sum_{t=0}^T \langle x^t - x^* - \gamma F(x^t), \omega_t \rangle \\ &\quad + \frac{\gamma^2}{T+1} \sum_{t=0}^T \|\omega_t\|^2 \end{aligned} \quad (161)$$

and

$$\|F(x^t)\| \stackrel{\text{(SC)}}{\leq} \ell \|x^t - x^*\| \stackrel{\text{(160)}}{\leq} \sqrt{2}\ell R \stackrel{\text{(156),(157)}}{\leq} \frac{\lambda}{2} \quad (162)$$

for all  $t = 0, 1, \dots, T$ . From (161) we have

$$R_{T+1}^2 \leq R_0^2 + 2\gamma \sum_{t=0}^T \langle x^t - x^* - \gamma F(x^t), \omega_t \rangle + \gamma^2 \sum_{t=0}^T \|\omega_t\|^2.$$

Next, we notice that

$$\begin{aligned} \|x^t - x^* - \gamma F(x^t)\| &\leq \|x^t - x^*\| + \gamma \|F(x^t)\| \stackrel{\text{(SC),(160)}}{\leq} 2R + \gamma\ell \|x^t - x^*\| \\ &\stackrel{\text{(160)}}{\leq} 2R + 2R\gamma\ell \stackrel{\text{(156)}}{\leq} 3R, \end{aligned} \quad (163)$$

for all  $t = 0, 1, \dots, T$ . Consider random vectors

$$\eta_t = \begin{cases} x^t - x^* - \gamma F(x^t), & \text{if } \|x^t - x^* - \gamma F(x^t)\| \leq 3R, \\ 0, & \text{otherwise,} \end{cases}$$

for all  $t = 0, 1, \dots, T$ . We notice that  $\eta_t$  is bounded with probability 1:

$$\|\eta_t\| \leq 3R \quad (164)$$

for all  $t = 0, 1, \dots, T$ . Moreover, in view of (163), probability event  $E_T$  implies  $\eta_t = x^t - x^* - \gamma F(x^t)$  for all  $t = 0, 1, \dots, T$ . Therefore,  $E_T$  implies

$$R_{T+1}^2 \leq R^2 + 2\gamma \sum_{t=0}^T \langle \eta_t, \omega_t \rangle + \gamma^2 \sum_{t=0}^T \|\omega_t\|^2.$$

To continue our derivation we introduce new notation:

$$\omega_t^u \stackrel{\text{def}}{=} \mathbb{E}_{\xi^t} \left[ \tilde{F}_{\xi^t}(x^t) \right] - \tilde{F}_{\xi^t}(x^t), \quad \omega_t^b \stackrel{\text{def}}{=} F(x^t) - \mathbb{E}_{\xi^t} \left[ \tilde{F}_{\xi^t}(x^t) \right] \quad (165)$$

By definition we have  $\omega_t = \omega_t^u + \omega_t^b$  for all  $t = 0, \dots, T$ . Using the introduced notation, we continue our derivation as follows:  $E_T$  implies

$$\begin{aligned} R_{T+1}^2 &\leq R^2 + 2\gamma \underbrace{\sum_{t=0}^T \langle \eta_t, \omega_t^u \rangle}_{\textcircled{1}} + 2\gamma \underbrace{\sum_{t=0}^T \langle \eta_t, \omega_t^b \rangle}_{\textcircled{2}} + 2\gamma^2 \underbrace{\sum_{t=0}^T (\mathbb{E}_{\xi^t} [\|\omega_t^u\|^2])}_{\textcircled{3}} \\ &\quad + 2\gamma^2 \underbrace{\sum_{t=0}^T (\|\omega_t^u\|^2 - \mathbb{E}_{\xi^t} [\|\omega_t^u\|^2])}_{\textcircled{4}} + 2\gamma^2 \underbrace{\sum_{t=0}^T (\|\omega_t^b\|^2)}_{\textcircled{5}}. \end{aligned} \quad (166)$$

We emphasize that the above inequality does not rely on monotonicity of  $F$ .

As we notice above,  $E_T$  implies  $x^t \in B_{2R}(x^*)$  for all  $t = 0, 1, \dots, T$ . This means that the assumptions of Lemma D.1 hold and we have that probability event  $E_T$  implies

$$\begin{aligned} 2\gamma(T+1)\text{Gap}_R(x_{\text{avg}}^T) &\leq \max_{u \in B_R(x^*)} \left\{ \|x^0 - u\|^2 + 2\gamma \sum_{t=0}^T \langle x^t - u - \gamma F(x^t), \omega_t \rangle \right\} \\ &\quad + \gamma^2 \sum_{t=0}^T (\|F(x^t)\|^2 + \|\omega_t\|^2), \\ &= \max_{u \in B_R(x^*)} \left\{ \|x^0 - u\|^2 + 2\gamma \sum_{t=0}^T \langle x^* - u, \omega_t \rangle \right\} \\ &\quad + 2\gamma \sum_{t=0}^T \langle x^t - x^* - \gamma F(x^t), \omega_t \rangle \\ &\quad + \gamma^2 \sum_{t=0}^T (\|F(x^t)\|^2 + \|\omega_t\|^2). \end{aligned}$$

We notice that  $E_T$  implies  $\eta_t = x^t - x^* - \gamma F(x^t)$  for all  $t = 0, 1, \dots, T$  as well as (161) and  $\gamma < 1/\ell$ . Therefore, probability event  $E_T$  implies

$$\begin{aligned} 2\gamma(T+1)\text{Gap}_R(x_{\text{avg}}^T) &\leq \max_{u \in B_R(x^*)} \left\{ \|x^0 - u\|^2 \right\} + 2\gamma \max_{u \in B_R(x^*)} \left\{ \sum_{t=0}^T \langle x^* - u, \omega_t \rangle \right\} \\ &\quad + 2\gamma \sum_{t=0}^T \langle \eta_t, \omega_t \rangle + \frac{\gamma}{\ell} \sum_{t=0}^T \|F(x^t)\|^2 + \gamma^2 \sum_{t=0}^T \|\omega_t\|^2 \\ &\leq 4R^2 + 2\gamma \max_{u \in B_R(x^*)} \left\{ \left\langle x^* - u, \sum_{t=0}^T \omega_t \right\rangle \right\} \\ &\quad + R^2 + 4\gamma \sum_{t=0}^T \langle \eta_t, \omega_t \rangle + 2\gamma^2 \sum_{t=0}^T \|\omega_t\|^2 \\ &\leq 5R^2 + 2\gamma R \left\| \sum_{t=0}^T \omega_t \right\| + 2 \cdot (\textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4} + \textcircled{5}), \end{aligned} \quad (167)$$

where  $\textcircled{1}, \textcircled{2}, \textcircled{3}, \textcircled{4}, \textcircled{5}$  are defined in (166).

The rest of the proof is based on deriving good enough upper bounds for  $\textcircled{1}, \textcircled{2}, \textcircled{3}, \textcircled{4}, \textcircled{5}$ , i.e., we want to prove that  $\textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4} + \textcircled{5} \leq R^2$  and  $2\gamma R \left\| \sum_{t=0}^T \omega_t \right\| \leq 2R^2$  with high probability.

Before we move on, we need to derive some useful inequalities for operating with  $\omega_t^u, \omega_t^b$ . First of all, Lemma B.2 implies that

$$\|\omega_t^u\| \leq 2\lambda \quad (168)$$

for all  $t = 0, 1, \dots, T$ . Next, since  $\{\xi^{i,t}\}_{i=1}^m$  are independently sampled from  $\mathcal{D}$ , we have  $\mathbb{E}_{\xi^t}[F_{\xi^t}(x^t)] = F(x^t)$ , and

$$\mathbb{E}_{\xi^t} [\|F_{\xi^t}(x^t) - F(x^t)\|^2] = \frac{1}{m^2} \sum_{i=1}^m \mathbb{E}_{\xi^{i,t}} [\|F_{\xi^{i,t}}(x^t) - F(x^t)\|^2] \stackrel{(1)}{\leq} \frac{\sigma^2}{m},$$

for all  $l = 0, 1, \dots, T$ . Therefore, in view of Lemma B.2,  $E_T$  implies that

$$\|\omega_t^b\| \leq \frac{4\sigma^2}{m\lambda}, \quad (169)$$

$$\mathbb{E}_{\xi^t} [\|\omega_t\|^2] \leq \frac{18\sigma^2}{m}, \quad (170)$$

$$\mathbb{E}_{\xi^t} [\|\omega_t^u\|^2] \leq \frac{18\sigma^2}{m} \quad (171)$$

for all  $l = 0, 1, \dots, T$ .

**Upper bound for ①.** Since  $\mathbb{E}_{\xi^t}[\omega_t^u] = 0$ , we have

$$\mathbb{E}_{\xi^t} [2\gamma \langle \eta_t, \omega_t^u \rangle] = 0.$$

Next, the summands in ① are bounded with probability 1:

$$|2\gamma \langle \eta_t, \omega_t^u \rangle| \leq 2\gamma \|\eta_t\| \cdot \|\omega_t^u\| \stackrel{(164),(168)}{\leq} 12\gamma R\lambda \stackrel{(157)}{\leq} \frac{R^2}{5 \ln \frac{6(K+1)}{\beta}} \stackrel{\text{def}}{=} c. \quad (172)$$

Moreover, these summands have bounded conditional variances  $\sigma_t^2 \stackrel{\text{def}}{=} \mathbb{E}_{\xi^t} [4\gamma^2 \langle \eta_t, \omega_t^u \rangle^2]$ :

$$\sigma_t^2 \leq \mathbb{E}_{\xi^t} [4\gamma^2 \|\eta_t\|^2 \cdot \|\omega_t^u\|^2] \stackrel{(164)}{\leq} 36\gamma^2 R^2 \mathbb{E}_{\xi^t} [\|\omega_t^u\|^2]. \quad (173)$$

That is, sequence  $\{2\gamma \langle \eta_t, \omega_t^u \rangle\}_{t \geq 0}$  is a bounded martingale difference sequence having bounded conditional variances  $\{\sigma_t^2\}_{t \geq 0}$ . Applying Bernstein's inequality (Lemma B.1) with  $X_t = 2\gamma \langle \eta_t, \omega_t^u \rangle$ ,  $c$  defined in (172),  $b = \frac{R^2}{5}$ ,  $G = \frac{R^4}{150 \ln \frac{6(K+1)}{\beta}}$ , we get that

$$\mathbb{P} \left\{ |\text{①}| > \frac{R^2}{5} \text{ and } \sum_{t=0}^T \sigma_t^2 \leq \frac{R^4}{150 \ln \frac{6(K+1)}{\beta}} \right\} \leq 2 \exp \left( -\frac{b^2}{2G + 2cb/3} \right) = \frac{\beta}{3(K+1)}.$$

In other words,  $\mathbb{P}\{E_{\text{①}}\} \geq 1 - \frac{\beta}{3(K+1)}$ , where probability event  $E_{\text{①}}$  is defined as

$$E_{\text{①}} = \left\{ \text{either } \sum_{t=0}^T \sigma_t^2 > \frac{R^4}{150 \ln \frac{6(K+1)}{\beta}} \text{ or } |\text{①}| \leq \frac{R^2}{5} \right\}. \quad (174)$$

Moreover, we notice here that probability event  $E_T$  implies that

$$\sum_{t=0}^T \sigma_t^2 \stackrel{(173)}{\leq} 36\gamma^2 R^2 \sum_{t=0}^T \mathbb{E}_{\xi^t} [\|\omega_t^u\|^2] \stackrel{(171), T \leq K+1}{\leq} \frac{648\gamma^2 R^2 \sigma^2 (K+1)}{m} \stackrel{(158)}{\leq} \frac{R^4}{150 \ln \frac{6(K+1)}{\beta}}. \quad (175)$$

**Upper bound for ②.** Probability event  $E_T$  implies

$$\begin{aligned} \text{②} &\leq 2\gamma \sum_{t=0}^T \|\eta_t\| \cdot \|\omega_t^b\| \stackrel{(164),(169), T \leq K+1}{\leq} \frac{24\gamma \sigma^2 R (K+1)}{m\lambda} \\ &\stackrel{(157)}{=} \frac{1440\gamma^2 \sigma^2 (K+1) \ln \frac{6(K+1)}{\beta}}{m} \stackrel{(158)}{\leq} \frac{R^2}{5}. \end{aligned} \quad (176)$$



**Upper bound for ③.** Probability event  $E_T$  implies

$$\textcircled{3} = 2\gamma^2 \sum_{t=0}^T \mathbb{E}_{\xi^t} [\|\omega_t^u\|^2] \stackrel{(171), T \leq K+1}{\leq} \frac{36\gamma^2 \sigma^2 (K+1)}{m} \stackrel{(158)}{\leq} \frac{R^2}{5}. \quad (177)$$

**Upper bound for ④.** We have

$$2\gamma^2 \mathbb{E}_{\xi^t} [\|\omega_t^u\|^2 - \mathbb{E}_{\xi^t} [\|\omega_t^u\|^2]] = 0.$$

Next, the summands in ④ are bounded with probability 1:

$$\begin{aligned} 2\gamma^2 \left| \|\omega_t^u\|^2 - \mathbb{E}_{\xi^t} [\|\omega_t^u\|^2] \right| &\leq 2\gamma^2 (\|\omega_t^u\|^2 + \mathbb{E}_{\xi^t} [\|\omega_t^u\|^2]) \stackrel{(168)}{\leq} 16\gamma^2 \lambda^2 \\ &\stackrel{(157)}{\leq} \frac{R^2}{225 \ln \frac{6(K+1)}{\beta}} \leq \frac{R^2}{5 \ln \frac{6(K+1)}{\beta}} \stackrel{\text{def}}{=} c. \end{aligned} \quad (178)$$

