

428 Appendix

429 430 Table of Contents

431	A Useful Lemmas	13
432	A.1 Lipschitz Continuity of $\omega(\theta)$	13
433	A.2 Lipschitz Continuity of $\nabla\omega(\theta)$	14
434	A.3 Smoothness of $J(\theta)$	16
435	B Non-asymptotic Analysis under the i.i.d. Setting	18
436	B.1 Tracking Error Analysis under the i.i.d. Setting	18
437	B.2 Proof under the i.i.d. Setting	22
438	C Non-asymptotic Analysis under the Markovian Setting	25
439	C.1 Tracking Error Analysis under the Markovian Setting	25
440	C.2 Proof under the Markovian Setting	32
441	D Experiments	35
442	D.1 Garnet Problem	35
443	D.2 Spiral Counter Example	35

447 We first introduce some notations. In the following proofs, $\|a\|$ denotes the ℓ_2 norm if a is a vector;
 448 and $\|A\|$ denotes the operator norm if A is a matrix.

449 In Appendix A, we prove the Lipschitz continuity of some important functions, including $\omega(\theta)$, $\nabla\omega(\theta)$
 450 and the gradient $\nabla J(\theta)$ of objective function. In Appendix B, we present the non-asymptotic analysis
 451 for the i.i.d. setting. In Appendix C, we present the non-asymptotic analysis for the Markovian
 452 setting. In appendix D, we present some numerical experiments.

453 A Useful Lemmas

454 A.1 Lipschitz Continuity of $\omega(\theta)$

455 In this section, we show that $\omega(\theta)$ is Lipschitz in θ .

456 **Lemma 1.** *For any $\theta, \theta' \in \mathbb{R}^N$, we have that*

$$\|\omega(\theta) - \omega(\theta')\| \leq L_\omega \|\theta - \theta'\|, \quad (21)$$

457 where $L_\omega = \frac{1}{\lambda_v} \left((1 + \gamma)C_\phi^2 + (r_{\max} + (1 + \gamma)C_v)D_v \right) + \frac{2C_\phi^2 D_v}{\lambda_v^2} (r_{\max} + (1 + \gamma)C_v).$

458 *Proof.* Recall that

$$\begin{aligned} \omega(\theta) &= \mathbb{E}_{\mu^{\pi_b}} [\phi_\theta(S)\phi_\theta(S)^\top]^{-1} \mathbb{E}_{\mu^{\pi_b}} [\rho(S, A)\delta_{S,A,S'}(\theta)\phi_\theta(S)] \\ &= A_\theta^{-1} \mathbb{E}_{\mu^{\pi_b}} [\rho(S, A)\delta_{S,A,S'}(\theta)\phi_\theta(S)], \end{aligned} \quad (22)$$

459 hence we can show the conclusion by showing that A_θ^{-1} and $\mathbb{E}_{\mu^{\pi_b}} [\rho(S, A)\delta_{S,A,S'}(\theta)\phi_\theta(S)]$ are both
 460 Lipschitz and bounded.

461 From Assumption 2, we know that

$$\|A_\theta^{-1}\| \leq \frac{1}{\lambda_v}. \quad (23)$$

462 We also show that

$$\|A_\theta^{-1} - A_{\theta'}^{-1}\|$$

$$\begin{aligned}
&= \|A_\theta^{-1} A_{\theta'} A_{\theta'}^{-1} - A_\theta^{-1} A_\theta A_{\theta'}^{-1}\| \\
&= \|A_\theta^{-1} (A_{\theta'} - A_\theta) A_{\theta'}^{-1}\| \\
&\leq \frac{2C_\phi D_v}{\lambda_v^2} \|\theta - \theta'\|,
\end{aligned} \tag{24}$$

463 which is from the fact that $\|A_\theta - A_{\theta'}\| = \|\mathbb{E}_{\mu^{\pi_b}}[\phi_\theta(S)\phi_\theta(S)^\top] - \mathbb{E}_{\mu^{\pi_b}}[\phi_{\theta'}(S)\phi_{\theta'}(S)^\top]\| \leq$
464 $2C_\phi D_v \|\theta - \theta'\|$.

465 By Assumption 1 and the boundedness of the reward function, it can be shown that for any $\theta \in \mathbb{R}^N$
466 and any $(s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$,

$$|\delta_{s,a,s'}(\theta)| = |r(s, a, s') + \gamma V_\theta(s') - V_\theta(s)| \leq r_{\max} + (1 + \gamma)C_v. \tag{25}$$

467 We then show that $\delta_{s,a,s'}(\theta)$ is Lipschitz, i.e., for any $\theta, \theta' \in \mathbb{R}^N$ and any $(s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$,

$$\begin{aligned}
&|\delta_{s,a,s'}(\theta) - \delta_{s,a,s'}(\theta')| \\
&= |\gamma V_\theta(s') - V_\theta(s) - \gamma V_{\theta'}(s') + V_{\theta'}(s)| \\
&\leq (\gamma + 1)C_\phi \|\theta - \theta'\|.
\end{aligned} \tag{26}$$

468 Hence, the function $\|\mathbb{E}_{\mu^{\pi_b}}[\rho(S, A)\delta_{S,A,S'}(\theta)\phi_\theta(S)]\|$ is Lipschitz:

$$\begin{aligned}
&\|\mathbb{E}_{\mu^{\pi_b}}[\rho(S, A)\delta_{S,A,S'}(\theta)\phi_\theta(S)] - \mathbb{E}_{\mu^{\pi_b}}[\rho(S, A)\delta_{S,A,S'}(\theta')\phi_{\theta'}(S)]\| \\
&= \|\mathbb{E}_{\mu^{\pi_b}}[\rho(S, A)\delta_{S,A,S'}(\theta)\phi_\theta(S)] - \mathbb{E}_{\mu^{\pi_b}}[\rho(S, A)\delta_{S,A,S'}(\theta')\phi_\theta(S)] \\
&\quad + \mathbb{E}_{\mu^{\pi_b}}[\rho(S, A)\delta_{S,A,S'}(\theta')\phi_\theta(S)] - \mathbb{E}_{\mu^{\pi_b}}[\rho(S, A)\delta_{S,A,S'}(\theta')\phi_{\theta'}(S)]\| \\
&\leq \|\mathbb{E}_{\mu^{\pi_b}}[\rho(S, A)\delta_{S,A,S'}(\theta)\phi_\theta(S)] - \mathbb{E}_{\mu^{\pi_b}}[\rho(S, A)\delta_{S,A,S'}(\theta')\phi_\theta(S)]\| \\
&\quad + \|\mathbb{E}_{\mu^{\pi_b}}[\rho(S, A)\delta_{S,A,S'}(\theta')\phi_\theta(S)] - \mathbb{E}_{\mu^{\pi_b}}[\rho(S, A)\delta_{S,A,S'}(\theta')\phi_{\theta'}(S)]\| \\
&\leq \mathbb{E}_{\mu^{\pi_b}}[\rho(S, A)|\delta_{S,A,S'}(\theta) - \delta_{S,A,S'}(\theta')| \|\phi_\theta(S)\|] \\
&\quad + \mathbb{E}_{\mu^{\pi_b}}[\rho(S, A)|\delta_{S,A,S'}(\theta')| \|\phi_\theta(S) - \phi_{\theta'}(S)\|] \\
&\stackrel{(a)}{\leq} (1 + \gamma)C_\phi^2 \|\theta - \theta'\| + (r_{\max} + (1 + \gamma)C_v)D_v \|\theta - \theta'\| \\
&= ((1 + \gamma)C_\phi^2 + (r_{\max} + (1 + \gamma)C_v)D_v) \|\theta - \theta'\|,
\end{aligned} \tag{27}$$

469 where (a) is from (26) and the fact that $\mathbb{E}_{\mu^{\pi_b}}[\rho(S, A)] = 1$. Also $\|\mathbb{E}_{\mu^{\pi_b}}[\rho(S, A)\delta_{S,A,S'}(\theta)\phi_\theta(S)]\|$
470 can be upper bounded as follows:

$$\|\mathbb{E}_{\mu^{\pi_b}}[\rho(S, A)\delta_{S,A,S'}(\theta)\phi_\theta(S)]\| \leq C_\phi(r_{\max} + (1 + \gamma)C_v). \tag{28}$$

471 Combining (23), (24), (28) and (27), we show that $\omega(\cdot)$ is Lipschitz in θ :

$$\begin{aligned}
&\|\omega(\theta) - \omega(\theta')\| \\
&\leq \left(\frac{1}{\lambda_v} ((1 + \gamma)C_\phi^2 + (r_{\max} + (1 + \gamma)C_v)D_v) + \frac{2C_\phi^2 D_v}{\lambda_v^2} (r_{\max} + (1 + \gamma)C_v) \right) \|\theta - \theta'\| \\
&\triangleq L_\omega \|\theta - \theta'\|,
\end{aligned} \tag{29}$$

472 where $L_\omega = \frac{1}{\lambda_v} ((1 + \gamma)C_\phi^2 + (r_{\max} + (1 + \gamma)C_v)D_v) + \frac{2C_\phi^2 D_v}{\lambda_v^2} (r_{\max} + (1 + \gamma)C_v)$. \square

473 A.2 Lipschitz Continuity of $\nabla\omega(\theta)$

474 In this section, we show that $\nabla\omega(\theta)$ is Lipschitz.

475 **Lemma 2.** For any $\theta, \theta' \in \mathbb{R}^N$, it follows that

$$\|\nabla\omega(\theta) - \nabla\omega(\theta')\| \leq D_\omega \|\theta - \theta'\|, \tag{30}$$

476 where

$$D_\omega = \left(\frac{(C_\phi L_v + 2D_v^2 + D_v C_\phi)}{\lambda_v^2} + \frac{8C_\phi^2 D_v^2}{\lambda_v^3} \right) C_\phi (r_{\max} + C_v + \gamma C_v)$$

$$\begin{aligned}
& + \frac{4C_\phi D_v}{\lambda_v^2} (C_\phi^2(1+\gamma) + D_v(r_{\max} + (1+\gamma)C_v)) \\
& + \frac{3C_\phi D_v(1+\gamma) + L_v(r_{\max} + (1+\gamma)C_v)}{\lambda_v}.
\end{aligned} \tag{31}$$

477 *Proof.* Recall the definition of $\omega(\theta) = A_\theta^{-1}\mathbb{E}_{\mu^{\pi_b}}[\rho(S, A)\delta_{S,A,S'}(\theta)\phi_\theta(S)]$, hence we have

$$\begin{aligned}
\nabla\omega(\theta) &= -A_\theta^{-1}(\nabla A_\theta)A_\theta^{-1}\mathbb{E}_{\mu^{\pi_b}}[\rho(S, A)\delta_{S,A,S'}(\theta)\phi_\theta(S)] \\
&\quad + A_\theta^{-1}\mathbb{E}_{\mu^{\pi_b}}[\nabla\rho(S, A)\delta_{S,A,S'}(\theta)\phi_\theta(S)],
\end{aligned} \tag{32}$$

478 where the tensor ∇A_θ can be equivalently viewed as an operator: $\mathbb{R}^N \rightarrow \mathbb{R}^{N \times N}$, i.e., $\nabla A_\theta(w) =$
479 $\nabla(A_\theta w)$ for any $w \in \mathbb{R}^N$.

480 We show that the operator norm of ∇A_θ is bounded as follows:

$$\begin{aligned}
\|\nabla A_\theta\| &= \sup_{\|w\|=1} \|\nabla A_\theta(w)\| \\
&= \sup_{\|w\|=1} \|\nabla(A_\theta w)\| \\
&= \sup_{\|w\|=1} \|\nabla\mathbb{E}_{\mu^{\pi_b}}[\phi_\theta(S)\phi_\theta(S)^\top w]\| \\
&= \sup_{\|w\|=1} \|\mathbb{E}_{\mu^{\pi_b}}[(\phi_\theta(S)^\top w)\nabla\phi_\theta(S)] + \mathbb{E}_{\mu^{\pi_b}}[\phi_\theta(S)(\nabla\phi_\theta(S)^\top w)^\top]\| \\
&\leq \sup_{\|w\|=1} 2C_\phi D_v \|w\| \\
&= 2C_\phi D_v.
\end{aligned} \tag{33}$$

481 The Lipschitz continuous of ∇A_θ can be shown as follows:

$$\begin{aligned}
&\|\nabla A_\theta - \nabla A_{\theta'}\| \\
&= \sup_{\|w\|=1} \|\nabla(A_\theta w) - \nabla(A_{\theta'} w)\| \\
&= \sup_{\|w\|=1} \|\mathbb{E}_{\mu^{\pi_b}}[\nabla\phi_\theta(S)(\phi_\theta(S)^\top w) + (\nabla\phi_\theta(S)^\top w)\phi_\theta(S)^\top] - \nabla\phi_{\theta'}(S)(\phi_{\theta'}(S)^\top w) \\
&\quad - (\nabla\phi_{\theta'}(S)^\top w)\phi_{\theta'}(S)^\top]\| \\
&\leq \sup_{\|w\|=1} (C_\phi L_v + 2D_v^2 + D_v C_\phi) \|\theta - \theta'\| \|w\| \\
&= (C_\phi L_v + 2D_v^2 + D_v C_\phi) \|\theta - \theta'\|.
\end{aligned} \tag{34}$$

482 Then we conclude that the operator norm of $-A_\theta^{-1}(\nabla A_\theta)$ is upper bounded by $\frac{2C_\phi D_v}{\lambda_v}$, and is
483 Lipschitz with constant $\frac{(C_\phi L_v + 2D_v^2 + D_v C_\phi)}{\lambda_v} + \frac{4C_\phi^2 D_v^2}{\lambda_v^2}$. It can be further seen that $-A_\theta^{-1}(\nabla A_\theta)A_\theta^{-1}$
484 is upper bounded by $\frac{2C_\phi D_v}{\lambda_v^2}$, and Lipschitz with constant $\frac{(C_\phi L_v + 2D_v^2 + D_v C_\phi)}{\lambda_v^2} + \frac{8C_\phi^2 D_v^2}{\lambda_v^3}$.

485 Recall that we have shown in (28) that

$$\begin{aligned}
&\|\mathbb{E}_{\mu^{\pi_b}}[\rho(S, A)\delta_{S,A,S'}(\theta)\phi_\theta(S)] - \mathbb{E}_{\mu^{\pi_b}}[\rho(S, A)\delta_{S,A,S'}(\theta')\phi_{\theta'}(S)]\| \\
&\leq ((1+\gamma)C_\phi^2 + (r_{\max} + (1+\gamma)C_v)D_v) \|\theta - \theta'\|,
\end{aligned} \tag{35}$$

486 and it is upper bounded by $C_\phi(r_{\max} + (1+\gamma)C_v)$. Hence we have that
487 $-A_\theta^{-1}(\nabla A_\theta)A_\theta^{-1}\mathbb{E}_{\mu^{\pi_b}}[\rho(S, A)\delta_{S,A,S'}(\theta)\phi_\theta(S)]$ can be upper bounded by $(r_{\max} + (1+\gamma)C_v)\frac{2C_\phi^2 D_v}{\lambda_v^2}$, and it is Lipschitz with constant $\left(\frac{(C_\phi L_v + 2D_v^2 + D_v C_\phi)}{\lambda_v^2} + \frac{8C_\phi^2 D_v^2}{\lambda_v^3}\right)C_\phi(r_{\max} + C_v + \gamma C_v) + \frac{2C_\phi D_v}{\lambda_v^2}((1+\gamma)C_\phi^2 + (r_{\max} + (1+\gamma)C_v)D_v) \triangleq L_A$.