Moreover, these summands have bounded conditional variances  $\tilde{\sigma}_t^2 \stackrel{\text{def}}{=} 4\gamma^4 \mathbb{E}_{\xi^t} \left[ \left( \|\omega_t^u\|^2 - \mathbb{E}_{\xi^t} [\|\omega_t^u\|^2] \right)^2 \right]$ :

$$\tilde{\sigma}_t^2 \stackrel{(178)}{\leq} \frac{2\gamma^2 R^2}{225 \ln \frac{6(K+1)}{\beta}} \mathbb{E}_{\xi^t} [\|\omega_t^u\|^2 - \mathbb{E}_{\xi^t} [\|\omega_t^u\|^2]] \leq \frac{4\gamma^2 R^2}{225 \ln \frac{6(K+1)}{\beta}} \mathbb{E}_{\xi^t} [\|\omega_t^u\|^2]. \quad (179)$$

That is, sequence  $\{\|\omega_t^u\|^2 - \mathbb{E}_{\xi^t} [\|\omega_t^u\|^2]\}_{t \geq 0}$  is a bounded martingale difference sequence having bounded conditional variances  $\{\tilde{\sigma}_t^2\}_{t \geq 0}$ . Applying Bernstein's inequality (Lemma B.1) with  $X_t = \|\omega_t^u\|^2 - \mathbb{E}_{\xi^t} [\|\omega_t^u\|^2]$ ,  $c$  defined in (178),  $b = \frac{R^2}{5}$ ,  $G = \frac{R^4}{150 \ln \frac{6(K+1)}{\beta}}$ , we get that

$$\mathbb{P} \left\{ |\textcircled{4}| > \frac{R^2}{5} \text{ and } \sum_{t=0}^T \tilde{\sigma}_t^2 \leq \frac{R^4}{150 \ln \frac{6(K+1)}{\beta}} \right\} \leq 2 \exp \left( -\frac{b^2}{2G + 2cb/3} \right) = \frac{\beta}{3(K+1)}.$$

In other words,  $\mathbb{P}\{E_{\textcircled{4}}\} \geq 1 - \frac{\beta}{3(K+1)}$ , where probability event  $E_{\textcircled{4}}$  is defined as

$$E_{\textcircled{4}} = \left\{ \text{either } \sum_{t=0}^T \tilde{\sigma}_t^2 > \frac{R^4}{150 \ln \frac{6(K+1)}{\beta}} \text{ or } |\textcircled{4}| \leq \frac{R^2}{5} \right\}. \quad (180)$$

Moreover, we notice here that probability event  $E_T$  implies that

$$\begin{aligned} \sum_{t=0}^T \tilde{\sigma}_t^2 &\stackrel{(179)}{\leq} \frac{4\gamma^2 R^2}{225 \ln \frac{6(K+1)}{\beta}} \sum_{t=0}^T \mathbb{E}_{\xi^t} [\|\omega_t^u\|^2] \stackrel{(171), T \leq K+1}{\leq} \frac{8\gamma^2 R^2 \sigma^2 (K+1)}{25m \ln \frac{6(K+1)}{\beta}} \\ &\stackrel{(158)}{\leq} \frac{R^4}{150 \ln \frac{6(K+1)}{\beta}}. \end{aligned} \quad (181)$$

**Upper bound for ⑤.** Probability event  $E_T$  implies

$$\begin{aligned} \textcircled{5} &= 2\gamma^2 \sum_{t=0}^T \|\omega_t^b\|^2 \stackrel{(169), T \leq K+1}{\leq} \frac{32\gamma^2 \sigma^4 (K+1)}{m^2 \lambda^2} \stackrel{(157)}{=} \frac{115200\gamma^4 \sigma^4 (K+1) \ln^2 \frac{6(K+1)}{\beta}}{m^2 R^2} \\ &\stackrel{(158)}{\leq} \frac{R^2}{5}. \end{aligned} \quad (182)$$

**Upper bound for  $\gamma \left\| \sum_{t=0}^T \omega_t \right\|$ .** To handle this term, we introduce new notation:

$$\zeta_t = \begin{cases} \gamma \sum_{r=0}^{l-1} \omega_r, & \text{if } \left\| \gamma \sum_{r=0}^{l-1} \omega_r \right\| \leq R, \\ 0, & \text{otherwise} \end{cases}$$

for  $l = 1, 2, \dots, T - 1$ . By definition, we have

$$\|\zeta_l\| \leq R. \quad (183)$$

Therefore, in view of (160), probability event  $E_T$  implies

$$\begin{aligned} \gamma \left\| \sum_{l=0}^T \omega_l \right\| &= \sqrt{\gamma^2 \left\| \sum_{l=0}^T \omega_l \right\|^2} \\ &= \sqrt{\gamma^2 \sum_{l=0}^T \|\omega_l\|^2 + 2\gamma \sum_{l=0}^T \left\langle \gamma \sum_{r=0}^{l-1} \omega_r, \omega_l \right\rangle} \\ &= \sqrt{\gamma^2 \sum_{l=0}^T \|\omega_l\|^2 + 2\gamma \sum_{l=0}^T \langle \zeta_l, \omega_l \rangle} \\ &\stackrel{(166)}{\leq} \sqrt{\textcircled{3} + \textcircled{4} + \textcircled{5} + \underbrace{2\gamma \sum_{l=0}^T \langle \zeta_l, \omega_l^u \rangle}_{\textcircled{6}} + \underbrace{2\gamma \sum_{l=0}^T \langle \zeta_l, \omega_l^b \rangle}_{\textcircled{7}}}. \end{aligned} \quad (184)$$

Following similar steps as before, we bound  $\textcircled{6}$  and  $\textcircled{7}$ .

**Upper bound for  $\textcircled{6}$ .** Since  $\mathbb{E}_{\xi^t}[\omega_t^u] = 0$ , we have

$$\mathbb{E}_{\xi^t} [2\gamma \langle \zeta_t, \omega_t^u \rangle] = 0.$$

Next, the summands in  $\textcircled{6}$  are bounded with probability 1:

$$|2\gamma \langle \zeta_t, \omega_t^u \rangle| \leq 2\gamma \|\zeta_t\| \cdot \|\omega_t^u\| \stackrel{(183),(168)}{\leq} 4\gamma R \lambda \stackrel{(157)}{\leq} \frac{R^2}{5 \ln \frac{6(K+1)}{\beta}} \stackrel{\text{def}}{=} c. \quad (185)$$

Moreover, these summands have bounded conditional variances  $\hat{\sigma}_t^2 \stackrel{\text{def}}{=} \mathbb{E}_{\xi^t} [4\gamma^2 \langle \zeta_t, \omega_t^u \rangle^2]$ :

$$\hat{\sigma}_t^2 \leq \mathbb{E}_{\xi^t} [4\gamma^2 \|\zeta_t\|^2 \cdot \|\omega_t^u\|^2] \stackrel{(164)}{\leq} 4\gamma^2 R^2 \mathbb{E}_{\xi^t} [\|\omega_t^u\|^2]. \quad (186)$$

That is, sequence  $\{2\gamma \langle \zeta_t, \omega_t^u \rangle\}_{t \geq 0}$  is a bounded martingale difference sequence having bounded conditional variances  $\{\hat{\sigma}_t^2\}_{t \geq 0}$ . Applying Bernstein's inequality (Lemma B.1) with  $X_t = 2\gamma \langle \zeta_t, \omega_t^u \rangle$ ,  $c$  defined in (172),  $b = \frac{R^2}{5}$ ,  $G = \frac{R^4}{150 \ln \frac{6(K+1)}{\beta}}$ , we get that

$$\mathbb{P} \left\{ |\textcircled{6}| > \frac{R^2}{5} \text{ and } \sum_{t=0}^T \hat{\sigma}_t^2 \leq \frac{R^4}{150 \ln \frac{6(K+1)}{\beta}} \right\} \leq 2 \exp \left( -\frac{b^2}{2G + 2cb/3} \right) = \frac{\beta}{3(K+1)}.$$

In other words,  $\mathbb{P}\{E_{\textcircled{6}}\} \geq 1 - \frac{\beta}{3(K+1)}$ , where probability event  $E_{\textcircled{6}}$  is defined as

$$E_{\textcircled{6}} = \left\{ \text{either } \sum_{t=0}^T \hat{\sigma}_t^2 > \frac{R^4}{150 \ln \frac{6(K+1)}{\beta}} \text{ or } |\textcircled{6}| \leq \frac{R^2}{5} \right\}. \quad (187)$$

Moreover, we notice here that probability event  $E_T$  implies that

$$\sum_{t=0}^T \hat{\sigma}_t^2 \stackrel{(186)}{\leq} 4\gamma^2 R^2 \sum_{t=0}^T \mathbb{E}_{\xi^t} [\|\omega_t^u\|^2] \stackrel{(171), T \leq K+1}{\leq} \frac{72\gamma^2 R^2 \sigma^2 (K+1)}{m} \stackrel{(158)}{\leq} \frac{R^4}{150 \ln \frac{6(K+1)}{\beta}}. \quad (188)$$

**Upper bound for  $\textcircled{7}$ .** Probability event  $E_T$  implies

$$\begin{aligned} \textcircled{7} &\leq 2\gamma \sum_{t=0}^T \|\zeta_t\| \cdot \|\omega_t^b\| \stackrel{(183),(169), T \leq K+1}{\leq} \frac{8\gamma \sigma^2 R (K+1)}{m\lambda} \\ &\stackrel{(157)}{=} \frac{480\gamma^2 \sigma^2 (K+1) \ln \frac{6(K+1)}{\beta}}{m} \stackrel{(158)}{\leq} \frac{R^2}{5}. \end{aligned} \quad (189)$$

**Final derivation.** Putting all bounds together, we get that  $E_T$  implies

$$\begin{aligned}
R_{T+1}^2 &\stackrel{(166)}{\leq} R^2 + \textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4} + \textcircled{5}, \\
2\gamma(T+1)\text{Gap}_R(x_{\text{avg}}^T) &\stackrel{(167)}{\leq} 5R^2 + 2\gamma R \left\| \sum_{t=0}^T \omega_t \right\| + 2 \cdot (\textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4} + \textcircled{5}), \\
\gamma \left\| \sum_{l=0}^T \omega_l \right\| &\stackrel{(184)}{\leq} \sqrt{\textcircled{3} + \textcircled{4} + \textcircled{5} + \textcircled{6} + \textcircled{7}}, \\
\textcircled{2} &\stackrel{(176)}{\leq} \frac{R^2}{5}, \quad \textcircled{3} \stackrel{(177)}{\leq} \frac{R^2}{5}, \quad \textcircled{5} \stackrel{(182)}{\leq} \frac{R^2}{5}, \quad \textcircled{7} \stackrel{(189)}{\leq} \frac{R^2}{5}, \\
\sum_{t=0}^T \sigma_t^2 &\stackrel{(175)}{\leq} \frac{R^4}{150 \ln \frac{6(K+1)}{\beta}}, \quad \sum_{t=0}^T \tilde{\sigma}_t^2 \stackrel{(181)}{\leq} \frac{R^4}{150 \ln \frac{6(K+1)}{\beta}}, \quad \sum_{t=0}^T \hat{\sigma}_t^2 \stackrel{(188)}{\leq} \frac{R^4}{150 \ln \frac{6(K+1)}{\beta}}.
\end{aligned}$$

Moreover, in view of (174), (180), (189), and our induction assumption, we have

$$\begin{aligned}
\mathbb{P}\{E_T\} &\geq 1 - \frac{T\beta}{K+1}, \\
\mathbb{P}\{E_{\textcircled{1}}\} &\geq 1 - \frac{\beta}{3(K+1)}, \quad \mathbb{P}\{E_{\textcircled{4}}\} \geq 1 - \frac{\beta}{3(K+1)}, \quad \mathbb{P}\{E_{\textcircled{6}}\} \geq 1 - \frac{\beta}{3(K+1)},
\end{aligned}$$

where probability events  $E_{\textcircled{1}}$ ,  $E_{\textcircled{4}}$ , and  $E_{\textcircled{6}}$  are defined as

$$\begin{aligned}
E_{\textcircled{1}} &= \left\{ \text{either } \sum_{t=0}^T \sigma_t^2 > \frac{R^4}{150 \ln \frac{6(K+1)}{\beta}} \text{ or } |\textcircled{1}| \leq \frac{R^2}{5} \right\}, \\
E_{\textcircled{4}} &= \left\{ \text{either } \sum_{t=0}^T \tilde{\sigma}_t^2 > \frac{R^4}{150 \ln \frac{6(K+1)}{\beta}} \text{ or } |\textcircled{4}| \leq \frac{R^2}{5} \right\}, \\
E_{\textcircled{6}} &= \left\{ \text{either } \sum_{t=0}^T \hat{\sigma}_t^2 > \frac{R^4}{150 \ln \frac{6(K+1)}{\beta}} \text{ or } |\textcircled{6}| \leq \frac{R^2}{5} \right\}.
\end{aligned}$$

Putting all of these inequalities together, we obtain that probability event  $E_T \cap E_{\textcircled{1}} \cap E_{\textcircled{4}} \cap E_{\textcircled{6}}$  implies

$$\begin{aligned}
R_{T+1}^2 &\leq R^2 + \textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4} + \textcircled{5} \leq 2R^2, \\
\gamma \left\| \sum_{l=0}^T \omega_l \right\| &\leq \sqrt{\textcircled{3} + \textcircled{4} + \textcircled{5} + \textcircled{6} + \textcircled{7}} \leq R, \\
2\gamma(T+1)\text{Gap}_R(x_{\text{avg}}^T) &\leq 5R^2 + 2\gamma R \left\| \sum_{t=0}^T \omega_t \right\| + 2 \cdot (\textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4} + \textcircled{5}) \\
&\leq 9R^2.
\end{aligned}$$

Moreover, union bound for the probability events implies

$$\mathbb{P}\{E_{T+1}\} \geq \mathbb{P}\{E_T \cap E_{\textcircled{1}} \cap E_{\textcircled{4}} \cap E_{\textcircled{6}}\} = 1 - \mathbb{P}\{\bar{E}_T \cup \bar{E}_{\textcircled{1}} \cup \bar{E}_{\textcircled{4}} \cup \bar{E}_{\textcircled{6}}\} \geq 1 - \frac{T\beta}{K+1}.$$

This is exactly what we wanted to prove (see the paragraph after inequality (160)). In particular,  $E_K$  implies

$$\text{Gap}_R(x_{\text{avg}}^K) \leq \frac{9R^2}{2\gamma(K+1)},$$

which finishes the proof.  $\square$

**Corollary D.1.** *Let the assumptions of Theorem D.1 hold. Then, the following statements hold.*

1. **Large stepsize/large batch.** The choice of stepsize and batchsize

$$\gamma = \frac{1}{170\ell \ln \frac{6(K+1)}{\beta}}, \quad m = \max \left\{ 1, \frac{972(K+1)\sigma^2}{289\ell^2 R^2 \ln \frac{6(K+1)}{\beta}} \right\} \quad (190)$$

satisfies conditions (156) and (158). With such choice of  $\gamma$ ,  $m$ , and the choice of  $\lambda$  as in (157), the iterates produced by **clipped-SGDA** after  $K$  iterations with probability at least  $1 - \beta$  satisfy

$$\text{Gap}(x_{\text{avg}}^K) \leq \frac{765\ell R^2 \ln \frac{6(K+1)}{\beta}}{K+1}. \quad (191)$$

In particular, to guarantee  $\text{Gap}(x_{\text{avg}}^K) \leq \varepsilon$  with probability at least  $1 - \beta$  for some  $\varepsilon > 0$  **clipped-SGDA** requires,

$$\mathcal{O} \left( \frac{\ell R^2}{\varepsilon} \ln \left( \frac{\ell R^2}{\varepsilon \beta} \right) \right) \text{ iterations}, \quad (192)$$

$$\mathcal{O} \left( \max \left\{ \frac{\ell R^2}{\varepsilon}, \frac{\sigma^2 R^2}{\varepsilon^2} \right\} \ln \left( \frac{\ell R^2}{\varepsilon \beta} \right) \right) \text{ oracle calls}. \quad (193)$$

2. **Small stepsize/small batch.** The choice of stepsize and batchsize

$$\gamma = \min \left\{ \frac{1}{170\ell \ln \frac{6(K+1)}{\beta}}, \frac{R}{180\sigma \sqrt{3(K+1) \ln \frac{6(K+1)}{\beta}}} \right\}, \quad m = 1 \quad (194)$$

satisfies conditions (156) and (158). With such choice of  $\gamma$ ,  $m$ , and the choice of  $\lambda$  as in (157), the iterates produced by **clipped-SGDA** after  $K$  iterations with probability at least  $1 - \beta$  satisfy

$$\text{Gap}(x_{\text{avg}}^K) \leq \max \left\{ \frac{765\ell R^2 \ln \frac{6(K+1)}{\beta}}{K+1}, \frac{810\sigma R \sqrt{3 \ln \frac{6(K+1)}{\beta}}}{\sqrt{K+1}} \right\}. \quad (195)$$

In particular, to guarantee  $\text{Gap}(x_{\text{avg}}^K) \leq \varepsilon$  with probability at least  $1 - \beta$  for some  $\varepsilon > 0$ , **clipped-SGDA** requires

$$\mathcal{O} \left( \max \left\{ \frac{\ell R^2}{\varepsilon} \ln \left( \frac{\ell R^2}{\varepsilon \beta} \right), \frac{\sigma^2 R^2}{\varepsilon^2} \ln \left( \frac{\sigma^2 R^2}{\varepsilon^2 \beta} \right) \right\} \right) \text{ iterations/oracle calls}. \quad (196)$$

*Proof.* 1. **Large stepsize/large batch.** First of all, we verify that the choice of  $\gamma$  and  $m$  from (190) satisfies conditions (156) and (158): (156) trivially holds and (158) holds since

$$m = \max \left\{ 1, \frac{972(K+1)\sigma^2}{289\ell^2 R^2 \ln \frac{6(K+1)}{\beta}} \right\} = \max \left\{ 1, \frac{97200(K+1)\gamma^2 \sigma^2 \ln \frac{6(K+1)}{\beta}}{R^2} \right\}.$$

Therefore, applying Theorem D.1, we derive that with probability at least  $1 - \beta$

$$\text{Gap}(x_{\text{avg}}^K) \leq \frac{9R^2}{2\gamma(K+1)} \stackrel{(190)}{\leq} \frac{765\ell R^2 \ln \frac{6(K+1)}{\beta}}{K+1}.$$

To guarantee  $\text{Gap}(x_{\text{avg}}^K) \leq \varepsilon$ , we choose  $K$  in such a way that the right-hand side of the above inequality is smaller than  $\varepsilon$  that gives

$$K = \mathcal{O} \left( \frac{\ell R^2}{\varepsilon} \ln \left( \frac{\ell R^2}{\varepsilon \beta} \right) \right).$$

The total number of oracle calls equals

$$\begin{aligned} m(K+1) &\stackrel{(190)}{=} \max \left\{ K+1, \frac{972(K+1)^2 \sigma^2}{289\ell^2 R^2 \ln \frac{6(K+1)}{\beta}} \right\} \\ &= \mathcal{O} \left( \max \left\{ \frac{\ell R^2}{\varepsilon}, \frac{\sigma^2 R^2}{\varepsilon^2} \right\} \ln \left( \frac{\ell R^2}{\varepsilon \beta} \right) \right). \end{aligned}$$

2. **Small stepsize/small batch.** First of all, we verify that the choice of  $\gamma$  and  $m$  from (194) satisfies conditions (156) and (158):

$$\gamma = \min \left\{ \frac{1}{170\ell \ln \frac{6(K+1)}{\beta}}, \frac{R}{180\sigma \sqrt{3(K+1) \ln \frac{6(K+1)}{\beta}}} \right\} \leq \frac{1}{170\ell \ln \frac{6(K+1)}{\beta}},$$

$$m = 1 \stackrel{(194)}{\geq} \frac{97200(K+1)\gamma^2\sigma^2 \ln \frac{6(K+1)}{\beta}}{R^2}.$$

Therefore, applying Theorem D.1, we derive that with probability at least  $1 - \beta$

$$\begin{aligned} \text{Gap}(x_{\text{avg}}^K) &\leq \frac{9R^2}{2\gamma(K+1)} \\ &\stackrel{(194)}{=} \max \left\{ \frac{765\ell R^2 \ln \frac{6(K+1)}{\beta}}{K+1}, \frac{810\sigma R \sqrt{3 \ln \frac{6(K+1)}{\beta}}}{\sqrt{K+1}} \right\}. \end{aligned}$$

To guarantee  $\text{Gap}(x_{\text{avg}}^K) \leq \varepsilon$ , we choose  $K$  in such a way that the right-hand side of the above inequality is smaller than  $\varepsilon$  that gives

$$K = \mathcal{O} \left( \max \left\{ \frac{\ell R^2}{\varepsilon} \ln \left( \frac{\ell R^2}{\varepsilon \beta} \right), \frac{\sigma^2 R^2}{\varepsilon^2} \ln \left( \frac{\sigma^2 R^2}{\varepsilon^2 \beta} \right) \right\} \right).$$

The total number of oracle calls equals  $K + 1$ .

□

## D.2 Star-Cocoercive Case

**Theorem D.2.** Let Assumptions 1.1, 1.6, hold for  $Q = B_{2R}(x^*)$ , where  $R \geq R_0 \stackrel{\text{def}}{=} \|x^0 - x^*\|$ , and

$$\gamma \leq \frac{1}{170\ell \ln \frac{4(K+1)}{\beta}}, \quad (197)$$

$$\lambda = \frac{R}{60\gamma \ln \frac{4(K+1)}{\beta}}, \quad (198)$$

$$m \geq \max \left\{ 1, \frac{97200(K+1)\gamma^2\sigma^2 \ln \frac{4(K+1)}{\beta}}{R^2} \right\}, \quad (199)$$

for some  $K \geq 0$  and  $\beta \in (0, 1]$  such that  $\ln \frac{4(K+1)}{\beta} \geq 1$ . Then, after  $K$  iterations the iterates produced by [clipped-SGDA](#) with probability at least  $1 - \beta$  satisfy

$$\frac{1}{K+1} \sum_{k=0}^K \|F(x^k)\|^2 \leq \frac{2\ell R^2}{\gamma(K+1)}. \quad (200)$$

*Proof.* We introduce new notation:  $R_k = \|x^k - x^*\|$  for all  $k \geq 0$ . The proof is based on deriving via induction that  $R_k^2 \leq CR^2$  for some numerical constant  $C > 0$ . In particular, for each  $k = 0, \dots, K+1$  we define probability event  $E_k$  as follows: inequalities

$$\|x^t - x^*\|^2 \leq 2R^2, \quad (201)$$

hold for  $t = 0, 1, \dots, k$  simultaneously. Our goal is to prove that  $\mathbb{P}\{E_k\} \geq 1 - k\beta/(K+1)$  for all  $k = 0, 1, \dots, K+1$ . We notice that inequalities (161) and (166) are derived without assuming monotonicity of  $F$ . Therefore, following exactly the same step as in the proof of Theorem D.1 (up to

the replacement of  $\ln \frac{6(K+1)}{\beta}$  by  $\ln \frac{4(K+1)}{\beta}$ ), we get that

$$\begin{aligned} R_{T+1}^2 &\stackrel{(166)}{\leq} R^2 + \textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4} + \textcircled{5}, \\ \textcircled{2} &\stackrel{(176)}{\leq} \frac{R^2}{5}, \quad \textcircled{3} \stackrel{(177)}{\leq} \frac{R^2}{5}, \quad \textcircled{5} \stackrel{(182)}{\leq} \frac{R^2}{5}, \\ \sum_{t=0}^T \sigma_t^2 &\stackrel{(175)}{\leq} \frac{R^4}{150 \ln \frac{4(K+1)}{\beta}}, \quad \sum_{t=0}^T \tilde{\sigma}_t^2 \stackrel{(181)}{\leq} \frac{R^4}{150 \ln \frac{4(K+1)}{\beta}}. \end{aligned}$$

Moreover, in view of (174), (180), and our induction assumption, we have

$$\begin{aligned} \mathbb{P}\{E_T\} &\geq 1 - \frac{T\beta}{K+1}, \\ \mathbb{P}\{E_{\textcircled{1}}\} &\geq 1 - \frac{\beta}{2(K+1)}, \quad \mathbb{P}\{E_{\textcircled{4}}\} \geq 1 - \frac{\beta}{2(K+1)}, \end{aligned}$$

where probability events  $E_{\textcircled{1}}$ , and  $E_{\textcircled{4}}$  are defined as

$$\begin{aligned} E_{\textcircled{1}} &= \left\{ \text{either } \sum_{t=0}^T \sigma_t^2 > \frac{R^4}{150 \ln \frac{4(K+1)}{\beta}} \quad \text{or } |\textcircled{1}| \leq \frac{R^2}{5} \right\}, \\ E_{\textcircled{4}} &= \left\{ \text{either } \sum_{t=0}^T \tilde{\sigma}_t^2 > \frac{R^4}{150 \ln \frac{4(K+1)}{\beta}} \quad \text{or } |\textcircled{4}| \leq \frac{R^2}{5} \right\}. \end{aligned}$$

Putting all of these inequalities together, we obtain that probability event  $E_{T-1} \cap E_{\textcircled{1}} \cap E_{\textcircled{4}}$  implies

$$R_{T+1}^2 \leq R^2 + \textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4} + \textcircled{5} \leq 2R^2.$$

Moreover, union bound for the probability events implies

$$\mathbb{P}\{E_{T+1}\} \geq \mathbb{P}\{E_T \cap E_{\textcircled{1}} \cap E_{\textcircled{4}}\} = 1 - \mathbb{P}\{\bar{E}_T \cup \bar{E}_{\textcircled{1}} \cup \bar{E}_{\textcircled{4}}\} \geq 1 - \frac{T\beta}{K+1}. \quad (202)$$

This is exactly what we wanted to prove (see the paragraph after inequality (201)). In particular,  $E_K$  implies

$$\begin{aligned} \frac{1}{K+1} \sum_{k=0}^K \|F(x^k)\|^2 &\stackrel{(161)}{\leq} \frac{\ell(R^2 - R_{K+1}^2)}{\gamma(K+1)} + \frac{\ell(\textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4} + \textcircled{5})}{\gamma(K+1)} \\ &\leq \frac{2\ell R^2}{\gamma(K+1)}. \end{aligned}$$

This finishes the proof. □

**Corollary D.2.** *Let the assumptions of Theorem D.2 hold. Then, the following statements hold.*

1. **Large stepsize/large batch.** *The choice of stepsize and batchsize*

$$\gamma = \frac{1}{170\ell \ln \frac{4(K+1)}{\beta}}, \quad m = \max \left\{ 1, \frac{972(K+1)\sigma^2}{289\ell^2 R^2 \ln \frac{4(K+1)}{\beta}} \right\} \quad (203)$$