489 For the second term of (32), we also show it is Lipschitz as follows. First
490 note that $\nabla\delta_{s,a,s'}(\theta)\phi_\theta(s) = \nabla\delta_{s,a,s'}(\theta)\phi_\theta(s)^\top + \delta_{s,a,s'}(\theta)\nabla\phi_\theta(s)$, hence we know

492 $\mathbb{E}_{\mu^{\pi_b}} [\nabla \rho(S, A) \delta_{S, A, S'}(\theta) \phi_\theta(S)]$ can be upper bounded by $C_\phi^2(1 + \gamma) + D_v(r_{\max} + (1 + \gamma)C_v)$,
493 and is Lipschitz with constant $3C_\phi D_v(1 + \gamma) + L_v(r_{\max} + (1 + \gamma)C_v)$. Finally we con-
494 clude that the second term in (32) $A_\theta^{-1} \mathbb{E}_{\mu^{\pi_b}} [\nabla \rho(S, A) \delta_{S, A, S'}(\theta) \phi_\theta(S)]$ is Lipschitz with constant
495 $\frac{2C_\phi D_v}{\lambda_v^2} \left(C_\phi^2(1 + \gamma) + D_v(r_{\max} + (1 + \gamma)C_v) \right) + \frac{3C_\phi D_v(1 + \gamma) + L_v(r_{\max} + (1 + \gamma)C_v)}{\lambda_v} \triangleq L'_A$.

496 Hence $\nabla \omega(\theta)$ is Lipschitz with constant $L_A + L'_A \triangleq D_\omega$, where

$$\begin{aligned} D_\omega = & \left(\frac{(C_\phi L_v + 2D_v^2 + D_v C_\phi)}{\lambda_v^2} + \frac{8C_\phi^2 D_v^2}{\lambda_v^3} \right) C_\phi(r_{\max} + C_v + \gamma C_v) \\ & + \frac{4C_\phi D_v}{\lambda_v^2} (C_\phi^2(1 + \gamma) + D_v(r_{\max} + (1 + \gamma)C_v)) \\ & + \frac{3C_\phi D_v(1 + \gamma) + L_v(r_{\max} + (1 + \gamma)C_v)}{\lambda_v}. \end{aligned} \quad (36)$$

497 \square

498 A.3 Smoothness of $J(\theta)$

499 In the following lemma, we show that the objective function $J(\theta)$ is L_J -smooth. We note that the
500 smoothness of $J(\theta)$ is assumed in [37] instead of being proved as in this paper.

501 **Lemma 3.** $J(\theta)$ is L_J -smooth, i.e., for any $\theta, \theta' \in \mathbb{R}^N$,

$$\|\nabla J(\theta) - \nabla J(\theta')\| \leq L_J \|\theta - \theta'\|, \quad (37)$$

502 where

$$\begin{aligned} L_J = & 2 \left((1 + \gamma)C_\phi^2 + (r_{\max} + (1 + \gamma)C_v)D_v \right) + 2\gamma \left(C_\phi^2 L_\omega + 2D_v \frac{C_\phi^2}{\lambda_v} (r_{\max} + (1 + \gamma)C_v) \right) \\ & + 2 \left((D_v R_\omega + C_\phi L_\omega + (1 + \gamma)C_\phi) D_v R_\omega \right. \\ & \left. + (R_\omega L_V + D_v L_\omega) ((r_{\max} + (1 + \gamma)C_v) + C_\phi R_\omega) \right). \end{aligned} \quad (38)$$

503 *Proof.* Before we prove the main statement, we first drive some boundedness and Lipschitz properties.
504 Recall that

$$\begin{aligned} -\frac{\nabla J(\theta)}{2} = & \mathbb{E}_{\mu^{\pi_b}} \left[(\rho(S, A) \delta_{S, A, S'}(\theta) \phi_\theta(S) - \gamma \rho(S, A) \phi_\theta(S') \phi_\theta(S)^\top \omega(\theta) \right. \\ & \left. - h_{S, A, S'}(\theta, \omega(\theta))) \right], \end{aligned} \quad (39)$$

$$\omega(\theta) = \mathbb{E}_{\mu^{\pi_b}} [\phi_\theta(S) \phi_\theta(S)^\top]^{-1} \mathbb{E}_{\mu^{\pi_b}} [\rho(S, A) \delta_{S, A, S'}(\theta) \phi_\theta(S)], \quad (40)$$

$$h_{s, a, s'}(\theta, \omega(\theta)) = (\rho(s, a) \delta_{s, a, s'}(\theta) - \phi_\theta(s)^\top \omega(\theta)) \nabla^2 V_\theta(s) \omega(\theta). \quad (41)$$

505 We have shown in Lemma 1 that for any $\theta \in \mathbb{R}^N$ and any $(s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$,

$$|\delta_{s, a, s'}(\theta)| = |r(s, a, s') + \gamma V_\theta(s') - V_\theta(s)| \leq r_{\max} + (1 + \gamma)C_v; \quad (42)$$

506 and that

$$\begin{aligned} & \|\mathbb{E}_{\mu^{\pi_b}} [\rho(S, A) \delta_{S, A, S'}(\theta) \phi_\theta(S)] - \mathbb{E}_{\mu^{\pi_b}} [\rho(S, A) \delta_{S, A, S'}(\theta') \phi_{\theta'}(S)]\| \\ & \leq ((1 + \gamma)C_\phi^2 + (r_{\max} + (1 + \gamma)C_v)D_v) \|\theta - \theta'\|. \end{aligned} \quad (43)$$

507 Also it is easy to see from the definition that

$$\|\omega(\theta)\| \leq \frac{C_\phi}{\lambda_v} (r_{\max} + (1 + \gamma)C_v) \triangleq R_\omega. \quad (44)$$

508 Hence the Lipschitz continuity of $\mathbb{E}_{\mu^{\pi_b}}[\rho(S, A)\phi_\theta(S')\phi_\theta(S)^\top]\omega(\theta)$ can be shown as follows

$$\begin{aligned}
& \|\mathbb{E}_{\mu^{\pi_b}}[\rho(S, A)\phi_\theta(S')\phi_\theta(S)^\top]\omega(\theta) - \mathbb{E}_{\mu^{\pi_b}}[\rho(S, A)\phi_{\theta'}(S')\phi_{\theta'}(S)^\top]\omega(\theta')\| \\
& \leq \|\mathbb{E}_{\mu^{\pi_b}}[\rho(S, A)\phi_\theta(S')\phi_\theta(S)^\top]\omega(\theta) - \mathbb{E}_{\mu^{\pi_b}}[\rho(S, A)\phi_\theta(S')\phi_\theta(S)^\top]\omega(\theta')\| \\
& \quad + \|\mathbb{E}_{\mu^{\pi_b}}[\rho(S, A)\phi_\theta(S')\phi_\theta(S)^\top]\omega(\theta') - \mathbb{E}_{\mu^{\pi_b}}[\rho(S, A)\phi_{\theta'}(S')\phi_{\theta'}(S)^\top]\omega(\theta')\| \\
& \stackrel{(a)}{\leq} C_\phi^2 L_\omega \|\theta - \theta'\| + 2C_\phi D_v R_\omega \|\theta - \theta'\| \\
& = \left(C_\phi^2 L_\omega + 2D_v \frac{C_\phi^2}{\lambda_v} (r_{\max} + (1 + \gamma)C_v) \right) \|\theta - \theta'\|,
\end{aligned} \tag{45}$$

509 where (a) is due to the fact that $\omega(\theta)$ is Lipschitz in (21) and the fact that

$$\|\mathbb{E}_{\mu^{\pi_b}}[\rho(S, A)\phi_\theta(S')\phi_\theta(S)^\top] - \mathbb{E}_{\mu^{\pi_b}}[\rho(S, A)\phi_{\theta'}(S')\phi_{\theta'}(S)^\top]\| \leq 2C_\phi D_v \|\theta - \theta'\|. \tag{46}$$

510 We then show that the function $h_{s,a,s'}(\theta, \omega(\theta))$ is Lipschitz in θ as follows. We first note that for any
511 $s \in \mathcal{S}$ and $\theta, \theta' \in \mathbb{R}^N$,

$$\begin{aligned}
& \|\phi_\theta(s)^\top \omega(\theta) - \phi_{\theta'}(s)^\top \omega(\theta')\| \\
& \leq \|\phi_\theta(s)^\top \omega(\theta) - \phi_{\theta'}(s)^\top \omega(\theta)\| + \|\phi_{\theta'}(s)^\top \omega(\theta) - \phi_{\theta'}(s)^\top \omega(\theta')\| \\
& \leq (D_v R_\omega + C_\phi L_\omega) \|\theta - \theta'\|.
\end{aligned} \tag{47}$$

512 This implies that for any $(s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$ and $\theta, \theta' \in \mathbb{R}^N$,

$$\begin{aligned}
& \|\rho(s, a)\delta_{s,a,s'}(\theta) - \phi_\theta(s)^\top \omega(\theta) - \rho(s, a)\delta_{s,a,s'}(\theta') + \phi_{\theta'}(s)^\top \omega(\theta')\| \\
& \leq (D_v R_\omega + C_\phi L_\omega + (1 + \gamma)C_\phi \rho(s, a)) \|\theta - \theta'\|.
\end{aligned} \tag{48}$$

513 We also show the following function is Lipschitz:

$$\begin{aligned}
& \|\nabla^2 V_\theta(s)\omega(\theta) - \nabla^2 V_{\theta'}(s)\omega(\theta')\| \\
& \leq \|\nabla^2 V_\theta(s)\omega(\theta) - \nabla^2 V_{\theta'}(s)\omega(\theta)\| + \|\nabla^2 V_{\theta'}(s)\omega(\theta) - \nabla^2 V_{\theta'}(s)\omega(\theta')\| \\
& \leq R_\omega L_V \|\theta - \theta'\| + D_v L_\omega \|\theta - \theta'\| \\
& = (R_\omega L_V + D_v L_\omega) \|\theta - \theta'\|.
\end{aligned} \tag{49}$$

514 Combining (48) and (49), it can be shown that $h_{s,a,s'}(\theta, \omega(\theta))$ is Lipschitz in θ as follows

$$\begin{aligned}
& \|h_{s,a,s'}(\theta, \omega(\theta)) - h_{s,a,s'}(\theta', \omega(\theta'))\| \\
& = \|(\rho(s, a)\delta_{s,a,s'}(\theta) - \phi_\theta(s)^\top \omega(\theta)) \nabla^2 V_\theta(s)\omega(\theta) \\
& \quad - (\rho(s, a)\delta_{s,a,s'}(\theta') - \phi_{\theta'}(s)^\top \omega(\theta')) \nabla^2 V_{\theta'}(s)\omega(\theta')\| \\
& \leq ((D_v R_\omega + C_\phi L_\omega + (1 + \gamma)C_\phi \rho(s, a)) D_v R_\omega) \|\theta - \theta'\| \\
& \quad + (R_\omega L_V + D_v L_\omega) (\rho(s, a)(r_{\max} + (1 + \gamma)C_v) + C_\phi R_\omega) \|\theta - \theta'\|.
\end{aligned} \tag{50}$$

515 From the results in (43), (45) and (50), it follows that

$$\begin{aligned}
& \|\nabla J(\theta) - \nabla J(\theta')\| \\
& \leq 2 \|\mathbb{E}_{\mu^{\pi_b}}[\rho(S, A)\delta_{S,A,S'}(\theta)\phi_\theta(S) - \rho(S, A)\delta_{S,A,S'}(\theta')\phi_{\theta'}(S)]\| \\
& \quad + 2\gamma \|\mathbb{E}_{\mu^{\pi_b}}[\rho(S, A)\phi_\theta(S')\phi_\theta(S)^\top \omega(\theta) - \rho(S, A)\phi_{\theta'}(S')\phi_{\theta'}(S)^\top \omega(\theta')]\| \\
& \quad + 2 \|\mathbb{E}_{\mu^{\pi_b}}[h_{S,A,S'}(\theta, \omega(\theta)) - h_{S,A,S'}(\theta', \omega(\theta'))]\| \\
& \leq 2 \left((1 + \gamma)C_\phi^2 + (r_{\max} + (1 + \gamma)C_v) D_v \right) \|\theta - \theta'\| \\
& \quad + 2\gamma \left(C_\phi^2 L_\omega + 2D_v \frac{C_\phi^2}{\lambda_v} (r_{\max} + (1 + \gamma)C_v) \right) \|\theta - \theta'\| \\
& \quad + 2\mathbb{E}_{\mu^{\pi_b}}[((D_v R_\omega + C_\phi L_\omega + (1 + \gamma)C_\phi \rho(S, A)) D_v R_\omega)] \|\theta - \theta'\| \\
& \quad + 2\mathbb{E}_{\mu^{\pi_b}}[(R_\omega L_V + D_v L_\omega) (\rho(S, A)(r_{\max} + (1 + \gamma)C_v) + C_\phi R_\omega)] \|\theta - \theta'\|
\end{aligned}$$

$$\begin{aligned}
&\stackrel{(a)}{\leq} 2 \left((1 + \gamma) C_\phi^2 + (r_{\max} + (1 + \gamma) C_v) D_v \right) \|\theta - \theta'\| \\
&\quad + 2\gamma \left(C_\phi^2 L_\omega + 2D_v \frac{C_\phi^2}{\lambda_v} (r_{\max} + (1 + \gamma) C_v) \right) \|\theta - \theta'\| \\
&\quad + 2 \left((D_v R_\omega + C_\phi L_\omega + (1 + \gamma) C_\phi) D_v R_\omega \right. \\
&\quad \left. + (R_\omega L_V + D_v L_\omega) ((r_{\max} + (1 + \gamma) C_v) + C_\phi R_\omega) \right) \|\theta - \theta'\| \\
&\triangleq L_J \|\theta - \theta'\|,
\end{aligned} \tag{51}$$

516 where (a) is due to the fact that $\mathbb{E}_{\mu^{\pi_b}} [\rho(S, A)] = 1$, and

$$\begin{aligned}
L_J = & 2 \left((1 + \gamma) C_\phi^2 + (r_{\max} + (1 + \gamma) C_v) D_v \right) + 2\gamma \left(C_\phi^2 L_\omega + 2D_v \frac{C_\phi^2}{\lambda_v} (r_{\max} + (1 + \gamma) C_v) \right) \\
& + 2 \left((D_v R_\omega + C_\phi L_\omega + (1 + \gamma) C_\phi) D_v R_\omega \right. \\
& \left. + (R_\omega L_V + D_v L_\omega) ((r_{\max} + (1 + \gamma) C_v) + C_\phi R_\omega) \right).
\end{aligned} \tag{52}$$

517 This completes the proof. \square

518 B Non-asymptotic Analysis under the i.i.d. Setting

519 First we introduce the off-policy TDC learning with non-linear function approximation algorithm
520 under the i.i.d. setting in Algorithm 2. We then bound the tracking error in Appendix B.1, and prove
the Theorem 1 under the i.i.d. setting in Appendix B.2.

Algorithm 2 Non-Linear Off-Policy TDC under the i.i.d. Setting

Input: $T, \alpha, \beta, \pi, \pi_b, \{V_\theta | \theta \in \mathbb{R}^N\}$

Initialization: θ_0, ω_0

- 1: Choose $W \sim \text{Uniform}(0, 1, \dots, T - 1)$
- 2: **for** $t = 0, 1, \dots, W - 1$ **do**
- 3: Sample $O_t = (s_t, a_t, r_t, s'_t)$ according to μ^{π_b}
- 4: $\rho_t = \frac{\pi(a_t | s_t)}{\pi_b(a_t | s_t)}$
- 5: $\delta_t(\theta_t) = r(s_t, a_t, s'_t) + \gamma V_{\theta_t}(s'_t) - V_{\theta_t}(s_t)$
- 6: $h_t(\theta_t, \omega_t) = (\rho_t \delta_t(\theta_t) - \phi_{\theta_t}(s_t)^\top \omega_t) \nabla^2 V_{\theta_t}(s_t) \omega_t$
- 7: $\omega_{t+1} = \Pi_{R_\omega} (\omega_t + \beta (-\phi_{\theta_t}(s_t) \phi_{\theta_t}(s_t)^\top \omega_t + \rho_t \delta_t(\theta_t) \phi_{\theta_t}(s_t)))$
- 8: $\theta_{t+1} = \theta_t + \alpha (\rho_t \delta_t(\theta_t) \phi_{\theta_t}(s_t) - \gamma \rho_t \phi_{\theta_t}(s'_t) \phi_{\theta_t}(s_t)^\top \omega_t - h_t(\theta_t, \omega_t))$
- 9: **end for**

Output: θ_W

521

522 We note that under the i.i.d. setting, it is assumed that at each time step t , a sample $O_t = (s_t, a_t, r_t, s'_t)$
523 is available, where $s_t \sim \mu^{\pi_b}(\cdot)$, $a_t \sim \pi_b(\cdot | s_t)$ and $s'_t \sim \mathbb{P}(\cdot | s_t, a_t)$.

524 B.1 Tracking Error Analysis under the i.i.d. Setting

525 Denote the tracking error by $z_t = \omega_t - \omega(\theta_t)$. Then by the update of ω_t , the update of z_t can be
526 written as

$$\begin{aligned}
z_{t+1} &= \omega_{t+1} - \omega(\theta_{t+1}) \\
&= \omega_t + \beta (-\phi_{\theta_t}(s_t) \phi_{\theta_t}(s_t)^\top \omega_t + \rho_t \delta_t(\theta_t) \phi_{\theta_t}(s_t)) - \omega(\theta_{t+1}) \\
&= z_t + \omega(\theta_t) - \omega(\theta_{t+1}) + \beta (-\phi_{\theta_t}(s_t) \phi_{\theta_t}(s_t)^\top (z_t + \omega(\theta_t)) + \rho_t \delta_t(\theta_t) \phi_{\theta_t}(s_t)) \\
&= z_t + \omega(\theta_t) - \omega(\theta_{t+1}) + \beta (-A_{\theta_t}(s_t) z_t - A_{\theta_t}(s_t) \omega(\theta_t) + \rho_t \delta_t(\theta_t) \phi_{\theta_t}(s_t)),
\end{aligned} \tag{53}$$

527 where $A_{\theta_t}(s_t) = \phi_{\theta_t}(s_t) \phi_{\theta_t}(s_t)^\top$. It then follows that

$$\|z_{t+1}\|^2$$

$$\begin{aligned}
&= \|z_t + \omega(\theta_t) - \omega(\theta_{t+1}) + \beta(-A_{\theta_t}(s_t)z_t - A_{\theta_t}(s_t)\omega(\theta_t) + \rho_t\delta_t(\theta_t)\phi_{\theta_t}(s_t))\|^2 \\
&= \|z_t\|^2 + \|\omega(\theta_t) - \omega(\theta_{t+1}) + \beta(-A_{\theta_t}(s_t)z_t - A_{\theta_t}(s_t)\omega(\theta_t) + \rho_t\delta_t(\theta_t)\phi_{\theta_t}(s_t))\|^2 \\
&\quad + 2\langle z_t, \omega(\theta_t) - \omega(\theta_{t+1}) \rangle - 2\beta\langle z_t, A_{\theta_t}(s_t)z_t \rangle + 2\beta\langle z_t, -A_{\theta_t}(s_t)\omega(\theta_t) + \rho_t\delta_t(\theta_t)\phi_{\theta_t}(s_t) \rangle \\
&\leq \|z_t\|^2 + \underbrace{2\beta^2\|(-A_{\theta_t}(s_t)z_t - A_{\theta_t}(s_t)\omega(\theta_t) + \rho_t\delta_t(\theta_t)\phi_{\theta_t}(s_t))\|^2}_{(a)} \\
&\quad + \underbrace{2\|\omega(\theta_t) - \omega(\theta_{t+1})\|^2}_{(b)} + \underbrace{2\langle z_t, \omega(\theta_t) - \omega(\theta_{t+1}) \rangle}_{(c)} - \underbrace{2\beta\langle z_t, A_{\theta_t}(s_t)z_t \rangle}_{(d)} \\
&\quad + 2\beta\langle z_t, -A_{\theta_t}(s_t)\omega(\theta_t) + \rho_t\delta_t(\theta_t)\phi_{\theta_t}(s_t) \rangle. \tag{54}
\end{aligned}$$

528 We then provide the bounds of the terms in (54) one by one. Their proofs can be found in Appendices B.1.1 to B.1.4.

530 **Term (a) can be bounded as follows:**

$$2\beta^2\|(-A_{\theta_t}(s_t)z_t - A_{\theta_t}(s_t)\omega(\theta_t) + \rho_t\delta_t(\theta_t)\phi_{\theta_t}(s_t))\|^2 \leq 4\beta^2C_\phi^2\|z_t\|^2 + 4\beta^2C_{g1}, \tag{55}$$

531 where $C_{g1} = \left(\frac{C_\phi^3}{\lambda_v}(r_{\max} + (1+\gamma)C_v) + \rho_{\max}C_\phi(r_{\max} + (1+\gamma)C_v)\right)^2$.

532 **Term (b) can be bounded as follows:**

$$2\|\omega(\theta_t) - \omega(\theta_{t+1})\|^2 \leq 4\alpha^2L_\omega^2L_g^2\|z_t\|^2 + 4\alpha^2C_g^2L_\omega^2, \tag{56}$$

533 where $C_g = \rho_{\max}C_\phi(r_{\max} + (1+\gamma)C_v) + \gamma\rho_{\max}R_\omega C_\phi^2 + D_vR_\omega(R_\omega C_\phi + \rho_{\max}(r_{\max} + C_v + \gamma C_v))$.

534 **Term (c) can be bounded as follows:**

$$\begin{aligned}
&2\langle z_t, \omega(\theta_t) - \omega(\theta_{t+1}) \rangle \\
&\leq 2(\alpha L_\omega L_g + \frac{1}{2}\alpha L_\omega + 4\alpha^2C_g L_g D_\omega)\|z_t\|^2 + \frac{\alpha L_\omega}{4}\|\nabla J(\theta_t)\|^2 + \frac{\alpha^2C_g^3D_\omega}{L_g} + 2\alpha\eta_G(\theta_t, z_t, O_t), \tag{57}
\end{aligned}$$

535 where $\eta_G(\theta_t, z_t, O_t) = -\left\langle z_t, \nabla\omega(\theta_t) \left(G_{t+1}(\theta_t, \omega(\theta_t)) + \frac{\nabla J(\theta_t)}{2}\right)\right\rangle$.

536 **Term (d) can be bounded as follows:**

$$-2\beta\langle z_t, A_{\theta_t}(s_t)z_t \rangle \leq -2\beta\lambda_v\|z_t\|^2 + 2\beta\langle z_t, (A_{\theta_t} - A_{\theta_t}(s_t))z_t \rangle, \tag{58}$$

537 where $A_\theta = \mathbb{E}_{\mu^{\pi_b}}[\phi_\theta(S)\phi_\theta(S)^\top]$ is the expectation of $A_\theta(S)$.

538 By plugging all the bounds from (55), (56), (57) and (58) in (54), it follows that

$$\begin{aligned}
&\|z_{t+1}\|^2 \\
&\leq (1 + 4\beta^2C_\phi^2 + 4\alpha^2L_\omega^2L_g^2 + 2\alpha L_w L_g + \alpha L_w + 8\alpha^2C_g L_g D_\omega - 2\beta\lambda_v)\|z_t\|^2 \\
&\quad + \frac{1}{4}\alpha L_\omega\|\nabla J(\theta_t)\|^2 + 4\beta^2C_{g1} + 4\alpha^2C_g^2L_\omega^2 + \frac{\alpha^2C_g^3D_\omega}{L_g} + 2\alpha\eta_G(\theta_t, z_t, O_t) \\
&\quad + 2\beta\langle z_t, (A_{\theta_t} - A_{\theta_t}(s_t))z_t \rangle + 2\beta\langle z_t, -A_{\theta_t}(s_t)\omega(\theta_t) + \rho_t\delta_t(\theta_t)\phi_{\theta_t}(s_t) \rangle \\
&\triangleq (1 - q)\|z_t\|^2 + \frac{\alpha L_\omega}{4}\|\nabla J(\theta_t)\|^2 + 4\beta^2C_{g1} + 4\alpha^2C_g^2L_\omega^2 + \frac{\alpha^2C_g^3D_\omega}{L_g} + 2\alpha\eta_G(\theta_t, z_t, O_t) \\
&\quad + 2\beta\langle z_t, (A_{\theta_t} - A_{\theta_t}(s_t))z_t \rangle + 2\beta\langle z_t, -A_{\theta_t}(s_t)\omega(\theta_t) + \rho_t\delta_t(\theta_t)\phi_{\theta_t}(s_t) \rangle, \tag{59}
\end{aligned}$$

539 where $q = 2\beta\lambda_v - 4\beta^2C_\phi^2 - 4\alpha^2L_\omega^2L_g^2 - 2\alpha L_w L_g - \alpha L_w - 8\alpha^2C_g L_g D_\omega$. Note that $q = \mathcal{O}(\beta - \beta^2 - \alpha - \alpha^2) = \mathcal{O}(\beta)$, hence we can choose α and β such that $q > 0$.

541 Note that under the i.i.d. setting,

$$\begin{aligned}
\mathbb{E}[\eta_G(\theta_t, z_t, O_t)] &= \mathbb{E}[\mathbb{E}[\eta_G(\theta_t, z_t, O_t)|\mathcal{F}_t]] \\
&= \mathbb{E}\left[-\left\langle z_t, \nabla\omega(\theta_t) \mathbb{E}\left[\left(G_{t+1}(\theta_t, \omega(\theta_t)) + \frac{\nabla J(\theta_t)}{2}\right) \middle| \mathcal{F}_t\right]\right\rangle\right]
\end{aligned}$$

$$= 0, \quad (60)$$

542 which is due to the fact that $\mathbb{E}_{\mu^{\pi_b}}[G_{t+1}(\theta, \omega(\theta))] = -\frac{\nabla J(\theta)}{2}$ when θ is fixed, and \mathcal{F}_t is the σ -field
543 generated by the randomness until θ_t and ω_t . Similarly, it can also be shown that

$$\mathbb{E}[\langle z_t, (A_{\theta_t} - A_{\theta_t}(s_t))z_t \rangle] = 0 \quad (61)$$

$$\mathbb{E}[\langle z_t, -A_{\theta_t}(s_t)\omega(\theta_t) + \rho_t\delta_t(\theta_t)\phi_{\theta_t}(s_t) \rangle] = 0. \quad (62)$$

544 Hence the tracking error in (59) can be further bounded as

$$\mathbb{E}[\|z_{t+1}\|^2] \leq (1-q)\mathbb{E}[\|z_t\|^2] + \frac{\alpha L_\omega}{4}\mathbb{E}[\|\nabla J(\theta_t)\|^2] + 4\beta^2 C_{g1} + 4\alpha^2 C_g^2 L_\omega^2 + \frac{\alpha^2 C_g^3 D_\omega}{L_g}. \quad (63)$$

545 Recursively applying the inequality in (63), it follows that

$$\begin{aligned} \mathbb{E}[\|z_t\|^2] &\leq (1-q)^t \|z_0\|^2 + \frac{\alpha L_\omega}{4} \sum_{i=0}^t (1-q)^{t-i} \mathbb{E}[\|\nabla J(\theta_i)\|^2] \\ &\quad + \frac{1}{q} \left(4\beta^2 C_{g1} + 4\alpha^2 C_g^2 L_\omega^2 + \frac{\alpha^2 C_g^3 D_\omega}{L_g} \right), \end{aligned} \quad (64)$$

546 and summing up w.r.t. t from 0 to $T-1$, it follows that

$$\begin{aligned} \frac{\sum_{t=0}^{T-1} \mathbb{E}[\|z_t\|^2]}{T} &\leq \frac{\sum_{t=0}^{T-1} (1-q)^t \|z_0\|^2}{T} + \frac{\alpha L_\omega}{4T} \sum_{t=0}^{T-1} \sum_{i=0}^t (1-q)^{t-i} \mathbb{E}[\|\nabla J(\theta_i)\|^2] \\ &\quad + \frac{1}{q} \left(4\beta^2 C_{g1} + 4\alpha^2 C_g^2 L_\omega^2 + \frac{\alpha^2 C_g^3 D_\omega}{L_g} \right) \\ &\stackrel{(a)}{\leq} \frac{\|z_0\|^2}{Tq} + \frac{\alpha L_\omega}{4q} \frac{\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]}{T} \\ &\quad + \frac{1}{q} \left(4\beta^2 C_{g1} + 4\alpha^2 C_g^2 L_\omega^2 + \frac{\alpha^2 C_g^3 D_\omega}{L_g} \right) \\ &= \mathcal{O}\left(\frac{1}{T\beta} + \frac{\alpha}{\beta} \frac{\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]}{T} + \beta\right), \end{aligned} \quad (65)$$

547 where (a) is due to the double-sum trick, i.e., for any $x_i \geq 0$, $\sum_{t=0}^{T-1} \sum_{i=0}^t (1-q)^{t-i} x_i \leq \sum_{t=0}^{T-1} (1-q)^t \sum_{t=0}^{T-1} x_t \leq \frac{1}{q} \sum_{t=0}^{T-1} x_t$, and the last step is because $q = \mathcal{O}(\beta)$.