*satisfies conditions (197) and (199). With such choice of  $\gamma$ ,  $m$ , and the choice of  $\lambda$  as in (198), the iterates produced by clipped-SGDA after  $K$  iterations with probability at least  $1 - \beta$  satisfy*

$$\frac{1}{K+1} \sum_{k=0}^K \|F(x^k)\|^2 \leq \frac{340\ell^2 R^2 \ln \frac{4(K+1)}{\beta}}{K+1}. \quad (204)$$

*In particular, to guarantee  $\frac{1}{K+1} \sum_{k=0}^K \|F(x^k)\|^2 \leq \varepsilon$  with probability at least  $1 - \beta$  for some  $\varepsilon > 0$  clipped-SGDA requires,*

$$\mathcal{O} \left( \frac{\ell^2 R^2}{\varepsilon} \ln \left( \frac{\ell^2 R^2}{\varepsilon \beta} \right) \right) \text{ iterations}, \quad (205)$$

$$\mathcal{O} \left( \max \left\{ \frac{\ell^2 R^2}{\varepsilon}, \frac{\ell^2 \sigma^2 R^2}{\varepsilon^2} \right\} \ln \left( \frac{\ell^2 R^2}{\varepsilon \beta} \right) \right) \text{ oracle calls}. \quad (206)$$

2. **Small stepsize/small batch.** The choice of stepsize and batchsize

$$\gamma = \min \left\{ \frac{1}{170\ell \ln \frac{4(K+1)}{\beta}}, \frac{R}{180\sigma \sqrt{3(K+1) \ln \frac{4(K+1)}{\beta}}} \right\}, \quad m = 1 \quad (207)$$

satisfies conditions (197) and (199). With such choice of  $\gamma$ ,  $m$ , and the choice of  $\lambda$  as in (198), the iterates produced by **clipped-SGDA** after  $K$  iterations with probability at least  $1 - \beta$  satisfy

$$\frac{1}{K+1} \sum_{k=0}^K \|F(x^k)\|^2 \leq \max \left\{ \frac{340\ell^2 R^2 \ln \frac{4(K+1)}{\beta}}{K+1}, \frac{360\ell\sigma R \sqrt{3 \ln \frac{4(K+1)}{\beta}}}{\sqrt{K+1}} \right\}. \quad (208)$$

In particular, to guarantee  $\frac{1}{K+1} \sum_{k=0}^K \|F(x^k)\|^2 \leq \varepsilon$  with probability at least  $1 - \beta$  for some  $\varepsilon > 0$ , **clipped-SGDA** requires

$$\mathcal{O} \left( \max \left\{ \frac{\ell^2 R^2}{\varepsilon} \ln \left( \frac{\ell^2 R^2}{\varepsilon \beta} \right), \frac{\ell^2 \sigma^2 R^2}{\varepsilon^2} \ln \left( \frac{\ell^2 \sigma^2 R^2}{\varepsilon^2 \beta} \right) \right\} \right) \text{ iterations/oracle calls.} \quad (209)$$

*Proof.* 1. **Large stepsize/large batch.** First of all, we verify that the choice of  $\gamma$  and  $m$  from (203) satisfies conditions (197) and (199): (197) trivially holds and (199) holds since

$$m = \max \left\{ 1, \frac{972(K+1)\sigma^2}{289\ell^2 R^2 \ln \frac{4(K+1)}{\beta}} \right\} = \max \left\{ 1, \frac{97200(K+1)\gamma^2 \sigma^2 \ln \frac{4(K+1)}{\beta}}{R^2} \right\}.$$

Therefore, applying Theorem D.2, we derive that with probability at least  $1 - \beta$

$$\frac{1}{K+1} \sum_{k=0}^K \|F(x^k)\|^2 \leq \frac{2\ell R^2}{\gamma(K+1)} \stackrel{(203)}{\leq} \frac{340\ell^2 R^2 \ln \frac{4(K+1)}{\beta}}{K+1}.$$

To guarantee  $\frac{1}{K+1} \sum_{k=0}^K \|F(x^k)\|^2 \leq \varepsilon$ , we choose  $K$  in such a way that the right-hand side of the above inequality is smaller than  $\varepsilon$  that gives

$$K = \mathcal{O} \left( \frac{\ell^2 R^2}{\varepsilon} \ln \left( \frac{\ell^2 R^2}{\varepsilon \beta} \right) \right).$$

The total number of oracle calls equals

$$\begin{aligned} m(K+1) &\stackrel{(203)}{=} \max \left\{ K+1, \frac{972(K+1)^2 \sigma^2}{289\ell^2 R^2 \ln \frac{4(K+1)}{\beta}} \right\} \\ &= \mathcal{O} \left( \max \left\{ \frac{\ell^2 R^2}{\varepsilon}, \frac{\ell^2 \sigma^2 R^2}{\varepsilon^2} \right\} \ln \left( \frac{\ell^2 R^2}{\varepsilon \beta} \right) \right). \end{aligned}$$

2. **Small stepsize/small batch.** First of all, we verify that the choice of  $\gamma$  and  $m$  from (207) satisfies conditions (197) and (199):

$$\gamma = \min \left\{ \frac{1}{170\ell \ln \frac{4(K+1)}{\beta}}, \frac{R}{180\sigma \sqrt{3(K+1) \ln \frac{4(K+1)}{\beta}}} \right\} \leq \frac{1}{170\ell \ln \frac{4(K+1)}{\beta}},$$

$$m = 1 \stackrel{(207)}{\geq} \frac{97200(K+1)\gamma^2 \sigma^2 \ln \frac{4(K+1)}{\beta}}{R^2}.$$

Therefore, applying Theorem D.2, we derive that with probability at least  $1 - \beta$

$$\begin{aligned} \frac{1}{K+1} \sum_{k=0}^K \|F(x^k)\|^2 &\leq \frac{2\ell R^2}{\gamma(K+1)} \\ &\stackrel{(207)}{=} \max \left\{ \frac{340\ell^2 R^2 \ln \frac{4(K+1)}{\beta}}{K+1}, \frac{360\ell\sigma R \sqrt{3 \ln \frac{4(K+1)}{\beta}}}{\sqrt{K+1}} \right\}. \end{aligned}$$

To guarantee  $\frac{1}{K+1} \sum_{k=0}^K \|F(x^k)\|^2 \leq \varepsilon$ , we choose  $K$  in such a way that the right-hand side of the above inequality is smaller than  $\varepsilon$  that gives

$$K = \mathcal{O} \left( \max \left\{ \frac{\ell^2 R^2}{\varepsilon} \ln \left( \frac{\ell^2 R^2}{\varepsilon \beta} \right), \frac{\ell^2 \sigma^2 R^2}{\varepsilon^2} \ln \left( \frac{\ell^2 \sigma^2 R^2}{\varepsilon^2 \beta} \right) \right\} \right).$$

The total number of oracle calls equals  $K + 1$ . □

### D.3 Quasi-Strongly Monotone Star-Cocoercive Case

**Lemma D.3.** *Let Assumptions 1.5, 1.6 hold for  $Q = B_{2R}(x^*)$ , where  $R \geq R_0 \stackrel{\text{def}}{=} \|x^0 - x^*\|$ , and  $0 < \gamma \leq 1/\ell$ . If  $x^k$  lies in  $B_{2R}(x^*)$  for all  $k = 0, 1, \dots, K$  for some  $K \geq 0$ , then the iterates produced by **clipped-SGDA** satisfy*

$$\begin{aligned} \|x^{K+1} - x^*\|^2 &\leq (1 - \gamma\mu)^{K+1} \|x^0 - x^*\|^2 + 2\gamma \sum_{k=0}^K (1 - \gamma\mu)^{K-k} \langle x^k - x^* - \gamma F(x^k), \omega_k \rangle \\ &\quad + \gamma^2 \sum_{k=0}^K (1 - \gamma\mu)^{K-k} \|\omega_k\|^2, \end{aligned} \quad (210)$$

where  $\omega_k$  is defined in (154).

*Proof.* Using the update rule of **clipped-SGDA**, we obtain

$$\begin{aligned} \|x^{k+1} - x^*\|^2 &= \|x^k - x^*\|^2 - 2\gamma \langle x^k - x^*, \tilde{F}_{\xi^k}(x^k) \rangle + \gamma^2 \|\tilde{F}_{\xi^k}(x^k)\|^2 \\ &= \|x^k - x^*\|^2 - 2\gamma \langle x^k - x^*, F(x^k) \rangle + 2\gamma \langle x^k - x^*, \omega_k \rangle \\ &\quad + \gamma^2 \|F(x^k)\|^2 - 2\gamma^2 \langle F(x^k), \omega_k \rangle + \gamma^2 \|\omega_k\|^2 \\ &= \|x^k - x^*\|^2 + 2\gamma \langle x^k - x^* - \gamma F(x^k), \omega_k \rangle \\ &\quad - 2\gamma \langle x^k - x^*, F(x^k) \rangle + \gamma^2 \|F(x^k)\|^2 + \gamma^2 \|\omega_k\|^2 \\ &\stackrel{\text{(SC)}}{\leq} \|x^k - x^*\|^2 + 2\gamma \langle x^k - x^* - \gamma F(x^k), \omega_k \rangle \\ &\quad - 2\gamma \langle x^k - x^*, F(x^k) \rangle + \gamma^2 \ell \langle x^k - x^*, F(x^k) \rangle + \gamma^2 \|\omega_k\|^2 \\ &= \|x^k - x^*\|^2 + 2\gamma \langle x^k - x^* - \gamma F(x^k), \omega_k \rangle \\ &\quad - 2\gamma \left( 1 - \frac{\gamma\ell}{2} \right) \langle x^k - x^*, F(x^k) \rangle + \gamma^2 \|\omega_k\|^2 \\ &\stackrel{\text{(QSM)}, \gamma \leq \frac{1}{\ell}}{\leq} \|x^k - x^*\|^2 + 2\gamma \langle x^k - x^* - \gamma F(x^k), \omega_k \rangle \\ &\quad - 2\gamma\mu \left( 1 - \frac{\gamma\ell}{2} \right) \|x^k - x^*\|^2 + \gamma^2 \|\omega_k\|^2 \\ &\stackrel{\gamma \leq \frac{1}{\ell}}{\leq} (1 - \gamma\mu) \|x^k - x^*\|^2 + 2\gamma \langle x^k - x^* - \gamma F(x^k), \omega_k \rangle + \gamma^2 \|\omega_k\|^2. \end{aligned}$$

Unrolling the recurrence, we obtain (210). □

**Theorem D.3.** *Let Assumptions 1.1, 1.5, 1.6 hold for  $Q = B_{2R}(x^*) = \{x \in \mathbb{R}^d \mid \|x - x^*\| \leq 2R\}$ , where  $R \geq R_0 \stackrel{\text{def}}{=} \|x^0 - x^*\|$ , and*

$$0 < \gamma \leq \frac{1}{400\ell \ln \frac{4(K+1)}{\beta}}, \quad (211)$$

$$\lambda_k = \frac{\exp(-\gamma\mu(1 + k/2))R}{120\gamma \ln \frac{4(K+1)}{\beta}}, \quad (212)$$

$$m_k \geq \max \left\{ 1, \frac{27000\gamma^2(K+1)\sigma^2 \ln \frac{4(K+1)}{\beta}}{\exp(-\gamma\mu k)R^2} \right\}, \quad (213)$$



for some  $K \geq 0$  and  $\beta \in (0, 1]$  such that  $\ln \frac{4(K+1)}{\beta} \geq 1$ . Then, after  $K$  iterations the iterates produced by **clipped-SGDA** with probability at least  $1 - \beta$  satisfy

$$\|x^{K+1} - x^*\|^2 \leq 2 \exp(-\gamma\mu(K+1))R^2. \quad (214)$$

*Proof.* As in the proof of Theorem D.1, we use the following notation:  $R_k = \|x^k - x^*\|^2$ ,  $k \geq 0$ . We will derive (214) by induction. In particular, for each  $k = 0, \dots, K+1$  we define probability event  $E_k$  as follows: inequalities

$$R_t^2 \leq 2 \exp(-\gamma\mu t)R^2 \quad (215)$$

hold for  $t = 0, 1, \dots, k$  simultaneously. Our goal is to prove that  $\mathbb{P}\{E_k\} \geq 1 - k\beta/(K+1)$  for all  $k = 0, 1, \dots, K+1$ . We use the induction to show this statement. For  $k = 0$  the statement is trivial since  $R_0^2 \leq 2R^2$  by definition. Next, assume that the statement holds for  $k = T \leq K$ , i.e., we have  $\mathbb{P}\{E_T\} \geq 1 - T\beta/(K+1)$ . We need to prove that  $\mathbb{P}\{E_{T+1}\} \geq 1 - (T+1)\beta/(K+1)$ . First of all, since  $R_t^2 \leq 2 \exp(-\gamma\mu t)R^2 \leq 2R^2$ , we have  $x^t \in B_{2R}(x^*)$ . Operator  $F$  is  $\ell$ -star-cocoercive on  $B_{2R}(x^*)$ . Therefore, probability event  $E_T$  implies

$$\|F(x^t)\| \stackrel{\text{(SC)}}{\leq} \ell \|x^t - x^*\| \stackrel{\text{(215)}}{\leq} \sqrt{2}\ell \exp(-\gamma\mu t/2)R \stackrel{\text{(211),(212)}}{\leq} \frac{\lambda_t}{2}. \quad (216)$$

and

$$\|\omega_t\|^2 \stackrel{\text{(5)}}{\leq} 2\|\tilde{F}_\xi(x^t)\|^2 + 2\|F(x^t)\|^2 \stackrel{\text{(216)}}{\leq} \frac{5}{2}\lambda_t^2 \stackrel{\text{(212)}}{\leq} \frac{\exp(-\gamma\mu t)R^2}{4\gamma^2} \quad (217)$$

for all  $t = 0, 1, \dots, T$ .