549 B.1.1 Bound on Term (a)

550 In this section we provide the detailed proof of the bound on term (a) in (55).

551 It can be shown that

$$\begin{aligned} &\|(-A_{\theta_t}(s_t)z_t - A_{\theta_t}(s_t)\omega(\theta_t) + \rho_t\delta_t(\theta_t)\phi_{\theta_t}(s_t))\|^2 \\ &\leq 2\| -A_{\theta_t}(s_t)z_t \|^2 + 2\| -A_{\theta_t}(s_t)\omega(\theta_t) + \rho_t\delta_t(\theta_t)\phi_{\theta_t}(s_t) \|^2 \\ &\stackrel{(a)}{\leq} 2C_\phi^2 \|z_t\|^2 + 2 \left(\frac{C_\phi^3}{\lambda_v} (r_{\max} + (1+\gamma)C_v) + \rho_{\max} C_\phi (r_{\max} + (1+\gamma)C_v) \right)^2, \end{aligned} \quad (66)$$

552 where (a) is from the fact that $\|A_\theta(s)\| = \|\phi_\theta(s)\phi_\theta(s)^\top\| \leq C_\phi^2$ and the bounds in (42) and (44).

553 B.1.2 Bound on Term (b)

554 In this section we provide the detailed proof of the bound on term (b) in (56).

555 We first show that $G_{t+1}(\theta, \omega)$ is Lipschitz in ω for any fixed θ . Specifically, for any $\theta, \omega_1, \omega_2 \in \mathbb{R}^N$,
556 it follows that

$$\begin{aligned}
& \|G_{t+1}(\theta, \omega_1) - G_{t+1}(\theta, \omega_2)\| \\
&= \|\rho_t \delta_t(\theta) \phi_\theta(s_t) - \gamma \rho_t \phi_\theta(s'_t) \phi_\theta(s_t)^\top \omega_1 - h_t(\theta, \omega_1) - \rho_t \delta_t(\theta) \phi_\theta(s_t) + \gamma \rho_t \phi_\theta(s'_t) \phi_\theta(s_t)^\top \omega_2 \\
&\quad + h_t(\theta, \omega_2)\| \\
&\leq \|h_t(\theta, \omega_1) - h_t(\theta, \omega_2)\| + \|\gamma \rho_t \phi_\theta(s'_t) \phi_\theta(s_t)^\top \omega_1 - \gamma \rho_t \phi_\theta(s'_t) \phi_\theta(s_t)^\top \omega_2\| \\
&\stackrel{(a)}{\leq} (C_\phi D_v R_\omega + D_v(C_\phi R_\omega + \rho_{\max}(r_{\max} + C_v + \gamma C_v)) + \gamma \rho_{\max} C_\phi^2) \|\omega_1 - \omega_2\| \\
&\triangleq L_g \|\omega_1 - \omega_2\|,
\end{aligned} \tag{67}$$

557 where $L_g = D_v(2C_\phi R_\omega + \rho_{\max}(r_{\max} + C_v + \gamma C_v)) + \gamma \rho_{\max} C_\phi^2$, and (a) is from the Lipschitz
558 continuous of $h_t(\theta, \cdot)$, i.e.,

$$\|h_t(\theta, \omega_1) - h_t(\theta, \omega_2)\| \leq \rho_{\max}(r_{\max} + (1 + \gamma)C_v) D_v \|\omega_1 - \omega_2\| + 2C_\phi D_v R_\omega \|\omega_1 - \omega_2\|. \tag{68}$$

559 We note that to show (67), we use the bound on ω_t , which is guaranteed by the projection step. And
560 this is the only step in our proof where the projection is used.

561 Then it follows that

$$\begin{aligned}
\|\theta_{t+1} - \theta_t\| &= \alpha \|G_{t+1}(\theta_t, \omega_t)\| \\
&\leq \alpha \|G_{t+1}(\theta_t, \omega_t) - G_{t+1}(\theta_t, \omega(\theta_t)) + G_{t+1}(\theta_t, \omega(\theta_t))\| \\
&\leq \alpha L_g \|z_t\| + \alpha \|G_{t+1}(\theta_t, \omega(\theta_t))\| \\
&\leq \alpha L_g \|z_t\| + \alpha C_g,
\end{aligned} \tag{69}$$

562 where $C_g = \rho_{\max} C_\phi (r_{\max} + (1 + \gamma)C_v) + \gamma \rho_{\max} R_\omega C_\phi^2 + D_v R_\omega (R_\omega C_\phi + \rho_{\max}(r_{\max} + C_v + \gamma C_v))$,
563 and the last step in (69) can be shown as follows

$$\begin{aligned}
&\|G_{t+1}(\theta_t, \omega(\theta_t))\| \\
&= \|\rho_t \delta_t(\theta) \phi_\theta(s_t) - \gamma \rho_t \phi_\theta(s'_t) \phi_\theta(s_t)^\top \omega(\theta) - h_t(\theta, \omega(\theta))\| \\
&\leq \rho_{\max} C_\phi (r_{\max} + (1 + \gamma)C_v) + \gamma \rho_{\max} R_\omega C_\phi^2 + D_v R_\omega (R_\omega C_\phi + \rho_{\max}(r_{\max} + C_v + \gamma C_v)).
\end{aligned} \tag{70}$$

564 Using (21) and (69), it follows that

$$\|\omega(\theta_t) - \omega(\theta_{t+1})\| \leq L_\omega \|\theta_{t+1} - \theta_t\| \leq \alpha L_\omega L_g \|z_t\| + \alpha C_g L_\omega, \tag{71}$$

565 and

$$\|\omega(\theta_t) - \omega(\theta_{t+1})\|^2 \leq 2\alpha^2 L_\omega^2 L_g^2 \|z_t\|^2 + 2\alpha^2 C_g^2 L_\omega^2. \tag{72}$$

566 This completes the proof for term (b).

567 B.1.3 Bound on Term (c)

568 In this section we provide the detailed proof of the bound on term (c) in (57).

569 Consider the inner product $\langle z_t, \omega(\theta_t) - \omega(\theta_{t+1}) \rangle$. By the Mean-Value Theorem, it follows that

$$\langle z_t, \omega(\theta_t) \rangle - \langle z_t, \omega(\theta_{t+1}) \rangle = \langle z_t, \omega(\theta_t) - \omega(\theta_{t+1}) \rangle = \langle z_t, \nabla \omega(\hat{\theta}_t)(\theta_t - \theta_{t+1}) \rangle, \tag{73}$$

570 where $\hat{\theta}_t = c\theta_t + (1 - c)\theta_{t+1}$ for some $c \in [0, 1]$. Thus, it follows that

$$\begin{aligned}
&\langle z_t, \omega(\theta_t) - \omega(\theta_{t+1}) \rangle \\
&= \langle z_t, \nabla \omega(\hat{\theta}_t)(\theta_t - \theta_{t+1}) \rangle \\
&= -\alpha \langle z_t, \nabla \omega(\hat{\theta}_t) G_{t+1}(\theta_t, \omega_t) \rangle \\
&= -\alpha \left\langle z_t, \nabla \omega(\hat{\theta}_t) \left(G_{t+1}(\theta_t, \omega_t) - G_{t+1}(\theta_t, \omega(\theta_t)) + G_{t+1}(\theta_t, \omega(\theta_t)) + \frac{\nabla J(\theta_t)}{2} \right) \right\rangle
\end{aligned}$$

$$\begin{aligned}
& + \alpha \left\langle z_t, \nabla \omega(\hat{\theta}_t) \frac{\nabla J(\theta_t)}{2} \right\rangle \\
& = -\alpha \left\langle z_t, \nabla \omega(\hat{\theta}_t) (G_{t+1}(\theta_t, \omega_t) - G_{t+1}(\theta_t, \omega(\theta_t))) \right\rangle + \alpha \left\langle z_t, \nabla \omega(\hat{\theta}_t) \frac{\nabla J(\theta_t)}{2} \right\rangle \\
& \quad - \alpha \left\langle z_t, \nabla \omega(\hat{\theta}_t) \left(G_{t+1}(\theta_t, \omega(\theta_t)) + \frac{\nabla J(\theta_t)}{2} \right) \right\rangle \\
& \stackrel{(a)}{\leq} \alpha L_\omega L_g \|z_t\|^2 + \alpha L_\omega \|z_t\| \left\| \frac{\nabla J(\theta_t)}{2} \right\| - \alpha \left\langle z_t, \nabla \omega(\theta_t) \left(G_{t+1}(\theta_t, \omega(\theta_t)) + \frac{\nabla J(\theta_t)}{2} \right) \right\rangle \\
& \quad + \alpha \left\langle z_t, (\nabla \omega(\theta_t) - \nabla \omega(\hat{\theta}_t)) \left(G_{t+1}(\theta_t, \omega(\theta_t)) + \frac{\nabla J(\theta_t)}{2} \right) \right\rangle \\
& \leq \alpha L_\omega L_g \|z_t\|^2 + \frac{1}{2} \alpha L_\omega \|z_t\|^2 + \frac{\alpha L_\omega}{8} \|\nabla J(\theta_t)\|^2 + \alpha \eta_G(\theta_t, z_t, O_t) \\
& \quad + \alpha \|z_t\| \|\nabla \omega(\theta_t) - \nabla \omega(\hat{\theta}_t)\| \left\| G_{t+1}(\theta_t, \omega(\theta_t)) + \frac{\nabla J(\theta_t)}{2} \right\| \\
& \stackrel{(b)}{\leq} \alpha L_\omega L_g \|z_t\|^2 + \frac{1}{2} \alpha L_\omega \|z_t\|^2 + \frac{\alpha L_\omega}{8} \|\nabla J(\theta_t)\|^2 + \alpha \eta_G(\theta_t, z_t, O_t) + 2\alpha C_g D_\omega \|z_t\| \|\theta_t - \hat{\theta}_t\| \\
& \stackrel{(c)}{\leq} \alpha L_\omega L_g \|z_t\|^2 + \frac{1}{2} \alpha L_\omega \|z_t\|^2 + \frac{\alpha L_\omega}{8} \|\nabla J(\theta_t)\|^2 + \alpha \eta_G(\theta_t, z_t, O_t) \\
& \quad + 2\alpha C_g D_\omega \|z_t\| \|\theta_t - \theta_{t+1}\| \\
& \stackrel{(d)}{\leq} \alpha L_\omega L_g \|z_t\|^2 + \frac{1}{2} \alpha L_\omega \|z_t\|^2 + \frac{\alpha L_\omega}{8} \|\nabla J(\theta_t)\|^2 + \alpha \eta_G(\theta_t, z_t, O_t) \\
& \quad + 2\alpha C_g D_\omega \|z_t\| (\alpha L_g \|z_t\| + \alpha C_g) \\
& \stackrel{(e)}{\leq} \alpha L_\omega L_g \|z_t\|^2 + \frac{1}{2} \alpha L_\omega \|z_t\|^2 + \frac{\alpha L_\omega}{8} \|\nabla J(\theta_t)\|^2 + \alpha \eta_G(\theta_t, z_t, O_t) \\
& \quad + 2\alpha^2 C_g D_\omega \left(2L_g \|z_t\|^2 + \frac{C_g^2}{4L_g} \right) \\
& \leq (\alpha L_\omega L_g + \frac{1}{2} \alpha L_\omega + 4\alpha^2 C_g L_g D_\omega) \|z_t\|^2 + \frac{\alpha L_\omega}{8} \|\nabla J(\theta_t)\|^2 + \frac{\alpha^2 C_g^3 D_\omega}{2L_g} + \alpha \eta_G(\theta_t, z_t, O_t),
\end{aligned} \tag{74}$$

571 where $\eta_G(\theta_t, z_t, O_t) = -\left\langle z_t, \nabla \omega(\theta_t) \left(G_{t+1}(\theta_t, \omega(\theta_t)) + \frac{\nabla J(\theta_t)}{2} \right) \right\rangle$, (a) is from the Lipschitz
572 continuity of $G_{t+1}(\theta, \cdot)$ proved in (67), (b) is from the Lipschitz continuity of $\nabla \omega(\theta)$, which is
573 shown in (30), (c) is from the fact that $\|\theta_t - \hat{\theta}_t\| = (1 - c)\|\theta_t - \theta_{t+1}\| \leq \|\theta_t - \theta_{t+1}\|$, (d) is from
574 the bound of $\|\theta_t - \theta_{t+1}\|$ in (69), and (e) is from the fact that $C_g \|z_t\| \leq L_g \|z_t\|^2 + \frac{C_g^2}{4L_g}$.

575 This completes the proof.

576 B.1.4 Bound on Term (d)

577 In this section we provide the detailed proof of the bound on term (d) in (58).

578 It can be shown that

$$\begin{aligned}
-2\beta \langle z_t, A_{\theta_t}(s_t) z_t \rangle & = -2\beta \langle z_t, A_{\theta_t} z_t \rangle + 2\beta \langle z_t, (A_{\theta_t} - A_{\theta_t}(s_t)) z_t \rangle \\
& \leq -2\beta \lambda_v \|z_t\|^2 + 2\beta \langle z_t, (A_{\theta_t} - A_{\theta_t}(s_t)) z_t \rangle,
\end{aligned} \tag{75}$$

579 where the inequality is due to the fact that $\langle z_t, A_{\theta_t} z_t \rangle = z_t^\top A_{\theta_t} z_t \geq \lambda_L(A_{\theta_t}) \|z_t\|^2 \geq \lambda_v \|z_t\|^2$.