Applying Lemma D.3 and  $(1 - \gamma\mu)^T \leq \exp(-\gamma\mu T)$ , we get that probability event  $E_T$  implies

$$\begin{aligned} R_T^2 &\leq \exp(-\gamma\mu T)R^2 + 2\gamma \sum_{t=0}^T (1 - \gamma\mu)^{T-t} \langle x^t - x^* - \gamma F(x^t), \omega_t \rangle \\ &\quad + \gamma^2 \sum_{t=0}^T (1 - \gamma\mu)^{T-t} \|\omega_t\|^2. \end{aligned}$$

To estimate the sums in the right-hand side, we introduce new vectors:

$$\eta_t = \begin{cases} x^t - x^* - \gamma F(x^t), & \text{if } \|x^t - x^* - \gamma F(x^t)\| \leq \sqrt{2}(1 + \gamma\ell) \exp(-\gamma\mu t/2)R, \\ 0, & \text{otherwise,} \end{cases} \quad (218)$$

for  $t = 0, 1, \dots, T$ . First of all, we point out that vector  $\eta_t$  is bounded with probability 1, i.e., with probability 1

$$\|\eta_t\| \leq \sqrt{2}(1 + \gamma\ell) \exp(-\gamma\mu t/2)R \quad (219)$$

for all  $t = 0, 1, \dots, T$ . Next, we notice that  $E_T$  implies  $\|F(x^t)\| \leq \sqrt{2}\ell \exp(-\gamma\mu t/2)R$  (due to (216)) for  $t = 0, 1, \dots, T$ , i.e., probability event  $E_T$  implies  $\eta_t = x^t - x^* - \gamma F(x^t)$  for all  $t = 0, 1, \dots, T$ . Therefore,  $E_T$  implies

$$\begin{aligned} R_T^2 &\leq \exp(-\gamma\mu T)R^2 + 2\gamma \sum_{t=0}^T (1 - \gamma\mu)^{T-t} \langle \eta_t, \omega_t \rangle \\ &\quad + \gamma^2 \sum_{t=0}^T (1 - \gamma\mu)^{T-t} \|\omega_t\|^2. \end{aligned}$$

As in the monotone case, to continue the derivation, we introduce vectors  $\omega_t^u, \omega_t^b$  defined as

$$\omega_t^u \stackrel{\text{def}}{=} \mathbb{E}_{\xi^t} [\tilde{F}_{\xi^t}(x^t)] - \tilde{F}_{\xi^t}(x^t), \quad \omega_t^b \stackrel{\text{def}}{=} F(x^t) - \mathbb{E}_{\xi^t} [\tilde{F}_{\xi^t}(x^t)], \quad (220)$$

for all  $t = 0, \dots, T$ . By definition we have  $\omega_t = \omega_t^u + \omega_t^b$  for all  $t = 0, \dots, T$ . Using the introduced notation, we continue our derivation as follows:  $E_T$  implies

$$\begin{aligned}
R_T^2 &\stackrel{(5)}{\leq} \underbrace{\exp(-\gamma\mu T)R^2 + 2\gamma \sum_{t=0}^T (1-\gamma\mu)^{T-t} \langle \eta_t, \omega_t^u \rangle}_{\textcircled{1}} + \underbrace{2\gamma \sum_{t=0}^T (1-\gamma\mu)^{T-t} \langle \eta_t, \omega_t^b \rangle}_{\textcircled{2}} \\
&\quad + \underbrace{2\gamma^2 \sum_{t=0}^T (1-\gamma\mu)^{T-t} \mathbb{E}_{\xi^t} [\|\omega_t^u\|^2]}_{\textcircled{3}} + \underbrace{2\gamma^2 \sum_{t=0}^T (1-\gamma\mu)^{T-t} (\|\omega_t^u\|^2 - \mathbb{E}_{\xi^t} [\|\omega_t^u\|^2])}_{\textcircled{4}} \\
&\quad + \underbrace{2\gamma^2 \sum_{t=0}^T (1-\gamma\mu)^{T-t} (\|\omega_t^b\|^2)}_{\textcircled{5}}. \tag{221}
\end{aligned}$$

The rest of the proof is based on deriving good enough upper bounds for  $\textcircled{1}$ ,  $\textcircled{2}$ ,  $\textcircled{3}$ ,  $\textcircled{4}$ ,  $\textcircled{5}$ , i.e., we want to prove that  $\textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4} + \textcircled{5} \leq \exp(-\gamma\mu T)R^2$  with high probability.

Before we move on, we need to derive some useful inequalities for operating with  $\omega_t^u, \omega_t^b$ . First of all, Lemma B.2 implies that

$$\|\omega_t^u\| \leq 2\lambda_t \tag{222}$$

for all  $t = 0, 1, \dots, T$ . Next, since  $\{\xi^{i,t}\}_{i=1}^{m_t}$  are independently sampled from  $\mathcal{D}$ , we have  $\mathbb{E}_{\xi^t} [F_{\xi^t}(x^t)] = F(x^t)$ , and

$$\mathbb{E}_{\xi^t} [\|F_{\xi^t}(x^t) - F(x^t)\|^2] = \frac{1}{m_t^2} \sum_{i=1}^{m_t} \mathbb{E}_{\xi^{i,t}} [\|F_{\xi^{i,t}}(x^t) - F(x^t)\|^2] \stackrel{(1)}{\leq} \frac{\sigma^2}{m_t},$$

for all  $l = 0, 1, \dots, T$ . Moreover, as we already derived, probability event  $E_T$  implies that  $\|F(x^t)\| \leq \lambda_t/2$  for all  $t = 0, 1, \dots, T$  (see (216)). Therefore, in view of Lemma B.2,  $E_T$  implies that

$$\|\omega_t^b\| \leq \frac{4\sigma^2}{m_t \lambda_t}, \tag{223}$$

$$\mathbb{E}_{\xi^t} [\|\omega_t\|^2] \leq \frac{18\sigma^2}{m_t}, \tag{224}$$

$$\mathbb{E}_{\xi^t} [\|\omega_t^u\|^2] \leq \frac{18\sigma^2}{m_t}, \tag{225}$$

for all  $l = 0, 1, \dots, T$ .

**Upper bound for  $\textcircled{1}$ .** Since  $\mathbb{E}_{\xi^t} [\omega_t^u] = 0$ , we have

$$\mathbb{E}_{\xi^t} [2\gamma(1-\gamma\mu)^{T-t} \langle \eta_t, \omega_t^u \rangle] = 0.$$

Next, the summands in  $\textcircled{1}$  are bounded with probability 1:

$$\begin{aligned}
|2\gamma(1-\gamma\mu)^{T-t} \langle \eta_t, \omega_t^u \rangle| &\leq 2\gamma \exp(-\gamma\mu(T-t)) \|\eta_t\| \cdot \|\omega_t^u\| \\
&\stackrel{(219),(222)}{\leq} 4\sqrt{2}\gamma(1+\gamma\ell) \exp(-\gamma\mu(T-t/2)) R\lambda_t \\
&\stackrel{(211),(212)}{\leq} \frac{\exp(-\gamma\mu T)R^2}{5 \ln \frac{4(K+1)}{\beta}} \stackrel{\text{def}}{=} c. \tag{226}
\end{aligned}$$

Moreover, these summands have bounded conditional variances  $\sigma_t^2 \stackrel{\text{def}}{=} \mathbb{E}_{\xi^t} [4\gamma^2(1-\gamma\mu)^{2T-2t} \langle \eta_t, \omega_t^u \rangle^2]$ :

$$\begin{aligned}
\sigma_t^2 &\leq \mathbb{E}_{\xi^t} [4\gamma^2 \exp(-\gamma\mu(2T-2t)) \|\eta_t\|^2 \cdot \|\omega_t^u\|^2] \\
&\stackrel{(219)}{\leq} 8\gamma^2(1+\gamma\ell)^2 \exp(-\gamma\mu(2T-t)) R^2 \mathbb{E}_{\xi^t} [\|\omega_t^u\|^2] \\
&\stackrel{(211)}{\leq} 10\gamma^2 \exp(-\gamma\mu(2T-t)) R^2 \mathbb{E}_{\xi^t} [\|\omega_t^u\|^2]. \tag{227}
\end{aligned}$$

That is, sequence  $\{2\gamma(1 - \gamma\mu)^{T-t}\langle \eta_t, \omega_t^u \rangle\}_{t \geq 0}$  is a bounded martingale difference sequence having bounded conditional variances  $\{\sigma_t^2\}_{t \geq 0}$ . Applying Bernstein's inequality (Lemma B.1) with  $X_t = 2\gamma(1 - \gamma\mu)^{T-t}\langle \eta_t, \omega_t^u \rangle$ ,  $c$  defined in (226),  $b = \frac{1}{5} \exp(-\gamma\mu T)R^2$ ,  $G = \frac{\exp(-2\gamma\mu T)R^4}{150 \ln \frac{4(K+1)}{\beta}}$ , we get that

$$\mathbb{P} \left\{ |\mathbb{1}| > \frac{1}{5} \exp(-\gamma\mu T)R^2 \text{ and } \sum_{t=0}^T \sigma_t^2 \leq \frac{\exp(-2\gamma\mu T)R^4}{150 \ln \frac{4(K+1)}{\beta}} \right\} \leq 2 \exp \left( -\frac{b^2}{2G + 2cb/3} \right) = \frac{\beta}{2(K+1)}.$$

In other words,  $\mathbb{P}\{E_{\mathbb{1}}\} \geq 1 - \frac{\beta}{2(K+1)}$ , where probability event  $E_{\mathbb{1}}$  is defined as

$$E_{\mathbb{1}} = \left\{ \text{either } \sum_{t=0}^T \sigma_t^2 > \frac{\exp(-2\gamma\mu T)R^4}{150 \ln \frac{4(K+1)}{\beta}} \text{ or } |\mathbb{1}| \leq \frac{1}{5} \exp(-\gamma\mu T)R^2 \right\}. \quad (228)$$

Moreover, we notice here that probability event  $E_T$  implies that

$$\begin{aligned} \sum_{t=0}^T \sigma_t^2 &\stackrel{(227)}{\leq} 10\gamma^2 \exp(-2\gamma\mu T)R^2 \sum_{t=0}^T \frac{\mathbb{E}_{\xi^t} [\|\omega_t^u\|^2]}{\exp(-\gamma\mu t)} \\ &\stackrel{(225), T \leq K+1}{\leq} 180\gamma^2 \exp(-2\gamma\mu T)R^2 \sigma^2 \sum_{t=0}^K \frac{1}{m_t \exp(-\gamma\mu t)} \\ &\stackrel{(213)}{\leq} \frac{\exp(-2\gamma\mu T)R^4}{150 \ln \frac{6(K+1)}{\beta}}. \end{aligned} \quad (229)$$

**Upper bound for ②.** Probability event  $E_T$  implies

$$\begin{aligned} \textcircled{2} &\leq 2\gamma \exp(-\gamma\mu T) \sum_{t=0}^T \frac{\|\eta_t\| \cdot \|\omega_t^b\|}{\exp(-\gamma\mu t)} \\ &\stackrel{(219), (223)}{\leq} 8\sqrt{2}\gamma(1 + \gamma\ell) \exp(-\gamma\mu T)R \sum_{t=0}^T \frac{\sigma^2}{m_t \lambda_t \exp(-\gamma\mu t/2)} \\ &\stackrel{(212)}{\leq} 960\sqrt{2}\gamma^2(1 + \gamma\ell) \exp(-\gamma\mu(T-1)) \sum_{t=0}^T \frac{\sigma^2 \ln \frac{4(K+1)}{\beta}}{m_t \exp(-\gamma\mu t)} \\ &\stackrel{(213), T \leq K+1}{\leq} \frac{1}{5} \exp(-\gamma\mu T)R^2. \end{aligned} \quad (230)$$

**Upper bound for ③.** Probability event  $E_T$  implies

$$\begin{aligned} \textcircled{3} &= 2\gamma^2 \exp(-\gamma\mu T) \sum_{t=0}^T \frac{\mathbb{E}_{\xi^t} [\|\omega_t^u\|^2]}{\exp(-\gamma\mu t)} \\ &\stackrel{(225)}{\leq} 36\gamma^2 \exp(-\gamma\mu T) \sum_{t=0}^T \frac{\sigma^2}{m_t \exp(-\gamma\mu t)} \\ &\stackrel{(213), T \leq K+1}{\leq} \frac{1}{5} \exp(-\gamma\mu T)R^2. \end{aligned} \quad (231)$$

**Upper bound for ④.** First of all, we have

$$2\gamma^2(1 - \gamma\mu)^{T-t} \mathbb{E}_{\xi^t} [\|\omega_t^u\|^2 - \mathbb{E}_{\xi^t} [\|\omega_t^u\|^2]] = 0.$$

Next, the summands in ④ are bounded with probability 1:

$$\begin{aligned} 2\gamma^2(1 - \gamma\mu)^{T-t} \left| \|\omega_t^u\|^2 - \mathbb{E}_{\xi^t} [\|\omega_t^u\|^2] \right| &\stackrel{(222)}{\leq} \frac{16\gamma^2 \exp(-\gamma\mu T) \lambda_t^2}{\exp(-\gamma\mu t)} \\ &\stackrel{(212)}{\leq} \frac{\exp(-\gamma\mu T)R^2}{5 \ln \frac{4(K+1)}{\beta}} \\ &\stackrel{\text{def}}{=} c. \end{aligned} \quad (232)$$

Moreover, these summands have bounded conditional variances  $\tilde{\sigma}_t^2 \stackrel{\text{def}}{=} \mathbb{E}_{\xi^t} \left[ 4\gamma^4(1-\gamma\mu)^{2T-2t} \left( \|\omega_t^u\|^2 - \mathbb{E}_{\xi^t} [\|\omega_t^u\|^2] \right)^2 \right]$ :

$$\begin{aligned} \tilde{\sigma}_t^2 &\stackrel{(232)}{\leq} \frac{2\gamma^2 \exp(-2\gamma\mu T) R^2}{5 \exp(-\gamma\mu t) \ln \frac{4(K+1)}{\beta}} \mathbb{E}_{\xi^t} \left[ \left| \|\omega_t^u\|^2 - \mathbb{E}_{\xi^t} [\|\omega_t^u\|^2] \right| \right] \\ &\leq \frac{4\gamma^2 \exp(-2\gamma\mu T) R^2}{5 \exp(-\gamma\mu t) \ln \frac{4(K+1)}{\beta}} \mathbb{E}_{\xi^t} [\|\omega_t^u\|^2]. \end{aligned} \quad (233)$$

That is, sequence  $\left\{ 2\gamma^2(1-\gamma\mu)^{T-t} \left( \|\omega_t^u\|^2 - \mathbb{E}_{\xi^t} [\|\omega_t^u\|^2] \right) \right\}_{t \geq 0}$  is a bounded martingale difference sequence having bounded conditional variances  $\{\tilde{\sigma}_t^2\}_{t \geq 0}$ . Applying Bernstein's inequality (Lemma B.1) with  $X_t = 2\gamma^2(1-\gamma\mu)^{T-t} \left( \|\omega_t^u\|^2 - \mathbb{E}_{\xi^t} [\|\omega_t^u\|^2] \right)$ ,  $c$  defined in (232),  $b = \frac{1}{5} \exp(-\gamma\mu T) R^2$ ,  $G = \frac{\exp(-2\gamma\mu T) R^4}{150 \ln \frac{4(K+1)}{\beta}}$ , we get that

$$\mathbb{P} \left\{ |\textcircled{4}| > \frac{1}{5} \exp(-\gamma\mu T) R^2 \text{ and } \sum_{t=0}^T \tilde{\sigma}_t^2 \leq \frac{\exp(-2\gamma\mu T) R^4}{150 \ln \frac{4(K+1)}{\beta}} \right\} \leq 2 \exp \left( -\frac{b^2}{2G + 2cb/3} \right) = \frac{\beta}{2(K+1)}.$$

In other words,  $\mathbb{P}\{E_{\textcircled{4}}\} \geq 1 - \frac{\beta}{2(K+1)}$ , where probability event  $E_{\textcircled{4}}$  is defined as

$$E_{\textcircled{4}} = \left\{ \text{either } \sum_{t=0}^T \tilde{\sigma}_t^2 > \frac{\exp(-2\gamma\mu T) R^4}{150 \ln \frac{4(K+1)}{\beta}} \text{ or } |\textcircled{4}| \leq \frac{1}{5} \exp(-\gamma\mu T) R^2 \right\}. \quad (234)$$

Moreover, we notice here that probability event  $E_T$  implies that

$$\begin{aligned} \sum_{t=0}^T \tilde{\sigma}_t^2 &\stackrel{(233)}{\leq} \frac{4\gamma^2 \exp(-2\gamma\mu T) R^2}{5 \ln \frac{4(K+1)}{\beta}} \sum_{t=0}^T \frac{\mathbb{E}_{\xi^t} [\|\omega_t^u\|^2]}{\exp(-\gamma\mu t)} \\ &\stackrel{(225), T \leq K+1}{\leq} \frac{72\gamma^2 \exp(-2\gamma\mu T) R^2 \sigma^2}{5 \ln \frac{4(K+1)}{\beta}} \sum_{t=0}^K \frac{1}{m_t \exp(-\gamma\mu t)} \\ &\stackrel{(213)}{\leq} \frac{\exp(-2\gamma\mu T) R^4}{150 \ln \frac{4(K+1)}{\beta}}. \end{aligned} \quad (235)$$

**Upper bound for ⑤.** Probability event  $E_T$  implies

$$\begin{aligned} \textcircled{5} &= 2\gamma^2 \sum_{t=0}^T \exp(-\gamma\mu(T-t)) \left( \|\omega_t^b\|^2 \right) \\ &\stackrel{(223)}{\leq} 32\gamma^2 \exp(-\gamma\mu T) \sum_{t=0}^T \frac{\sigma^4}{m_t^2 \lambda_t^2 \exp(-\gamma\mu t)} \\ &\stackrel{(212)}{=} 460800\gamma^4 \exp(-\gamma\mu(T-2)) \sum_{t=0}^T \frac{\sigma^4 \ln^2 \frac{4(K+1)}{\beta}}{m_t^2 R^2 \exp(-2\gamma\mu t)} \\ &\stackrel{(213), T \leq K+1}{\leq} \frac{1}{5} \exp(-\gamma\mu T) R^2. \end{aligned} \quad (236)$$

**Final derivation.** Putting all bounds together, we get that  $E_T$  implies

$$\begin{aligned} R_T^2 &\stackrel{(221)}{\leq} \exp(-\gamma\mu T) R^2 + \textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4} + \textcircled{5}, \\ &\stackrel{(230)}{\leq} \frac{1}{5} \exp(-\gamma\mu T) R^2, \\ \textcircled{3} &\stackrel{(231)}{\leq} \frac{1}{5} \exp(-\gamma\mu T) R^2, \quad \textcircled{5} \stackrel{(236)}{\leq} \frac{1}{5} \exp(-\gamma\mu T) R^2, \\ \sum_{t=0}^T \sigma_t^2 &\stackrel{(229)}{\leq} \frac{\exp(-2\gamma\mu T) R^4}{150 \ln \frac{4(K+1)}{\beta}}, \quad \sum_{t=0}^T \tilde{\sigma}_t^2 \stackrel{(235)}{\leq} \frac{\exp(-2\gamma\mu T) R^4}{150 \ln \frac{4(K+1)}{\beta}}. \end{aligned}$$

Moreover, in view of (228), (234), and our induction assumption, we have

$$\begin{aligned} \mathbb{P}\{E_T\} &\geq 1 - \frac{T\beta}{K+1}, \\ \mathbb{P}\{E_{\textcircled{1}}\} &\geq 1 - \frac{\beta}{2(K+1)}, \quad \mathbb{P}\{E_{\textcircled{4}}\} \geq 1 - \frac{\beta}{2(K+1)}, \end{aligned}$$

where probability events  $E_{\textcircled{1}}$ , and  $E_{\textcircled{4}}$  are defined as

$$\begin{aligned} E_{\textcircled{1}} &= \left\{ \text{either } \sum_{t=0}^T \sigma_t^2 > \frac{\exp(-2\gamma\mu T)R^4}{150 \ln \frac{4(K+1)}{\beta}} \text{ or } |\textcircled{1}| \leq \frac{1}{5} \exp(-\gamma\mu T)R^2 \right\}, \\ E_{\textcircled{4}} &= \left\{ \text{either } \sum_{t=0}^T \tilde{\sigma}_t^2 > \frac{\exp(-2\gamma\mu T)R^4}{150 \ln \frac{4(K+1)}{\beta}} \text{ or } |\textcircled{4}| \leq \frac{1}{5} \exp(-\gamma\mu T)R^2 \right\}. \end{aligned}$$

Putting all of these inequalities together, we obtain that probability event  $E_T \cap E_{\textcircled{1}} \cap E_{\textcircled{4}}$  implies

$$\begin{aligned} R_T^2 &\stackrel{(221)}{\leq} \exp(-\gamma\mu T)R^2 + \textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4} + \textcircled{5} \\ &\leq 2 \exp(-\gamma\mu T)R^2. \end{aligned}$$

Moreover, union bound for the probability events implies

$$\mathbb{P}\{E_{T+1}\} \geq \mathbb{P}\{E_T \cap E_{\textcircled{1}} \cap E_{\textcircled{4}}\} = 1 - \mathbb{P}\{\bar{E}_T \cup \bar{E}_{\textcircled{1}} \cup \bar{E}_{\textcircled{4}}\} \geq 1 - \frac{(T+1)\beta}{K+1}. \quad (237)$$

This is exactly what we wanted to prove (see the paragraph after inequality (215)). In particular, with probability at least  $1 - \beta$  we have

$$\|x^{K+1} - x^*\|^2 \leq 2 \exp(-\gamma\mu(K+1))R^2,$$

which finishes the proof.  $\square$

**Corollary D.3.** *Let the assumptions of Theorem D.3 hold. Then, the following statements hold.*

1. **Large stepsize/large batch.** *The choice of stepsize and batchsize*

$$\gamma = \frac{1}{400\ell \ln \frac{4(K+1)}{\beta}}, \quad m_k = \max \left\{ 1, \frac{27000\gamma^2(K+1)\sigma^2 \ln \frac{4(K+1)}{\beta}}{\exp(-\gamma\mu k)R^2} \right\} \quad (238)$$

*satisfies conditions (211) and (213). With such choice of  $\gamma$ ,  $m_k$ , and the choice of  $\lambda_k$  as in (212), the iterates produced by clipped-SGDA after  $K$  iterations with probability at least  $1 - \beta$  satisfy*

$$\|x^{K+1} - x^*\|^2 \leq 2 \exp \left( -\frac{\mu(K+1)}{400\ell \ln \frac{4(K+1)}{\beta}} \right) R^2. \quad (239)$$

*In particular, to guarantee  $\|x^{K+1} - x^*\|^2 \leq \varepsilon$  with probability at least  $1 - \beta$  for some  $\varepsilon > 0$  clipped-SGDA requires*

$$\mathcal{O} \left( \frac{\ell}{\mu} \ln \left( \frac{R^2}{\varepsilon} \right) \ln \left( \frac{\ell}{\mu\beta} \ln \left( \frac{R^2}{\varepsilon} \right) \right) \right) \text{ iterations}, \quad (240)$$

$$\mathcal{O} \left( \max \left\{ \frac{\ell}{\mu}, \frac{\sigma^2}{\mu^2\varepsilon} \right\} \ln \left( \frac{R^2}{\varepsilon} \right) \ln \left( \frac{\ell}{\mu\beta} \ln \left( \frac{R^2}{\varepsilon} \right) \right) \right) \text{ oracle calls}. \quad (241)$$

2. **Small stepsize/small batch.** *The choice of stepsize and batchsize*

$$\gamma = \min \left\{ \frac{1}{400\ell \ln \frac{4(K+1)}{\beta}}, \frac{\ln(B_K)}{\mu(K+1)} \right\}, \quad m_k \equiv 1 \quad (242)$$

satisfies conditions (211) and (213), where  $B_K = \max \left\{ 2, \frac{(K+1)\mu^2 R^2}{27000\sigma^2 \ln\left(\frac{4(K+1)}{\beta}\right) \ln^2(B_K)} \right\} = \mathcal{O} \left( \max \left\{ 2, \frac{(K+1)\mu^2 R^2}{27000\sigma^2 \ln\left(\frac{4(K+1)}{\beta}\right) \ln^2 \left( \max \left\{ 2, \frac{(K+1)\mu^2 R^2}{27000\sigma^2 \ln\left(\frac{4(K+1)}{\beta}\right)} \right\} \right)} \right\} \right)$ . With such choice of

$\gamma, m_k$ , and the choice of  $\lambda_k$  as in (212), the iterates produced by clipped-SGDA after  $K$  iterations with probability at least  $1 - \beta$  satisfy

$$\|x^{K+1} - x^*\|^2 \leq \max \left\{ 2 \exp \left( -\frac{\mu(K+1)}{400\ell \ln \frac{4(K+1)}{\beta}} \right) R^2, \frac{54000\sigma^2 \ln \left( \frac{4(K+1)}{\beta} \right) \ln^2(B_K)}{\mu^2(K+1)} \right\}. \quad (243)$$

In particular, to guarantee  $\|x^{K+1} - x^*\|^2 \leq \varepsilon$  with probability at least  $1 - \beta$  for some  $\varepsilon > 0$  clipped-SGDA requires

$$\mathcal{O} \left( \max \left\{ \frac{\ell}{\mu} \ln \left( \frac{R^2}{\varepsilon} \right) \ln \left( \frac{\ell}{\mu\beta} \ln \left( \frac{R^2}{\varepsilon} \right) \right), \frac{\sigma^2}{\mu^2\varepsilon} \ln \left( \frac{\sigma^2}{\mu^2\varepsilon\beta} \right) \ln^2(B_\varepsilon) \right\} \right) \quad (244)$$

iterations/oracle calls, where

$$B_\varepsilon = \max \left\{ 2, \frac{R^2}{\varepsilon \ln \left( \frac{\sigma^2}{\mu^2\varepsilon\beta} \right) \ln^2 \left( \max \left\{ 2, \frac{R^2}{\varepsilon \ln \left( \frac{\sigma^2}{\mu^2\varepsilon\beta} \right)} \right\} \right)} \right\}.$$