580 B.2 Proof under the i.i.d. Setting

581 In this section we provide the proof of Theorem 1 under the i.i.d. setting.

582 From Lemma 3, we know that the objective function $J(\theta)$ is L_J -smooth, hence it follows that

$$\begin{aligned}
J(\theta_{t+1}) &\leq J(\theta_t) + \langle \nabla J(\theta_t), \theta_{t+1} - \theta_t \rangle + \frac{L_J}{2} \|\theta_{t+1} - \theta_t\|^2 \\
&= J(\theta_t) + \alpha \langle \nabla J(\theta_t), G_{t+1}(\theta_t, \omega_t) \rangle + \frac{L_J}{2} \alpha^2 \|G_{t+1}(\theta_t, \omega_t)\|^2 \\
&= J(\theta_t) - \alpha \left\langle \nabla J(\theta_t), -G_{t+1}(\theta_t, \omega_t) - \frac{\nabla J(\theta_t)}{2} + G_{t+1}(\theta_t, \omega(\theta_t)) - G_{t+1}(\theta_t, \omega(\theta_t)) \right\rangle \\
&\quad - \frac{\alpha}{2} \|\nabla J(\theta_t)\|^2 + \frac{L_J}{2} \alpha^2 \|G_{t+1}(\theta_t, \omega_t)\|^2 \\
&= J(\theta_t) - \alpha \langle \nabla J(\theta_t), -G_{t+1}(\theta_t, \omega_t) + G_{t+1}(\theta_t, \omega(\theta_t)) \rangle \\
&\quad + \alpha \left\langle \nabla J(\theta_t), \frac{\nabla J(\theta_t)}{2} + G_{t+1}(\theta_t, \omega(\theta_t)) \right\rangle - \frac{\alpha}{2} \|\nabla J(\theta_t)\|^2 + \frac{L_J}{2} \alpha^2 \|G_{t+1}(\theta_t, \omega_t)\|^2 \\
&\stackrel{(a)}{\leq} J(\theta_t) + \alpha L_g \|\nabla J(\theta_t)\| \|\omega(\theta_t) - \omega_t\| + \alpha \left\langle \nabla J(\theta_t), \frac{\nabla J(\theta_t)}{2} + G_{t+1}(\theta_t, \omega(\theta_t)) \right\rangle \\
&\quad - \frac{\alpha}{2} \|\nabla J(\theta_t)\|^2 + \frac{L_J}{2} \alpha^2 \|G_{t+1}(\theta_t, \omega_t)\|^2 \\
&\stackrel{(b)}{\leq} J(\theta_t) + \alpha L_g \|\nabla J(\theta_t)\| \|z_t\| + \alpha \left\langle \nabla J(\theta_t), \frac{\nabla J(\theta_t)}{2} + G_{t+1}(\theta_t, \omega(\theta_t)) \right\rangle \\
&\quad - \frac{\alpha}{2} \|\nabla J(\theta_t)\|^2 + \frac{L_J}{2} \alpha^2 (2L_g^2 \|z_t\|^2 + 2C_g^2), \tag{76}
\end{aligned}$$

583 where (a) is from (67) and (b) is because $\|\theta_{t+1} - \theta_t\| = \alpha \|G_{t+1}(\theta_t, \omega_t)\| \leq \alpha L_g \|z_t\| + \alpha C_g$, whose
584 detailed proof is provided in (69). Thus by re-arranging the terms, taking expectation and summing
585 up w.r.t. t from 0 to $T-1$, it follows that

$$\begin{aligned}
&\frac{\alpha}{2} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2] \\
&\leq -\mathbb{E}[J(\theta_T)] + J(\theta_0) + \alpha L_g \sqrt{\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]} \sqrt{\sum_{t=0}^{T-1} \mathbb{E}[\|z_t\|^2]} + \alpha^2 L_J L_g^2 \sum_{t=0}^{T-1} \mathbb{E}[\|z_t\|^2] \\
&\quad + \alpha^2 C_g^2 L_J T, \tag{77}
\end{aligned}$$

586 which is due to the fact that under the i.i.d. setting,

$$\begin{aligned}
&\mathbb{E} \left[\left\langle \nabla J(\theta_t), \frac{\nabla J(\theta_t)}{2} + G_{t+1}(\theta_t, \omega(\theta_t)) \right\rangle \right] \\
&= \mathbb{E} \left[\left\langle \nabla J(\theta_t), \mathbb{E} \left[\frac{\nabla J(\theta_t)}{2} + G_{t+1}(\theta_t, \omega(\theta_t)) \middle| \mathcal{F}_t \right] \right\rangle \right] = 0, \tag{78}
\end{aligned}$$

587 and the Cauchy's inequality

$$\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\| \|z_t\|] \leq \sqrt{\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]} \sqrt{\sum_{t=0}^{T-1} \mathbb{E}[\|z_t\|^2]}. \tag{79}$$

588 Thus dividing both sides by $\frac{\alpha T}{2}$, it follows that

$$\begin{aligned}
&\frac{\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]}{T} \\
&\leq \frac{2(J(\theta_0) - J^*)}{T\alpha} + 2L_g \sqrt{\frac{\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]}{T}} \sqrt{\frac{\sum_{t=0}^{T-1} \mathbb{E}[\|z_t\|^2]}{T}} \\
&\quad + 2\alpha L_J L_g^2 \frac{\sum_{t=0}^{T-1} \mathbb{E}[\|z_t\|^2]}{T} + 2\alpha C_g^2 L_J, \tag{80}
\end{aligned}$$

589 where $J^* \triangleq \min_{\theta} J(\theta)$.

590 Recall the tracking error in (65):

$$\begin{aligned} & \frac{\sum_{t=0}^{T-1} \mathbb{E}[\|z_t\|^2]}{T} \\ & \leq \frac{\|z_0\|^2}{Tq} + \frac{\alpha L_\omega}{4q} \frac{\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]}{T} + \frac{1}{q} \left(4\beta^2 C_{g1} + 4\alpha^2 C_g^2 L_\omega^2 + \frac{\alpha^2 C_g^3 D_\omega}{L_g} \right). \end{aligned} \quad (81)$$

591 We then plug in the tracking error and obtain that

$$\begin{aligned} & \frac{\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]}{T} \\ & \leq \frac{2(J(\theta_0) - J^*)}{T\alpha} + 2\alpha C_g^2 L_J + 2L_g \sqrt{\frac{\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]}{T}} \\ & \quad \times \sqrt{\frac{\|z_0\|^2}{Tq} + \alpha L_\omega \frac{1}{4q} \frac{\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]}{T} + \frac{1}{q} \left(4\beta^2 C_{g1} + 4\alpha^2 C_g^2 L_\omega^2 + \frac{\alpha^2 C_g^3 D_\omega}{L_g} \right)} \\ & \quad + 2\alpha L_J L_g^2 \left(\frac{\|z_0\|^2}{Tq} + \alpha L_\omega \frac{1}{4q} \frac{\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]}{T} \right. \\ & \quad \left. + \frac{1}{q} \left(4\beta^2 C_{g1} + 4\alpha^2 C_g^2 L_\omega^2 + \frac{\alpha^2 C_g^3 D_\omega}{L_g} \right) \right) \\ & \leq \frac{2(J(\theta_0) - J^*)}{T\alpha} + 2\alpha C_g^2 L_J + L_g \sqrt{\frac{\alpha L_\omega}{q} \frac{\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]}{T}} \\ & \quad + 2L_g \sqrt{\frac{\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]}{T}} \sqrt{\frac{\|z_0\|^2}{Tq} + \frac{1}{q} \left(4\beta^2 C_{g1} + 4\alpha^2 C_g^2 L_\omega^2 + \frac{\alpha^2 C_g^3 D_\omega}{L_g} \right)} \\ & \quad + 2\alpha L_J L_g^2 \left(\frac{\|z_0\|^2}{Tq} + \alpha L_\omega \frac{1}{4q} \frac{\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]}{T} \right. \\ & \quad \left. + \frac{1}{q} \left(4\beta^2 C_{g1} + 4\alpha^2 C_g^2 L_\omega^2 + \frac{\alpha^2 C_g^3 D_\omega}{L_g} \right) \right), \end{aligned} \quad (82)$$

592 where the last step is from the fact that $\sqrt{x+y} \leq \sqrt{x} + \sqrt{y}$ for any $x, y \geq 0$. Re-arranging the
593 terms, it follows that

$$\begin{aligned} & \left(1 - L_g \sqrt{\frac{\alpha L_\omega}{q}} - \frac{\alpha^2 L_J L_g^2 L_\omega}{2q} \right) \frac{\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]}{T} \\ & \leq \frac{2(J(\theta_0) - J^*)}{T\alpha} + 2\alpha C_g^2 L_J + 2\alpha L_J L_g^2 \left(\frac{\|z_0\|^2}{Tq} + \frac{1}{q} \left(4\beta^2 C_{g1} + 4\alpha^2 C_g^2 L_\omega^2 + \frac{\alpha^2 C_g^3 D_\omega}{L_g} \right) \right) \\ & \quad + 2L_g \sqrt{\frac{\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]}{T}} \sqrt{\frac{\|z_0\|^2}{Tq} + \frac{1}{q} \left(4\beta^2 C_{g1} + 4\alpha^2 C_g^2 L_\omega^2 + \frac{\alpha^2 C_g^3 D_\omega}{L_g} \right)}. \end{aligned} \quad (83)$$

594 Note that $\left(L_g \sqrt{\frac{\alpha L_\omega}{q}} + \frac{\alpha^2 L_J L_g^2 L_\omega}{2q} \right) = \mathcal{O} \left(\sqrt{\frac{\alpha}{\beta}} + \frac{\alpha^2}{\beta} \right)$, hence we can choose α and β such that

595 $\left(1 - L_g \sqrt{\frac{\alpha L_\omega}{q}} - \frac{\alpha^2 L_J L_g^2 L_\omega}{2q} \right) \geq \frac{1}{2}$. Thus (83) implies that

$$\begin{aligned} & \frac{\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]}{T} \\ & \leq \frac{4(J(\theta_0) - J^*)}{T\alpha} + 4\alpha C_g^2 L_J + 4\alpha L_J L_g^2 \left(\frac{\|z_0\|^2}{Tq} + \frac{1}{q} \left(4\beta^2 C_{g1} + 4\alpha^2 C_g^2 L_\omega^2 + \frac{\alpha^2 C_g^3 D_\omega}{L_g} \right) \right) \end{aligned}$$

$$+ 4L_g \sqrt{\frac{\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]}{T}} \sqrt{\frac{\|z_0\|^2}{Tq} + \frac{1}{q} \left(4\beta^2 C_{g1} + 4\alpha^2 C_g^2 L_\omega^2 + \frac{\alpha^2 C_g^3 D_\omega}{L_g} \right)}. \quad (84)$$

596 Denote $U = \frac{4(J(\theta_0) - J^*)}{T\alpha} + 4\alpha C_g^2 L_J + 4\alpha L_J L_g^2 \left(\frac{\|z_0\|^2}{Tq} + \frac{1}{q} \left(4\beta^2 C_{g1} + 4\alpha^2 C_g^2 L_\omega^2 + \frac{\alpha^2 C_g^3 D_\omega}{L_g} \right) \right)$,
597 and $V = 4L_g \sqrt{\frac{\|z_0\|^2}{Tq} + \frac{1}{q} \left(4\beta^2 C_{g1} + 4\alpha^2 C_g^2 L_\omega^2 + \frac{\alpha^2 C_g^3 D_\omega}{L_g} \right)}$. Then it follows that

$$\frac{\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]}{T} \leq V \sqrt{\frac{\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]}{T}} + U, \quad (85)$$

598 which further implies that

$$\begin{aligned} & \frac{\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]}{T} \\ & \leq V^2 + 2U \\ & = 16L_g^2 \left(\frac{\|z_0\|^2}{Tq} + \frac{1}{q} \left(4\beta^2 C_{g1} + 4\alpha^2 C_g^2 L_\omega^2 + \frac{\alpha^2 C_g^3 D_\omega}{L_g} \right) \right) + \frac{8(J(\theta_0) - J^*)}{T\alpha} \\ & \quad + 8\alpha C_g^2 L_J + 8\alpha L_J L_g^2 \left(\frac{\|z_0\|^2}{Tq} + \frac{1}{q} \left(4\beta^2 C_{g1} + 4\alpha^2 C_g^2 L_\omega^2 + \frac{\alpha^2 C_g^3 D_\omega}{L_g} \right) \right) \\ & = (16L_g^2 + 8\alpha L_J L_g^2) \left(\frac{\|z_0\|^2}{Tq} + \frac{1}{q} \left(4\beta^2 C_{g1} + 4\alpha^2 C_g^2 L_\omega^2 + \frac{\alpha^2 C_g^3 D_\omega}{L_g} \right) \right) \\ & \quad + \frac{8(J(\theta_0) - J^*)}{T\alpha} + 8\alpha C_g^2 L_J \\ & = \mathcal{O} \left(\frac{1}{T\beta} + \beta + \frac{1}{T\alpha} \right) \\ & = \mathcal{O} \left(\frac{1}{T^{1-a}} + \frac{1}{T^b} + \frac{1}{T^{1-b}} \right). \end{aligned} \quad (86)$$

599 This completes the proof.

600 C Non-asymptotic Analysis under the Markovian Setting

601 In this section we provide the proof of Theorem 1 under that Markovian setting. In Appendix C.1 we
602 develop the finite-time analysis of the tracking error and in Appendix C.2 we prove Theorem 1.

603 C.1 Tracking Error Analysis under the Markovian Setting

604 We first define the mixing time $\tau_\beta = \inf \{t : m\kappa^t \leq \beta\}$ (Assumption 4). It can be shown that for
605 any bounded function $\|f(O_t)\| \leq C_f$, for any $t \geq \tau_\beta$, $\|\mathbb{E}[f(O_t)] - \mathbb{E}_{O \sim \mu^{\pi_b}}[f(O)]\| \leq C_f \beta$ and
606 $\tau_\beta = \mathcal{O}(-\log \beta)$. We note that $\beta \tau_\beta \rightarrow 0$ as $\beta \rightarrow 0$, and we assume that $\beta \tau_\beta C_\phi^2 \leq \frac{1}{4}$.

607 From (53), the update of the tracking error z_t can be written as

$$z_{t+1} = z_t + \beta(-A_{\theta_t}(s_t)z_t + b_t(\theta_t)) + \omega(\theta_t) - \omega(\theta_{t+1}), \quad (87)$$

608 where $A_{\theta_t}(s_t) = \phi_{\theta_t}(s_t)\phi_{\theta_t}(s_t)^\top$ and $b_t(\theta_t) = -A_{\theta_t}(s_t)\omega(\theta_t) + \rho_t \delta_t(\theta_t)\phi_{\theta_t}(s_t)$. Note that for
609 any $\theta \in \mathbb{R}^N$ and any sample $O_t = (s_t, a_t, r_t, s_{t+1}) \in \mathcal{S} \times \mathcal{A} \times \mathbb{R} \times \mathcal{S}$, $\|b_t(\theta_t)\| \leq C_\phi^2 R_\omega +$
610 $\rho_{\max} C_\phi (r_{\max} + C_v + \gamma C_v) \triangleq b_{\max}$.

611 Then it can be shown that

$$\begin{aligned} & \mathbb{E} [\|z_{t+1}\|^2 - \|z_t\|^2] \\ & = \mathbb{E} [2z_t^\top (z_{t+1} - z_t) + \|z_{t+1} - z_t\|^2] \end{aligned}$$

$$\begin{aligned}
&= \mathbb{E} [2z_t^\top (z_{t+1} - z_t + \beta A_{\theta_t} z_t)] + \mathbb{E} [\|z_{t+1} - z_t\|^2] + \beta \mathbb{E} [2z_t^\top (-A_{\theta_t}) z_t] \\
&\leq \underbrace{\mathbb{E} [\|z_{t+1} - z_t\|^2]}_{(a)} + \underbrace{\mathbb{E} [2z_t^\top (z_{t+1} - z_t + \beta A_{\theta_t} z_t)] - 2\beta \lambda_v \mathbb{E} [\|z_t\|^2]}_{(b)}, \tag{88}
\end{aligned}$$

612 where the last inequality is due to the fact that $\lambda_L(A_{\theta_t}) \geq \lambda_v$. We first provide the bounds on terms
613 (a) and (b) as follows, and their detailed proof can be found in Appendices C.1.1 and C.1.2.

614 **Term (a) can be bounded as follows:**

615 For any $t \geq 0$, we have that

$$\|z_{t+1} - z_t\|^2 \leq 2\beta^2 C_\phi^4 \|z_t\|^2 + 2\beta^2 (b_{\max} + L_\omega C_g)^2. \tag{89}$$

616 **Term (b) can be bounded as follows:**

617 For any $t \geq \tau_\beta$, we have that

$$\begin{aligned}
&\left| \mathbb{E} \left[z_t^\top \left(-A_{\theta_t} z_t - \frac{1}{\beta} (z_{t+1} - z_t) \right) \right] \right| \\
&\leq (R_1 + R_3 + P_1 + P_2 + P_3) \mathbb{E} [\|z_t\|^2] + (Q_1 + Q_2 + Q_3 + P_1 + P_2 + P_3) \\
&\quad + \frac{\alpha}{8\beta} L_\omega \mathbb{E} [\|\nabla J(\theta_t)\|^2], \tag{90}
\end{aligned}$$

618 where the definition of P_i, Q_i and R_i , $i = 1, 2, 3$, can be found in (111), (114) and (117).