*Proof.* **1. Large stepsize/large batch.** First of all, it is easy to see that the choice of  $\gamma$  and  $m_k$  from (238) satisfies conditions (211) and (213). Therefore, applying Theorem D.3, we derive that with probability at least  $1 - \beta$

$$\|x^{K+1} - x^*\|^2 \leq 2 \exp(-\gamma\mu(K+1))R^2 \stackrel{(238)}{=} 2 \exp \left( -\frac{\mu(K+1)}{400\ell \ln \frac{4(K+1)}{\beta}} \right) R^2.$$

To guarantee  $\|x^{K+1} - x^*\|^2 \leq \varepsilon$ , we choose  $K$  in such a way that the right-hand side of the above inequality is smaller than  $\varepsilon$  that gives

$$K = \mathcal{O} \left( \frac{\ell}{\mu} \ln \left( \frac{R^2}{\varepsilon} \right) \ln \left( \frac{\ell}{\mu\beta} \ln \left( \frac{R^2}{\varepsilon} \right) \right) \right).$$

The total number of oracle calls equals

$$\begin{aligned} \sum_{k=0}^K m_k &\stackrel{(238)}{=} \sum_{k=0}^K \max \left\{ 1, \frac{27000\gamma^2(K+1)\sigma^2 \ln \frac{4(K+1)}{\beta}}{\exp(-\gamma\mu k)R^2} \right\} \\ &= \mathcal{O} \left( \max \left\{ K, \frac{\gamma(K+1) \exp(\gamma\mu(K+1))\sigma^2 \ln \frac{4(K+1)}{\beta}}{\mu R^2} \right\} \right) \\ &= \mathcal{O} \left( \max \left\{ \frac{\ell}{\mu}, \frac{\sigma^2}{\mu^2\varepsilon} \right\} \ln \left( \frac{R^2}{\varepsilon} \right) \ln \left( \frac{\ell}{\mu\beta} \ln \left( \frac{R^2}{\varepsilon} \right) \right) \right). \end{aligned}$$

**2. Small stepsize/small batch.** First of all, we verify that the choice of  $\gamma$  and  $m_k$  from (242) satisfies conditions (211) and (213): (211) trivially holds and (213) holds since for all  $k = 0, \dots, K$

$$\begin{aligned} \frac{27000\gamma^2(K+1)\sigma^2 \ln \frac{4(K+1)}{\beta}}{\exp(-\gamma\mu k)R^2} &\leq \frac{27000\gamma^2(K+1)\sigma^2 \ln \frac{4(K+1)}{\beta}}{\exp(-\gamma\mu(K+1))R^2} \\ &\stackrel{(242)}{\leq} \frac{27000 \ln^2(B_K) \exp(\gamma\mu(K+1))\sigma^2 \ln \frac{4(K+1)}{\beta}}{\mu^2(K+1)R^2} \\ &\stackrel{(242)}{\leq} 1. \end{aligned}$$

Therefore, applying Theorem D.3, we derive that with probability at least  $1 - \beta$

$$\begin{aligned}
\|x^{K+1} - x^*\|^2 &\leq 2 \exp(-\gamma\mu(K+1))R^2 \\
&\stackrel{(242)}{=} \max \left\{ 2 \exp \left( -\frac{\mu(K+1)}{400\ell \ln \frac{4(K+1)}{\beta}} \right) R^2, \frac{2R^2}{B_K} \right\} \\
&= \max \left\{ 2 \exp \left( -\frac{\mu(K+1)}{400\ell \ln \frac{4(K+1)}{\beta}} \right) R^2, \frac{54000\sigma^2 \ln \left( \frac{4(K+1)}{\beta} \right) \ln^2(B_K)}{\mu^2(K+1)} \right\}.
\end{aligned}$$

To guarantee  $\|x^{K+1} - x^*\|^2 \leq \varepsilon$ , we choose  $K$  in such a way that the right-hand side of the above inequality is smaller than  $\varepsilon$  that gives  $K$  of the order

$$\mathcal{O} \left( \max \left\{ \frac{\ell}{\mu} \ln \left( \frac{R^2}{\varepsilon} \right) \ln \left( \frac{\ell}{\mu\beta} \ln \left( \frac{R^2}{\varepsilon} \right) \right), \frac{\sigma^2}{\mu^2\varepsilon} \ln \left( \frac{\sigma^2}{\mu^2\varepsilon\beta} \right) \ln^2(B_\varepsilon) \right\} \right),$$

where

$$B_\varepsilon = \max \left\{ 2, \frac{R^2}{\varepsilon \ln \left( \frac{\sigma^2}{\mu^2\varepsilon\beta} \right) \ln^2 \left( \max \left\{ 2, \frac{R^2}{\varepsilon \ln \left( \frac{\sigma^2}{\mu^2\varepsilon\beta} \right)} \right\} \right)} \right\}.$$

The total number of oracle calls equals  $\sum_{k=0}^K m_k = (K+1)$ .

□

## E Extra Experiments

In this section, we provide more details for the experiments done in § 4, as well as additional tables, figures, and image samples from some of our trained models.

### E.1 WGAN-GP

In all cases, everything in the experimental setup other than learning rates and clip values remained constant. We use the same ResNet architectures and training parameters specified in Gulrajani et al. [2017]: the gradient penalty coefficient  $\lambda_{GP} = 10$ ,  $n_{dis} = 5$  where  $n_{dis}$  is the number of discriminator steps for every generator step, and a learning rate decayed linearly to 0 over 100k steps. The only exception is we double the feature map of the generator from 128 to 256 dimensions. For all stochastic extragradient (SEG) methods, we use the ExtraSGD implementation provided by Gidel et al. [2019a]. We alternate between exploration and update steps and do not treat the exploration steps as “free” – this means we only have 50k parameter updates as opposed to 100k for all SGDA methods (we decay the learning rate twice as fast such that it still reaches 0 after 50k parameter updates).

All of the hyperparameter sweeps performed for SGDA, clipped-SGDA, clipped-SEG, clipped-SGDA (coordinate), and clipped-SEG (coordinate), as well as the associated best FID score obtained within the first 35k training steps, can be found in Tables 2, 6, 5, 6, and 7 respectively. **Bold** rows denote the hyperparameters that were trained for the full 100k steps and are henceforth referred to as the “*best models*”. For each of the methods tested, additional samples for the best models trained can be found in Figures 7, 8, 9, 10, & 11. We also plot the evolution of the gradient noise histograms in Figures 12, 13, 14, 15, & 16. We emphasize that our goal is not to get the best possible FID score (e.g. are often able to obtain marginally better FIDs by training for longer), but rather to compare the systematic differences in performance between the various unclipped and clipped methods. Therefore, log-space hyperparameter sweeps are appropriate for our experiments and we do not tune further.

### E.2 StyleGAN2

We train on FFHQ downsampled to  $128 \times 128$  pixels, and use the recommended StyleGAN2 hyperparameter configuration for this resolution: batch size = 32,  $\gamma = 0.1024$ , map depth = 2, and channel multiplier = 16384. For both SGDA and clipped-SGDA, we sweep over a (roughly) log-scale of learning rates and clipping values; a summary of the hyperparameters and best FID scores obtained Table 8 and Table 9 respectively.

Based on the results in Table 9, the best hyperparameters are lr=0.35 and clip=0.0025 which we then used to train our “*best model*”. We trained for longer, and decayed the learning rate twice (by a factor of  $\times 10$ ) when the FID plateaued or worsened. The best schedule we found was to scale the learning rate by  $\times 0.1$  after 6000 kimgs (thousands of real images shown to the discriminator), by another  $\times 0.1$  after 3600 kimgs, and then train until the FID begins increasing (for another 8000 kimgs) – for a total of 17600 kimgs. We did not explore different scale factors or other schedules (such as cosine annealing). Additional samples for this model can be found in Figure 18(a).

In general, we observe that every SGDA-trained model for the wide range of learning rates we tested failed to improve the FID, while models trained with clipped-SGDA (with appropriately set hyperparameters) are generally able to learn some meaningful features and improve the FID. We show this behaviour in Figure 17 – the FID scores for SGDA-trained models fluctuate around 320 and only generate noise such as the samples shown in Figure 18(b), which is in contrast to models trained with clipped-SGDA. Note that the range of the hyperparameter sweep is fairly narrow and favourable for clipped-SGDA, while being quite wide for SGDA. The purpose for these parameter ranges is not to directly compare the parameter sweeps (which would unfairly favour clipped-SGDA), but to show that in general SGDA fails, while clipped-SGDA is capable of learning.



Table 2: SGDA hyperparameter sweep, and the best FID score obtained in 35k training steps.

G-LR	D-LR	FID
6e-06	6e-06	233.3
2e-05	2e-05	177.2
2e-05	4e-05	183.4
2e-05	8e-05	187.3
0.0002	0.0002	85.6
<b>0.0002</b>	<b>0.0004</b>	<b>82.8</b>
0.0002	0.0008	NaN
0.002	0.002	NaN
0.02	0.02	NaN
0.2	0.2	NaN

Table 3: SEG hyperparameter sweep, and the best FID score obtained in 35k training steps.

G-LR	D-LR	FID
6e-06	6e-06	236.1
2e-05	2e-05	208.6
2e-05	4e-05	213.7
<b>4e-05</b>	<b>4e-05</b>	<b>176.5</b>
4e-05	0.0001	NaN
0.0002	0.0002	NaN
0.0002	0.0004	NaN
0.0002	0.0008	NaN
0.002	0.002	NaN
0.02	0.02	NaN
0.2	0.2	NaN
2	2	NaN

Table 4: clipped-SGDA (norm) hyperparameter sweep, and the best FID score obtained in 35k training steps.

G-LR	D-LR	G-clip	D-clip	FID
0.002	0.002	0.1	0.1	257.6
0.002	0.002	1	1	121.6
0.002	0.002	10	10	145.4
0.02	0.02	0.1	0.1	115.4
0.02	0.02	1	1	141.8
0.02	0.02	10	10	27.4
0.2	0.2	0.1	0.1	133.0
<b>0.2</b>	<b>0.2</b>	<b>1</b>	<b>1</b>	<b>26.3</b>
2	2	0.1	0.1	26.1

Table 5: clipped-SEG (norm) hyperparameter sweep, and the best FID score obtained in 35k training steps (17.5k parameter updates).

G-LR	D-LR	G-clip	D-clip	FID
0.002	0.002	0.1	0.1	232.5
0.002	0.002	1	1	150.5
0.002	0.002	10	10	192.7
0.02	0.02	0.1	0.1	161.0
0.02	0.02	1	1	160.3
0.02	0.02	10	10	39.3
0.2	0.2	0.1	0.1	160.0
<b>0.2</b>	<b>0.2</b>	<b>1</b>	<b>1</b>	<b>36.3</b>
2	2	0.1	0.1	37.7

Table 6: clipped-SGDA (coordinate) hyperparameter sweep, and the best FID score obtained in 35k training steps.

G-LR	D-LR	G-clip	D-clip	FID
0.0002	0.0002	0.001	0.001	292.2
0.0002	0.0002	0.01	0.01	108.6
0.0002	0.0002	0.1	0.1	91.5
0.002	0.002	0.001	0.001	76.5
0.002	0.002	0.01	0.01	43.5
0.002	0.002	0.1	0.1	45.1
0.02	0.02	0.001	0.001	37.3
<b>0.02</b>	<b>0.02</b>	<b>0.01</b>	<b>0.01</b>	<b>26.7</b>
0.02	0.02	0.1	0.1	34.7

Table 7: clipped-SEG (coordinate) hyperparameter sweep, and the best FID score obtained in 35k training steps (17.5k parameter updates).

G-LR	D-LR	G-clip	D-clip	FID
0.0002	0.0002	0.001	0.001	298.7
0.0002	0.0002	0.01	0.01	146.5
0.0002	0.0002	0.1	0.1	158.4
0.002	0.002	0.001	0.001	112.8
0.002	0.002	0.01	0.01	52.7
0.002	0.002	0.1	0.1	66.5
0.02	0.02	0.001	0.001	43.5
<b>0.02</b>	<b>0.02</b>	<b>0.01</b>	<b>0.01</b>	<b>36.2</b>
0.02	0.02	0.1	0.1	75.3

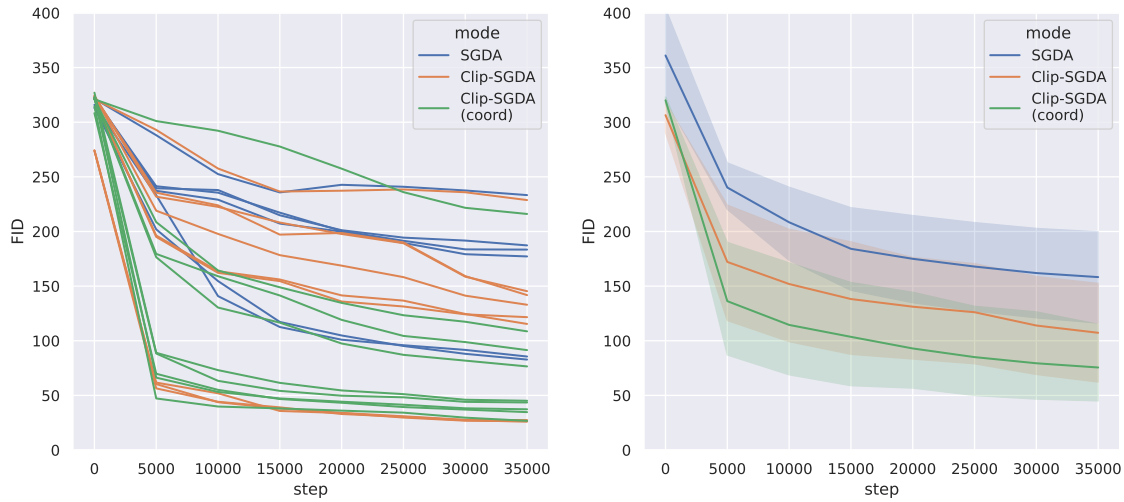


Figure 5: FID curves when training WGAN-GP for 35k steps with SGDA, clipped-SGDA (norm), and clipped-SGDA (coordinate), corresponding to the hyperparameters in Tables 2, 4 & 6 respectively. The left figure is the individual runs for each choice of hyperparameters, and the right is the mean and 95% confidence interval of these runs. Note that 4 of 10 runs diverged (NaN loss) for SGDA, which is not reflected in the mean FID for the right figure beyond the first step.