619 From (88), it can be shown that for any $t \geq \tau_\beta$,

$$\begin{aligned}
&\mathbb{E} [\|z_{t+1}\|^2 - \|z_t\|^2] \\
&\leq 2\beta(R_1 + R_3 + P_1 + P_2 + P_3) \mathbb{E} [\|z_t\|^2] + 2\beta(Q_1 + Q_2 + Q_3 + P_1 + P_2 + P_3) \\
&\quad + \frac{\alpha}{4} L_\omega \mathbb{E} [\|\nabla J(\theta_t)\|^2] + 2\beta^2 C_\phi^4 \mathbb{E} [\|z_t\|^2] + 2\beta^2 (b_{\max} + L_\omega C_g)^2 - 2\beta \lambda_v \mathbb{E} [\|z_t\|^2]. \tag{91}
\end{aligned}$$

620 Thus by re-arranging the terms we obtain that

$$\begin{aligned}
&\mathbb{E} [\|z_{t+1}\|^2] \\
&\leq (1 - 2\beta \lambda_v + 2\beta(R_1 + R_3 + P_1 + P_2 + P_3) + 2\beta^2 C_\phi^4) \mathbb{E} [\|z_t\|^2] + \frac{\alpha}{4} L_\omega \mathbb{E} [\|\nabla J(\theta_t)\|^2] \\
&\quad + 2\beta(Q_1 + Q_2 + Q_3 + P_1 + P_2 + P_3) + 2\beta^2 (b_{\max} + L_\omega C_g)^2 \\
&\triangleq (1 - q) \mathbb{E} [\|z_t\|^2] + \frac{\alpha}{4} L_\omega \mathbb{E} [\|\nabla J(\theta_t)\|^2] + p, \tag{92}
\end{aligned}$$

621 where $q = 2\beta \lambda_v - 2\beta(R_1 + R_3 + P_1 + P_2 + P_3) - 2\beta^2 C_\phi^4 = \mathcal{O}(\beta)$ and $p = 2\beta(Q_1 + Q_2 + Q_3 + P_1 + P_2 + P_3) + 2\beta^2 (b_{\max} + L_\omega C_g)^2 = \mathcal{O}(\beta^2 \tau_\beta)$. Then by recursively using the previous inequality,
622 it follows that for any $t \geq \tau_\beta$,

$$\mathbb{E} [\|z_t\|^2] \leq (1 - q)^{t - \tau_\beta} \mathbb{E} [\|z_{\tau_\beta}\|^2] + \frac{\alpha L_\omega}{4} \sum_{j=0}^{t-\tau_\beta} (1 - q)^{t-j} \mathbb{E} [\|\nabla J(\theta_j)\|^2] + \frac{p}{q}, \tag{93}$$

624 and hence

$$\begin{aligned}
&\frac{\sum_{t=0}^{T-1} \mathbb{E} [\|z_t\|^2]}{T} \\
&= \frac{\sum_{t=\tau_\beta}^{T-1} \mathbb{E} [\|z_t\|^2]}{T} + \frac{\sum_{t=0}^{\tau_\beta-1} \mathbb{E} [\|z_t\|^2]}{T} \\
&\leq \frac{\mathbb{E} [\|z_{\tau_\beta}\|^2]}{Tq} + \frac{\tau_\beta (2\|z_0\| + 2\beta \tau_\beta (b_{\max} + L_\omega C_g))^2}{T} + \frac{\alpha L_\omega}{4q} \frac{\sum_{t=0}^{T-1} \mathbb{E} [\|\nabla J(\theta_t)\|^2]}{T} + \frac{p}{q} \\
&\leq (2\|z_0\| + 2\beta \tau_\beta (b_{\max} + L_\omega C_g))^2 \left(\frac{1}{Tq} + \frac{\tau_\beta}{T} \right) + \frac{\alpha L_\omega}{4q} \frac{\sum_{t=0}^{T-1} \mathbb{E} [\|\nabla J(\theta_t)\|^2]}{T} + \frac{p}{q} \\
&= \mathcal{O} \left(\frac{1}{T\beta} + \frac{\alpha}{\beta} \frac{\sum_{t=0}^{T-1} \mathbb{E} [\|\nabla J(\theta_t)\|^2]}{T} + \beta \tau_\beta \right), \tag{94}
\end{aligned}$$

625 where the last step is because $q = \mathcal{O}(\beta)$ and $p = \mathcal{O}(\beta^2 \tau_\beta)$.

626 **C.1.1 Bound on Term (a)**

627 In this section we provide the detailed proof of the bound on term (a) in (88).

628 We first note that from the update of z_t in (87), term $\|z_{t+1} - z_t\|$ can be bounded as follows

$$\begin{aligned}\|z_{t+1} - z_t\| &\leq \|\beta(-A_{\theta_t}(s_t)z_t + b_t(\theta_t))\| + \|\omega(\theta_t) - \omega(\theta_{t+1})\| \\ &\leq \beta C_\phi^2 \|z_t\| + \beta b_{\max} + L_\omega \|\theta_t - \theta_{t+1}\| \\ &\stackrel{(a)}{\leq} \beta C_\phi^2 \|z_t\| + \beta b_{\max} + \alpha L_\omega C_g \\ &\leq \beta C_\phi^2 \|z_t\| + \beta(b_{\max} + L_\omega C_g),\end{aligned}\tag{95}$$

629 where (a) is due to the fact $\|G_{t+1}(\theta_t, \omega_t)\| \leq C_g$ for any $t \geq 0$, and where the last inequality is from
630 the fact that $\alpha \leq \beta$. Hence term (a) can be bounded as follows

$$\|z_{t+1} - z_t\|^2 \leq 2\beta^2 C_\phi^4 \|z_t\|^2 + 2\beta^2(b_{\max} + L_\omega C_g)^2.\tag{96}$$

631 This completes the proof.

632 **C.1.2 Bound on Term (b)**

633 In this section we provide the detailed proof of the bound on term (b) in (88).

634 From (95), it follows that

$$\begin{aligned}\|z_{t+1}\| &\leq (1 + \beta C_\phi^2) \|z_t\| + \beta b_{\max} + \alpha L_\omega C_g \\ &\leq (1 + \beta C_\phi^2) \|z_t\| + \beta(b_{\max} + L_\omega C_g).\end{aligned}\tag{97}$$

635 By applying (97) recursively, it follows that

$$\begin{aligned}\|z_t\| &\leq (1 + \beta C_\phi^2)^t \|z_0\| + \beta(b_{\max} + L_\omega C_g) \frac{(1 + \beta C_\phi^2)^t - 1}{\beta C_\phi^2} \\ &= (1 + \beta C_\phi^2)^t \|z_0\| + (b_{\max} + L_\omega C_g) \frac{(1 + \beta C_\phi^2)^t - 1}{C_\phi^2}.\end{aligned}\tag{98}$$

636 We first show the following lemma which bounds the update $\|z_t - z_{t-\tau_\beta}\|$ by $\|z_t\|$.

637 **Lemma 4.** For any $t \geq \tau_\beta$ and $t \geq j \geq t - \tau_\beta$, we have that

$$\|z_j\| \leq 2\|z_{t-\tau_\beta}\| + 2\beta\tau_\beta(b_{\max} + L_\omega C_g);\tag{99}$$

$$\|z_t - z_{t-\tau_\beta}\| \leq 2\beta\tau_\beta C_\phi^2 \|z_{t-\tau_\beta}\| + 2\beta\tau_\beta(b_{\max} + L_\omega C_g),\tag{100}$$

$$\|z_t - z_{t-\tau_\beta}\| \leq 4\beta\tau_\beta C_\phi^2 \|z_t\| + 4\beta\tau_\beta(b_{\max} + L_\omega C_g).\tag{101}$$

638 *Proof.* From (97), it follows that

$$\|z_{t+1}\| \leq (1 + \beta C_\phi^2) \|z_t\| + \beta(b_{\max} + L_\omega C_g).\tag{102}$$

639 First note that $\beta C_\phi^2 \tau_\beta \leq \frac{1}{4}$ and hence $\beta C_\phi^2 \leq \frac{1}{4\tau_\beta} \leq \frac{\log 2}{\tau_\beta - 1}$. This implies that

$$(1 + \beta C_\phi^2)^{\tau_\beta} \leq 1 + 2\tau_\beta \beta C_\phi^2,\tag{103}$$

640 which is because $(1 + x)^k \leq 1 + 2kx$ for $x \leq \frac{\log 2}{k-1}$.

641 Applying inequality (102) recursively, it follows that

$$\begin{aligned}\|z_j\| &\leq (1 + \beta C_\phi^2)^{j-t+\tau_\beta} \|z_{t-\tau_\beta}\| + (b_{\max} + L_\omega C_g) \frac{(1 + \beta C_\phi^2)^{\tau_\beta} - 1}{C_\phi^2} \\ &\leq (1 + \beta C_\phi^2)^{\tau_\beta} \|z_{t-\tau_\beta}\| + (b_{\max} + L_\omega C_g) \frac{(1 + \beta C_\phi^2)^{\tau_\beta} - 1}{C_\phi^2} \\ &\stackrel{(a)}{\leq} (1 + 2\tau_\beta \beta C_\phi^2) \|z_{t-\tau_\beta}\| + 2\beta\tau_\beta(b_{\max} + L_\omega C_g)\end{aligned}$$

$$\stackrel{(b)}{\leq} 2\|z_{t-\tau_\beta}\| + 2\beta\tau_\beta(b_{\max} + L_\omega C_g), \quad (104)$$

642 where (a) is from (103), and (b) is from the fact that $\beta\tau_\beta C_\phi^2 \leq \frac{1}{4}$.

643 To prove (100) and (101), first note that

$$\begin{aligned} \|z_t - z_{t-\tau_\beta}\| &\leq \sum_{j=t-\tau_\beta}^{t-1} \|z_{j+1} - z_j\| \\ &\stackrel{(a)}{\leq} \sum_{j=t-\tau_\beta}^{t-1} \beta C_\phi^2 \|z_j\| + \beta\tau_\beta(b_{\max} + L_\omega C_g) \\ &\stackrel{(b)}{\leq} \sum_{j=t-\tau_\beta}^{t-1} \beta C_\phi^2 (2\|z_{t-\tau_\beta}\| + 2\beta\tau_\beta(b_{\max} + L_\omega C_g)) + \beta\tau_\beta(b_{\max} + L_\omega C_g) \\ &\leq \beta\tau_\beta C_\phi^2 (2\|z_{t-\tau_\beta}\| + 2\beta\tau_\beta(b_{\max} + L_\omega C_g)) + \beta\tau_\beta(b_{\max} + L_\omega C_g) \\ &= 2\beta\tau_\beta C_\phi^2 \|z_{t-\tau_\beta}\| + (2\beta^2\tau_\beta^2 C_\phi^2 + \beta\tau_\beta)(b_{\max} + L_\omega C_g) \\ &\stackrel{(c)}{\leq} 2\beta\tau_\beta C_\phi^2 \|z_{t-\tau_\beta}\| + 2\beta\tau_\beta(b_{\max} + L_\omega C_g), \end{aligned} \quad (105)$$

644 where (a) is from (95), (b) is from (104) and (c) is due to the fact that $\beta\tau_\beta C_\phi^2 \leq \frac{1}{4}$. Moreover, it
645 can be further shown that

$$\begin{aligned} \|z_t - z_{t-\tau_\beta}\| &\leq 2\beta\tau_\beta C_\phi^2 (\|z_t\| + \|z_t - z_{t-\tau_\beta}\|) + 2\beta\tau_\beta(b_{\max} + L_\omega C_g) \\ &\leq 2\beta\tau_\beta C_\phi^2 \|z_t\| + \frac{1}{2} \|z_t - z_{t-\tau_\beta}\| + 2\beta\tau_\beta(b_{\max} + L_\omega C_g), \end{aligned} \quad (106)$$

646 where the last step is because $\beta\tau_\beta C_\phi^2 \leq \frac{1}{4}$. Hence

$$\|z_t - z_{t-\tau_\beta}\| \leq 4\beta\tau_\beta C_\phi^2 \|z_t\| + 4\beta\tau_\beta(b_{\max} + L_\omega C_g). \quad (107)$$

647 \square

648 The bound on term (b) in (88) is straightforward from the following lemma.

649 **Lemma 5.** For any $t \geq \tau_\beta$, it follows that

$$\begin{aligned} &\left| \mathbb{E} \left[z_t^\top \left(-A_{\theta_t} z_t - \frac{1}{\beta} (z_{t+1} - z_t) \right) \right] \right| \\ &\leq (R_1 + R_3 + P_1 + P_2 + P_3) \mathbb{E} [\|z_t\|^2] + (Q_1 + Q_2 + Q_3 + P_1 + P_2 + P_3) \\ &\quad + \frac{\alpha}{8\beta} L_\omega \mathbb{E} [\|\nabla J(\theta_t)\|^2], \end{aligned} \quad (108)$$

650 where the definition of P_i, Q_i and R_i , $i = 1, 2, 3$, can be found in (111), (114) and (117).

651 *Proof.* We only prove the case $t = \tau_\beta$ here. The proof for the general case with $t > \tau_\beta$ is similar,
652 and thus is omitted here. First note that

$$\begin{aligned} &\mathbb{E} \left[z_{\tau_\beta}^\top \left(-A_{\theta_{\tau_\beta}} z_{\tau_\beta} - \frac{1}{\beta} (z_{\tau_\beta+1} - z_{\tau_\beta}) \right) \right] \\ &= \mathbb{E} \left[z_{\tau_\beta}^\top \left(-A_{\theta_{\tau_\beta}} + A_{\theta_{\tau_\beta}}(s_{\tau_\beta}) \right) z_{\tau_\beta} \right] - \mathbb{E} \left[z_{\tau_\beta}^\top b_{\tau_\beta} \right] - \mathbb{E} \left[z_{\tau_\beta}^\top \frac{\omega(\theta_{\tau_\beta}) - \omega(\theta_{\tau_\beta+1})}{\beta} \right]. \end{aligned} \quad (109)$$

653 We then bound the terms in (109) one by one. First, it can be shown that

$$\begin{aligned} &\left| \mathbb{E} \left[z_{\tau_\beta}^\top \left(-A_{\theta_{\tau_\beta}} + A_{\theta_{\tau_\beta}}(s_{\tau_\beta}) \right) z_{\tau_\beta} \right] \right| \\ &\leq \left| \mathbb{E} \left[z_0^\top \left(-A_{\theta_{\tau_\beta}} + A_{\theta_{\tau_\beta}}(s_{\tau_\beta}) \right) z_0 \right] \right| + \left| \mathbb{E} \left[(z_{\tau_\beta} - z_0)^\top \left(-A_{\theta_{\tau_\beta}} + A_{\theta_{\tau_\beta}}(s_{\tau_\beta}) \right) (z_{\tau_\beta} - z_0) \right] \right| \end{aligned}$$

$$\begin{aligned}
& + 2 \left| \mathbb{E} \left[(z_{\tau_\beta} - z_0)^\top (-A_{\theta_{\tau_\beta}} + A_{\theta_{\tau_\beta}}(s_{\tau_\beta})) z_0 \right] \right| \\
& \leq \|z_0\|^2 \left\| \mathbb{E} \left[-A_{\theta_{\tau_\beta}} + A_{\theta_{\tau_\beta}}(s_{\tau_\beta}) \right] \right\| + 2C_\phi^2 \mathbb{E} [\|z_{\tau_\beta} - z_0\|^2] + 4\|z_0\|C_\phi^2 \mathbb{E} [\|z_{\tau_\beta} - z_0\|] \\
& \leq \|z_0\|^2 \left\| \mathbb{E} \left[-A_{\theta_0} + A_{\theta_0}(s_{\tau_\beta}) \right] \right\| + \|z_0\|^2 \left\| \mathbb{E} \left[-A_{\theta_0} + A_{\theta_{\tau_\beta}} \right] \right\| \\
& \quad + \|z_0\|^2 \left\| \mathbb{E} \left[-A_{\theta_{\tau_\beta}}(s_{\tau_\beta}) + A_{\theta_0}(s_{\tau_\beta}) \right] \right\| + 2C_\phi^2 \mathbb{E} [\|z_{\tau_\beta} - z_0\|^2] + 4\|z_0\|C_\phi^2 \mathbb{E} [\|z_{\tau_\beta} - z_0\|] \\
& \stackrel{(a)}{\leq} (\beta C_\phi^2 + 4C_\phi D_v C_g \alpha \tau_\beta) \|z_0\|^2 + 2C_\phi^2 \mathbb{E} [\|z_{\tau_\beta} - z_0\|^2] + 4\|z_0\|C_\phi^2 \mathbb{E} [\|z_{\tau_\beta} - z_0\|], \tag{110}
\end{aligned}$$

654 where (a) is due to the facts that $\|\mathbb{E}[-A_{\theta_0} + A_{\theta_0}(s_{\tau_\beta})]\| \leq C_\phi^2 \beta$ from the uniform ergodicity of the
655 MDP, both A_θ and $A_\theta(s_{\tau_\beta})$ are Lipschitz with constant $2C_\phi D_v$, and $\|\theta_0 - \theta_{\tau_\beta}\| \leq \sum_{j=0}^{\tau_\beta-1} \|\theta_{j+1} - \theta_j\| \leq \alpha \tau_\beta C_g$.

657 We then plug in the results from Lemma 4, and hence we have that

$$\begin{aligned}
& \left| \mathbb{E} \left[z_{\tau_\beta}^\top (-A_{\theta_{\tau_\beta}} + A_{\theta_{\tau_\beta}}(s_{\tau_\beta})) z_{\tau_\beta} \right] \right| \\
& \leq (\beta C_\phi^2 + 4C_\phi D_v C_g \alpha \tau_\beta) \|z_0\|^2 + 2C_\phi^2 \mathbb{E} [\|z_{\tau_\beta} - z_0\|^2] + 4\|z_0\|C_\phi^2 \mathbb{E} [\|z_{\tau_\beta} - z_0\|] \\
& \stackrel{(a)}{\leq} (\beta C_\phi^2 + 4C_\phi D_v C_g \alpha \tau_\beta) \left(2(1 + 4\beta \tau_\beta C_\phi^2)^2 \mathbb{E} [\|z_{\tau_\beta}\|^2] + 32\beta^2 \tau_\beta^2 (b_{\max} + L_\omega C_g)^2 \right) \\
& \quad + 2C_\phi^2 \left(32\beta^2 \tau_\beta^2 C_\phi^4 \mathbb{E} [\|z_{\tau_\beta}\|^2] + 32\beta^2 \tau_\beta^2 (b_{\max} + L_\omega C_g)^2 \right) \\
& \quad + 4C_\phi^2 \left(4\beta \tau_\beta C_\phi^2 (1 + 4\beta \tau_\beta C_\phi^2) \mathbb{E} [\|z_{\tau_\beta}\|^2] + 4\beta \tau_\beta (b_{\max} + L_\omega C_g) (1 + 8\beta \tau_\beta C_\phi^2) \mathbb{E} [\|z_{\tau_\beta}\|] \right) \\
& \quad + 64C_\phi^2 \beta^2 \tau_\beta^2 (b_{\max} + L_\omega C_g)^2 \\
& \triangleq R_1 \mathbb{E} [\|z_{\tau_\beta}\|^2] + P_1 \mathbb{E} [\|z_{\tau_\beta}\|] + Q_1, \tag{111}
\end{aligned}$$

658 where (a) is from (101) and the fact that

$$\|z_0\| \leq \|z_{\tau_\beta} - z_0\| + \|z_{\tau_\beta}\| \leq (1 + 4\beta \tau_\beta C_\phi^2) \|z_{\tau_\beta}\| + 4\beta \tau_\beta (b_{\max} + L_\omega C_g); \tag{112}$$

659 and $R_1 = 2(1 + 4\beta \tau_\beta C_\phi^2)^2 (\beta C_\phi^2 + 4C_\phi D_v C_g \alpha \tau_\beta) + 64\beta^2 \tau_\beta^2 C_\phi^6 + 16\beta \tau_\beta C_\phi^4 (1 + 4\beta \tau_\beta C_\phi^2) =$
660 $\mathcal{O}(\beta \tau_\beta)$, $P_1 = 16C_\phi^2 \beta \tau_\beta (b_{\max} + L_\omega C_g) (1 + 8\beta \tau_\beta C_\phi^2) = \mathcal{O}(\beta \tau_\beta)$ and $Q_1 =$
661 $(\beta C_\phi^2 + 4C_\phi D_v C_g \alpha \tau_\beta) 32\beta^2 \tau_\beta^2 (b_{\max} + L_\omega C_g)^2 + 64C_\phi^2 \beta^2 \tau_\beta^2 (b_{\max} + L_\omega C_g)^2 + 64C_\phi^2 \beta^2 \tau_\beta^2 (b_{\max} + L_\omega C_g)^2 = \mathcal{O}(\beta^2 \tau^2)$.

663 Similarly, the second term in (109) can be bounded as follows

$$\begin{aligned}
\left| \mathbb{E} \left[z_{\tau_\beta}^\top b_{\tau_\beta}(\theta_{\tau_\beta}) \right] \right| & \leq \left| \mathbb{E} \left[(z_{\tau_\beta} - z_0)^\top b_{\tau_\beta}(\theta_{\tau_\beta}) \right] \right| + \left| \mathbb{E} \left[z_0^\top b_{\tau_\beta}(\theta_0) \right] \right| \\
& \quad + \left| \mathbb{E} \left[z_0^\top (b_{\tau_\beta}(\theta_{\tau_\beta}) - b_{\tau_\beta}(\theta_0)) \right] \right| \\
& \leq b_{\max} \mathbb{E} [\|z_{\tau_\beta} - z_0\|] + \beta b_{\max} \|z_0\| + \alpha \tau_\beta C_g L_b \|z_0\|, \tag{113}
\end{aligned}$$

664 where $L_b = 2C_\phi D_v R_\omega + L_\omega C_\phi^2 + \rho_{\max}((1 + \gamma)C_\phi^2 + D_v(r_{\max} + (1 + \gamma)C_v))$ is the Lipschitz
665 constant of $b_t(\theta)$. Again applying Lemma 4 implies that

$$\begin{aligned}
& \left| \mathbb{E} \left[z_{\tau_\beta}^\top b_{\tau_\beta}(\theta_{\tau_\beta}) \right] \right| \\
& \leq b_{\max} \mathbb{E} [\|z_{\tau_\beta} - z_0\|] + \beta b_{\max} \|z_0\| + \alpha \tau_\beta C_g L_b \|z_0\| \\
& \leq b_{\max} (4\beta \tau_\beta C_\phi^2 \mathbb{E} [\|z_{\tau_\beta}\|] + 4\beta \tau_\beta (b_{\max} + L_\omega C_g)) \\
& \quad + (\beta b_{\max} + \alpha \tau_\beta C_g L_b) ((1 + 4\beta \tau_\beta C_\phi^2) \mathbb{E} [\|z_{\tau_\beta}\|] + 4\beta \tau_\beta (b_{\max} + L_\omega C_g)) \\
& \triangleq P_2 \mathbb{E} [\|z_{\tau_\beta}\|] + Q_2, \tag{114}
\end{aligned}$$

666 where $P_2 = 4\beta \tau_\beta b_{\max} C_\phi^2 + (\beta b_{\max} + \alpha \tau_\beta C_g L_b) (1 + 4\beta \tau_\beta C_\phi^2) = \mathcal{O}(\beta \tau_\beta)$ and $Q_2 = 4\beta \tau_\beta (b_{\max} + L_\omega C_g) (b_{\max} + \beta b_{\max} + \alpha \tau_\beta C_g L_b) = \mathcal{O}(\beta \tau_\beta)$.

668 We then bound the last term in (109) as follows

$$\begin{aligned}
& \left| \mathbb{E} \left[z_{\tau_\beta}^\top \frac{\omega(\theta_{\tau_\beta}) - \omega(\theta_{\tau_\beta+1})}{\beta} \right] \right| \\
& \stackrel{(a)}{=} \left| \frac{1}{\beta} \mathbb{E} [z_{\tau_\beta}^\top \nabla \omega(\hat{\theta}_{\tau_\beta})(\theta_{\tau_\beta+1} - \theta_{\tau_\beta})] \right| \\
& = \left| \frac{\alpha}{\beta} \mathbb{E} [z_{\tau_\beta}^\top \nabla \omega(\hat{\theta}_{\tau_\beta}) G_{\tau_\beta+1}(\theta_{\tau_\beta}, \omega_{\tau_\beta})] \right| \\
& = \left| \frac{\alpha}{\beta} \mathbb{E} \left[z_{\tau_\beta}^\top \nabla \omega(\hat{\theta}_{\tau_\beta}) \left(G_{\tau_\beta+1}(\theta_{\tau_\beta}, \omega_{\tau_\beta}) - G_{\tau_\beta+1}(\theta_{\tau_\beta}, \omega(\theta_{\tau_\beta})) + G_{\tau_\beta+1}(\theta_{\tau_\beta}, \omega(\theta_{\tau_\beta})) \right. \right. \right. \\
& \quad \left. \left. \left. + \frac{\nabla J(\theta_{\tau_\beta})}{2} - \frac{\nabla J(\theta_{\tau_\beta})}{2} \right) \right] \right| \\
& = \left| \frac{\alpha}{\beta} \mathbb{E} \left[z_{\tau_\beta}^\top \nabla \omega(\hat{\theta}_{\tau_\beta}) (G_{\tau_\beta+1}(\theta_{\tau_\beta}, \omega_{\tau_\beta}) - G_{\tau_\beta+1}(\theta_{\tau_\beta}, \omega(\theta_{\tau_\beta}))) \right] \right| \\
& \quad + \left| \frac{\alpha}{\beta} \mathbb{E} \left[z_{\tau_\beta}^\top \nabla \omega(\hat{\theta}_{\tau_\beta}) \left(G_{\tau_\beta+1}(\theta_{\tau_\beta}, \omega(\theta_{\tau_\beta})) + \frac{\nabla J(\theta_{\tau_\beta})}{2} \right) \right] \right| \\
& \quad + \left| \frac{\alpha}{\beta} \mathbb{E} \left[z_{\tau_\beta}^\top \nabla \omega(\hat{\theta}_{\tau_\beta}) \left(-\frac{\nabla J(\theta_{\tau_\beta})}{2} \right) \right] \right| \\
& \stackrel{(b)}{\leq} \frac{\alpha}{\beta} L_\omega L_g \mathbb{E} [\|z_{\tau_\beta}\|^2] + \frac{\alpha}{2\beta} L_\omega \mathbb{E} [\|z_{\tau_\beta}\|^2] + \frac{\alpha}{8\beta} L_\omega \mathbb{E} [\|\nabla J(\theta_{\tau_\beta})\|^2] \\
& \quad + \frac{\alpha}{\beta} \left| \mathbb{E} \left[z_{\tau_\beta}^\top \nabla \omega(\theta_{\tau_\beta}) \left(G_{\tau_\beta+1}(\theta_{\tau_\beta}, \omega(\theta_{\tau_\beta})) + \frac{\nabla J(\theta_{\tau_\beta})}{2} \right) \right] \right| \\
& \quad + \frac{\alpha}{\beta} \left| \mathbb{E} \left[z_{\tau_\beta}^\top (\nabla \omega(\hat{\theta}_{\tau_\beta}) - \nabla \omega(\theta_{\tau_\beta})) \left(G_{\tau_\beta+1}(\theta_{\tau_\beta}, \omega(\theta_{\tau_\beta})) + \frac{\nabla J(\theta_{\tau_\beta})}{2} \right) \right] \right| \\
& \leq \frac{\alpha}{\beta} L_\omega L_g \mathbb{E} [\|z_{\tau_\beta}\|^2] + \frac{\alpha}{2\beta} L_\omega \mathbb{E} [\|z_{\tau_\beta}\|^2] + \frac{\alpha}{8\beta} L_\omega \mathbb{E} [\|\nabla J(\theta_{\tau_\beta})\|^2] \\
& \quad + \frac{\alpha}{\beta} \left| \mathbb{E} \left[z_0^\top \nabla \omega(\theta_{\tau_\beta}) \left(G_{\tau_\beta+1}(\theta_{\tau_\beta}, \omega(\theta_{\tau_\beta})) + \frac{\nabla J(\theta_{\tau_\beta})}{2} \right) \right] \right| \\
& \quad + \frac{\alpha}{\beta} \left| \mathbb{E} \left[(z_{\tau_\beta} - z_0)^\top \nabla \omega(\theta_{\tau_\beta}) \left(G_{\tau_\beta+1}(\theta_{\tau_\beta}, \omega(\theta_{\tau_\beta})) + \frac{\nabla J(\theta_{\tau_\beta})}{2} \right) \right] \right| \\
& \quad + \frac{2\alpha}{\beta} C_g D_\omega \mathbb{E} [\|z_{\tau_\beta}\| \|\theta_{\tau_\beta} - \theta_{\tau_\beta+1}\|] \\
& \leq \frac{\alpha}{\beta} L_\omega L_g \mathbb{E} [\|z_{\tau_\beta}\|^2] + \frac{\alpha}{2\beta} L_\omega \mathbb{E} [\|z_{\tau_\beta}\|^2] + \frac{\alpha}{8\beta} L_\omega \mathbb{E} [\|\nabla J(\theta_{\tau_\beta})\|^2] \\
& \quad + \frac{\alpha}{\beta} \left| \mathbb{E} \left[z_0^\top \nabla \omega(\theta_0) \left(G_{\tau_\beta+1}(\theta_0, \omega(\theta_0)) + \frac{\nabla J(\theta_0)}{2} \right) \right] \right| \\
& \quad + \frac{\alpha}{\beta} \left| \mathbb{E} \left[z_0^\top \left(\nabla \omega(\theta_{\tau_\beta}) \left(G_{\tau_\beta+1}(\theta_{\tau_\beta}, \omega(\theta_{\tau_\beta})) + \frac{\nabla J(\theta_{\tau_\beta})}{2} \right) \right. \right. \right. \\
& \quad \left. \left. \left. - \nabla \omega(\theta_0) \left(G_{\tau_\beta+1}(\theta_0, \omega(\theta_0)) + \frac{\nabla J(\theta_0)}{2} \right) \right) \right] \right| \\
& \quad + \frac{\alpha}{\beta} \left| \mathbb{E} \left[(z_{\tau_\beta} - z_0)^\top \nabla \omega(\theta_{\tau_\beta}) \left(G_{\tau_\beta+1}(\theta_{\tau_\beta}, \omega(\theta_{\tau_\beta})) + \frac{\nabla J(\theta_{\tau_\beta})}{2} \right) \right] \right| \\
& \quad + \frac{2\alpha}{\beta} C_g D_\omega \mathbb{E} [\|z_{\tau_\beta}\| \|\theta_{\tau_\beta} - \theta_{\tau_\beta+1}\|] \\
& \leq \frac{\alpha}{\beta} L_\omega L_g \mathbb{E} [\|z_{\tau_\beta}\|^2] + \frac{\alpha}{8\beta} L_\omega \mathbb{E} [\|\nabla J(\theta_{\tau_\beta})\|^2]
\end{aligned}$$