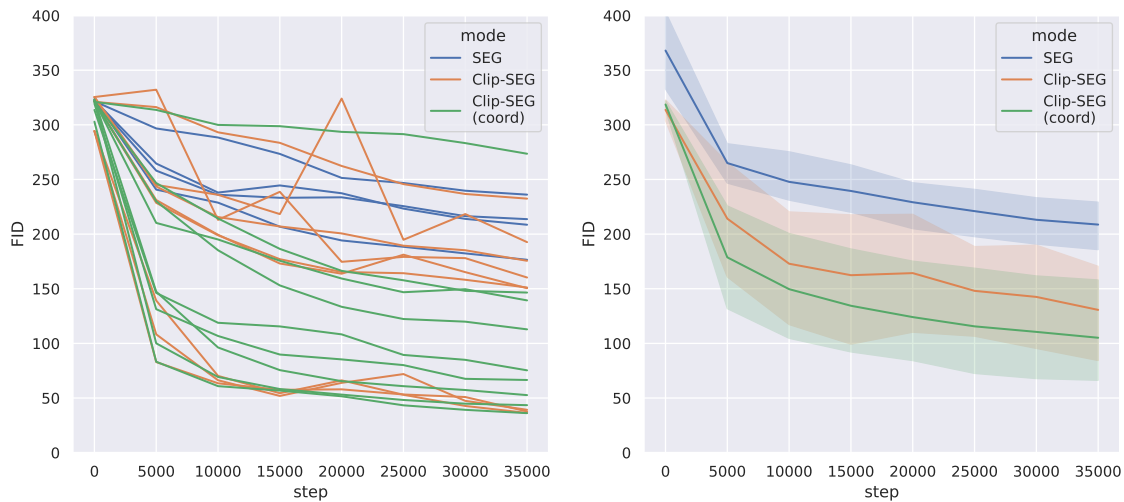


Figure 6: FID curves when training WGAN-GP for 35k steps with SEG, clipped-SEG (norm), and clipped-SEG (coordinate), corresponding to the hyperparameters in Tables 3, 5 & 7 respectively. The left figure is the individual runs for each choice of hyperparameters, and the right is the mean and 95% confidence interval of these runs. Note that 8 of 12 runs diverged (NaN loss) for SEG, which is not reflected in the mean FID for the right figure beyond the first step.



Figure 7: Samples generated from the best WGAN-GP model trained with SGDA.



Figure 8: Samples generated from the best WGAN-GP model trained with clipped-SGDA.



Figure 9: Samples generated from the best WGAN-GP model trained with clipped-SEG.





Figure 10: Samples generated from the best WGAN-GP model trained with clipped-SGDA (coordinate clipping).

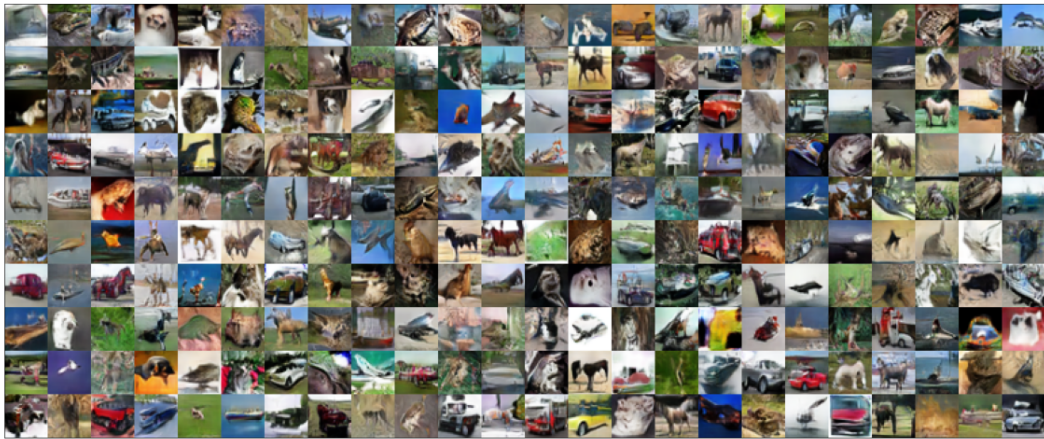


Figure 11: Samples generated from the best WGAN-GP model trained with clipped-SEG (coordinate clipping).

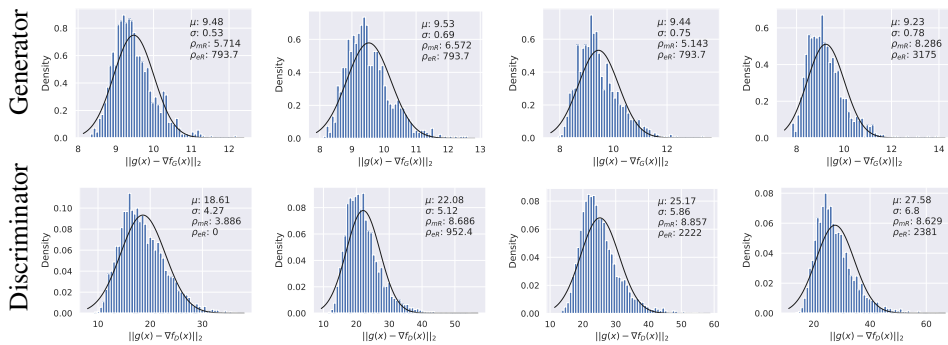


Figure 12: Evolution of gradient noise histograms for the best WGAN-GP model trained with SGDA.

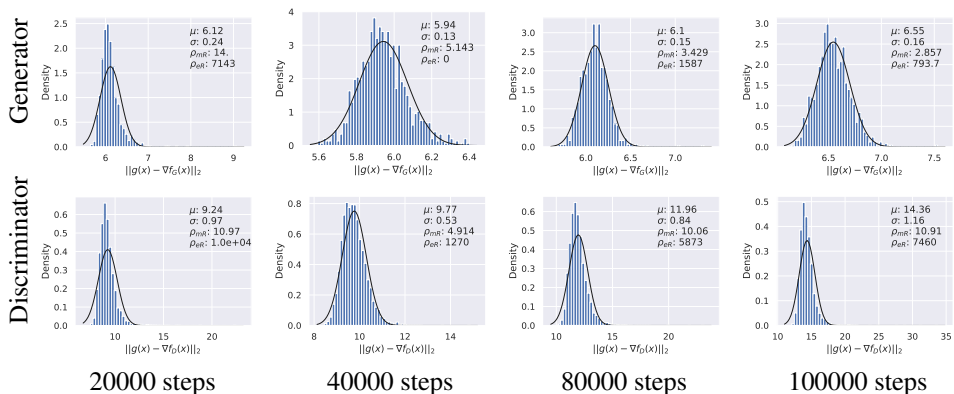


Figure 13: Evolution of gradient noise histograms for the best WGAN-GP model trained with clipped-SGDA.

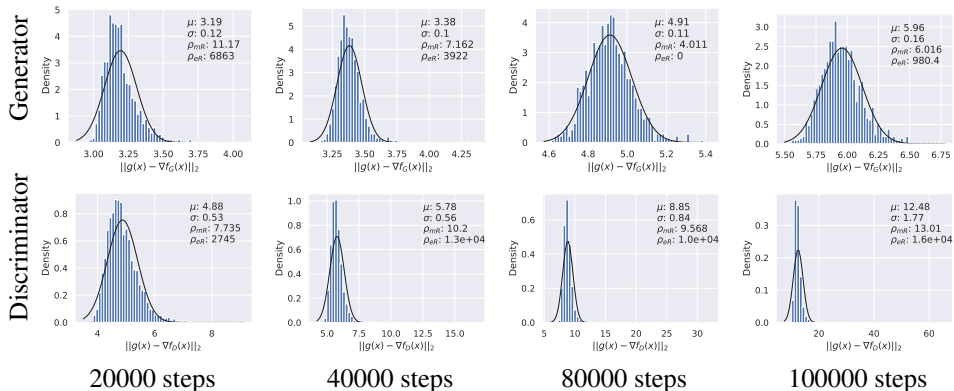


Figure 14: Evolution of gradient noise histograms for the best WGAN-GP model trained with clipped-SEG.

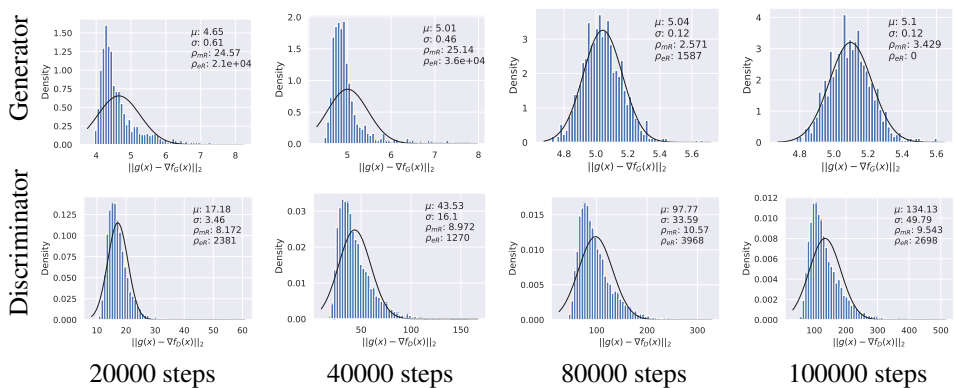


Figure 15: Evolution of gradient noise histograms for the best WGAN-GP model trained with clipped-SGDA (coordinate clipping).

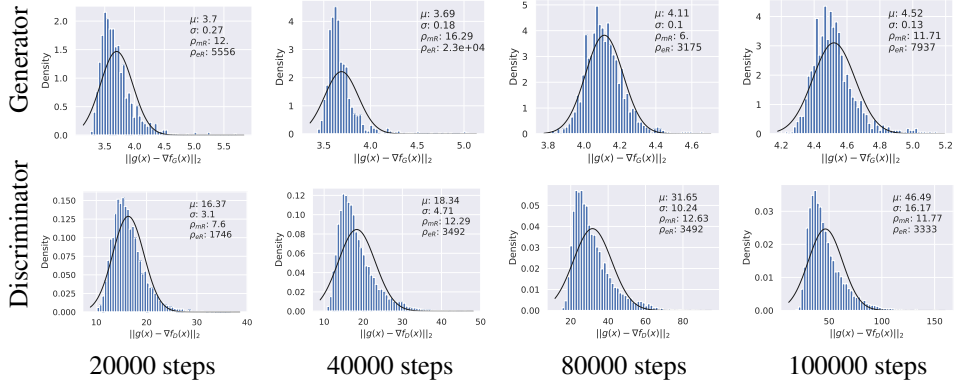


Figure 16: Evolution of gradient noise histograms for the best WGAN-GP model trained with clipped-SEG (coordinate clipping).

Table 8: StyleGAN2 SGDA hyperparameter sweep, and the best FID score obtained in 2600 kimgs.

G-LR	D-LR	FID
0.003	0.003	319.7
0.0075	0.0075	318.5
0.01	0.01	317.7
0.035	0.035	301.9
0.05	0.05	300.3
0.075	0.075	299.6
0.1	0.1	308.5
0.35	0.35	342.6
0.5	0.5	346.6

Table 9: StyleGAN2 clipped-SGDA (coordinate) hyperparameter sweep, and the best FID score obtained in 2600 kimgs. Bold row denotes the best run which was trained to convergence.

G-LR	D-LR	G-clip	D-clip	FID
0.2	0.2	0.001	0.001	243.5
0.3	0.3	0.001	0.001	169.5
0.35	0.35	0.0005	0.0005	192.9
0.35	0.35	0.001	0.001	148.6
<b>0.35</b>	<b>0.35</b>	<b>0.0025</b>	<b>0.0025</b>	<b>104.9</b>
0.35	0.35	0.005	0.005	149.1
0.35	0.5	0.01	0.01	170.6
0.4	0.4	0.001	0.001	155.8
0.5	0.5	0.0001	0.0001	289.8
0.5	0.5	0.001	0.001	136.1



Figure 17: FID curves when training StyleGAN2 for 2600 kings (thousands of images seen by the discriminator) with SGDA and clipped-SGDA (coordinate), corresponding to the hyperparameters in Tables 8 & 9 respectively. The left figure is the individual runs for each choice of hyperparameters, and the right is the mean and 95% confidence interval of these runs. Every SGDA-trained model for the wide range of learning rates we tried failed to improve the FID, while models trained with clipped-SGDA (with appropriately set hyperparameters) are generally able to learn some meaningful features and improve the FID.





(a) Additional samples generated from our best model trained with clipped-SGDA ( $\text{lr}=0.35$ ,  $\text{clip}=0.0025$ ).



(b) Additional samples generated from several different SGDA trained models, all of which failed to generate meaningful features. Each row corresponds to a model trained with different learning rates.

Figure 18: More StyleGAN2 samples.