$$\begin{aligned}
& + \frac{\alpha}{\beta} \|z_0\| L_\omega \left\| \mathbb{E} \left[G_{\tau_\beta+1}(\theta_0, \omega(\theta_0)) + \frac{\nabla J(\theta_{\tau_\beta})}{2} \right] \right\| \\
& + \frac{\alpha}{2\beta} L_\omega \mathbb{E} [\|z_{\tau_\beta}\|^2] + \frac{\alpha}{\beta} \|z_0\| L_k \mathbb{E} [\|\theta_{\tau_\beta} - \theta_0\|] + \frac{2\alpha}{\beta} L_\omega C_g \mathbb{E} [\|z_{\tau_\beta} - z_0\|] \\
& + \frac{2\alpha}{\beta} C_g D_\omega \mathbb{E} [\|z_{\tau_\beta}\| \|\theta_{\tau_\beta} - \theta_{\tau_\beta+1}\|] \\
& \stackrel{(c)}{\leq} \frac{\alpha}{\beta} L_\omega L_g \mathbb{E} [\|z_{\tau_\beta}\|^2] + \frac{\alpha}{2\beta} L_\omega \mathbb{E} [\|z_{\tau_\beta}\|^2] + \frac{\alpha}{8\beta} L_\omega \mathbb{E} [\|\nabla J(\theta_{\tau_\beta})\|^2] + \frac{\alpha}{\beta} \|z_0\| L_\omega C_g \beta \\
& + \frac{\alpha^2}{\beta} \tau_\beta L_k \|z_0\| C_g + \frac{2\alpha}{\beta} L_\omega C_g \mathbb{E} [\|z_{\tau_\beta} - z_0\|] + \frac{2\alpha^2}{\beta} C_g^2 D_\omega \mathbb{E} [\|z_{\tau_\beta}\|] \\
& = \left(\frac{\alpha}{\beta} L_\omega L_g + \frac{\alpha}{2\beta} L_\omega \right) \mathbb{E} [\|z_{\tau_\beta}\|^2] + \frac{2\alpha^2}{\beta} C_g^2 D_\omega \mathbb{E} [\|z_{\tau_\beta}\|] + \frac{\alpha}{8\beta} L_\omega \mathbb{E} [\|\nabla J(\theta_{\tau_\beta})\|^2] \\
& + \left(\alpha L_\omega C_g + \frac{\alpha^2}{\beta} \tau_\beta L_k C_g \right) \|z_0\| + \frac{2\alpha}{\beta} L_\omega C_g \mathbb{E} [\|z_{\tau_\beta} - z_0\|], \tag{115}
\end{aligned}$$

669 where (a) is from the Mean-Value theorem and $\hat{\theta}_{\tau_\beta} = c\theta_{\tau_\beta} + (1-c)\theta_{\tau_\beta+1}$ for some $c \in [0, 1]$, (b)
670 is from Lemmas 1 and 2, (c) is due to the fact that $\left\| \mathbb{E} \left[G_{t+1}(\theta_0, \omega(\theta_0)) + \frac{\nabla J(\theta_0)}{2} \right] \right\| \leq C_g \beta$ for any
671 $t \geq \tau_\beta$ and $\|\theta_{\tau_\beta} - \theta_0\| \leq \alpha \tau_\beta C_g$, and $L_k = 2C_g D_\omega + \left(L_J + \frac{L'_g}{2} \right) L_\omega$ is the Lipschitz constant of
672 $\nabla \omega(\theta) \left(G_{t+1}(\theta, \omega(\theta)) + \frac{\nabla J(\theta)}{2} \right)$, and L'_g is the Lipschitz constant of $G_{t+1}(\theta, \omega(\theta))$.

673 Our next step is to rewrite the bound in (115) using $\|z_{\tau_\beta}\|$. Note that from Lemma 4, we have that

$$\|z_0\| \leq \|z_{\tau_\beta} - z_0\| + \|z_{\tau_\beta}\| \leq (1 + 4\beta \tau_\beta C_\phi^2) \|z_{\tau_\beta}\| + 4\beta \tau_\beta (b_{\max} + L_\omega C_g). \tag{116}$$

674 Plugging in (115), it follows that

$$\begin{aligned}
& \left\| \mathbb{E} \left[z_{\tau_\beta}^\top \frac{\omega(\theta_{\tau_\beta}) - \omega(\theta_{\tau_\beta+1})}{\beta} \right] \right\| \\
& \leq \left(\frac{\alpha}{\beta} L_\omega L_g + \frac{\alpha}{2\beta} L_\omega \right) \mathbb{E} [\|z_{\tau_\beta}\|^2] + \frac{2\alpha^2}{\beta} C_g^2 D_\omega \mathbb{E} [\|z_{\tau_\beta}\|] + \frac{\alpha}{8\beta} L_\omega \mathbb{E} [\|\nabla J(\theta_{\tau_\beta})\|^2] \\
& + \left(\alpha L_\omega C_g + \frac{\alpha^2}{\beta} \tau_\beta L_k C_g \right) \|z_0\| + \frac{2\alpha}{\beta} L_\omega C_g \mathbb{E} [\|z_{\tau_\beta} - z_0\|] \\
& \leq \left(\frac{\alpha}{\beta} L_\omega L_g + \frac{\alpha}{2\beta} L_\omega \right) \mathbb{E} [\|z_{\tau_\beta}\|^2] + \frac{2\alpha^2}{\beta} C_g^2 D_\omega \mathbb{E} [\|z_{\tau_\beta}\|] + \frac{\alpha}{8\beta} L_\omega \mathbb{E} [\|\nabla J(\theta_{\tau_\beta})\|^2] \\
& + \left(\alpha L_\omega C_g + \frac{\alpha^2}{\beta} \tau_\beta L_k C_g \right) ((1 + 4\beta \tau_\beta C_\phi^2) \mathbb{E} [\|z_{\tau_\beta}\|] + 4\beta \tau_\beta (b_{\max} + L_\omega C_g)) \\
& + \frac{2\alpha}{\beta} L_\omega C_g (\mathbb{E} [4\beta \tau_\beta C_\phi^2 \|z_{\tau_\beta}\|] + 4\beta \tau_\beta (b_{\max} + L_\omega C_g)) \\
& = \left(\frac{\alpha}{\beta} L_\omega L_g + \frac{\alpha}{2\beta} L_\omega \right) \mathbb{E} [\|z_{\tau_\beta}\|^2] \\
& + \left(\frac{2\alpha^2}{\beta} C_g^2 D_\omega + \left(\alpha L_\omega C_g + \frac{\alpha^2}{\beta} \tau_\beta L_k C_g \right) (1 + 4\beta \tau_\beta C_\phi^2) + 8\alpha \tau_\beta L_\omega C_g C_\phi^2 \right) \mathbb{E} [\|z_{\tau_\beta}\|] \\
& + \frac{\alpha}{8\beta} L_\omega \mathbb{E} [\|\nabla J(\theta_{\tau_\beta})\|^2] + \left(\alpha L_\omega C_g + \frac{\alpha^2}{\beta} \tau_\beta L_k C_g \right) (4\beta \tau_\beta (b_{\max} + L_\omega C_g)) \\
& + 8\alpha \tau_\beta L_\omega C_g (b_{\max} + L_\omega C_g) \\
& \triangleq R_3 \mathbb{E} [\|z_{\tau_\beta}\|^2] + P_3 \mathbb{E} [\|z_{\tau_\beta}\|] + Q_3 + \frac{\alpha}{8\beta} L_\omega \mathbb{E} [\|\nabla J(\theta_{\tau_\beta})\|^2], \tag{117}
\end{aligned}$$

675 where $R_3 = \left(\frac{\alpha}{\beta} L_\omega L_g + \frac{\alpha}{2\beta} L_\omega \right) = \mathcal{O}\left(\frac{\alpha}{\beta}\right)$, $P_3 = \left(\frac{2\alpha^2}{\beta} C_g^2 D_\omega \right) +$
676 $\left(\alpha L_\omega C_g + \frac{\alpha^2}{\beta} \tau_\beta L_k C_g \right) (1 + 4\beta \tau_\beta C_\phi^2) + 8\alpha \tau_\beta L_\omega C_g C_\phi^2 \right) = \mathcal{O}(\alpha \tau_\beta)$ and $Q_3 =$
677 $\left(\alpha L_\omega C_g + \frac{\alpha^2}{\beta} \tau_\beta L_k C_g \right) (4\beta \tau_\beta (b_{\max} + L_\omega C_g)) + 8\alpha \tau_\beta L_\omega C_g (b_{\max} + L_\omega C_g) = \mathcal{O}(\alpha \tau_\beta)$.

678 Then we combine all three bounds in (111), (114) and (117), and it follows that

$$\begin{aligned} & \left| \mathbb{E} \left[z_{\tau_\beta}^\top \left(-A_{\theta_{\tau_\beta}} z_{\tau_\beta} - \frac{1}{\beta} (z_{\tau_\beta+1} - z_{\tau_\beta}) \right) \right] \right| \\ & \leq (R_1 + R_3) \mathbb{E} [\|z_{\tau_\beta}\|^2] + (P_1 + P_2 + P_3) \mathbb{E} [\|z_{\tau_\beta}\|] + (Q_1 + Q_2 + Q_3) \\ & \quad + \frac{\alpha}{8\beta} L_\omega \mathbb{E} [\|\nabla J(\theta_{\tau_\beta})\|^2], \end{aligned} \tag{118}$$

679 Finally due to the fact that $x \leq x^2 + 1, \forall x \in \mathbb{R}$, it follows that

$$\begin{aligned} & \left| \mathbb{E} \left[z_{\tau_\beta}^\top \left(-A_{\theta_{\tau_\beta}} z_{\tau_\beta} - \frac{1}{\beta} (z_{\tau_\beta+1} - z_{\tau_\beta}) \right) \right] \right| \\ & \leq (R_1 + R_3 + P_1 + P_2 + P_3) \mathbb{E} [\|z_{\tau_\beta}\|^2] + (Q_1 + Q_2 + Q_3 + P_1 + P_2 + P_3) \\ & \quad + \frac{\alpha}{8\beta} L_\omega \mathbb{E} [\|\nabla J(\theta_{\tau_\beta})\|^2]. \end{aligned} \tag{119}$$

680 This completes the proof. \square

681 C.2 Proof under the Markovian Setting

682 In this section, we prove Theorem 1 under the Markovian setting.

683 From the L_J -smoothness of $J(\theta)$, it follows that

$$\begin{aligned} J(\theta_{t+1}) & \leq J(\theta_t) + \langle \nabla J(\theta_t), \theta_{t+1} - \theta_t \rangle + \frac{L_J}{2} \|\theta_{t+1} - \theta_t\|^2 \\ & = J(\theta_t) + \alpha \langle \nabla J(\theta_t), G_{t+1}(\theta_t, \omega_t) \rangle + \frac{L_J}{2} \alpha^2 \|G_{t+1}(\theta_t, \omega_t)\|^2 \\ & = J(\theta_t) - \alpha \left\langle \nabla J(\theta_t), -G_{t+1}(\theta_t, \omega_t) - \frac{\nabla J(\theta_t)}{2} + G_{t+1}(\theta_t, \omega(\theta_t)) - G_{t+1}(\theta_t, \omega(\theta_t)) \right\rangle \\ & \quad - \frac{\alpha}{2} \|\nabla J(\theta_t)\|^2 + \frac{L_J}{2} \alpha^2 \|G_{t+1}(\theta_t, \omega_t)\|^2 \\ & = J(\theta_t) - \alpha \langle \nabla J(\theta_t), -G_{t+1}(\theta_t, \omega_t) + G_{t+1}(\theta_t, \omega(\theta_t)) \rangle \\ & \quad + \alpha \left\langle \nabla J(\theta_t), \frac{\nabla J(\theta_t)}{2} + G_{t+1}(\theta_t, \omega(\theta_t)) \right\rangle - \frac{\alpha}{2} \|\nabla J(\theta_t)\|^2 + \frac{L_J}{2} \alpha^2 \|G_{t+1}(\theta_t, \omega_t)\|^2 \\ & \leq J(\theta_t) + \alpha L_g \|\nabla J(\theta_t)\| \|\omega(\theta_t) - \omega_t\| + \alpha \left\langle \nabla J(\theta_t), \frac{\nabla J(\theta_t)}{2} + G_{t+1}(\theta_t, \omega(\theta_t)) \right\rangle \\ & \quad - \frac{\alpha}{2} \|\nabla J(\theta_t)\|^2 + \frac{L_J}{2} \alpha^2 \|G_{t+1}(\theta_t, \omega_t)\|^2 \\ & \stackrel{(a)}{\leq} J(\theta_t) + \alpha L_g \|\nabla J(\theta_t)\| \|z_t\| + \alpha \left\langle \nabla J(\theta_t), \frac{\nabla J(\theta_t)}{2} + G_{t+1}(\theta_t, \omega(\theta_t)) \right\rangle \\ & \quad - \frac{\alpha}{2} \|\nabla J(\theta_t)\|^2 + \frac{L_J}{2} \alpha^2 C_g^2, \end{aligned} \tag{120}$$

684 where (a) is from the fact that $\|\theta_{t+1} - \theta_t\| \leq \alpha C_g$. Thus by re-arranging the terms, taking expectation
685 and summing up w.r.t. t from 0 to $T-1$, it follows that

$$\frac{\alpha}{2} \sum_{t=0}^{T-1} \mathbb{E} [\|\nabla J(\theta_t)\|^2]$$

$$\begin{aligned} &\leq -\mathbb{E}[J(\theta_T)] + J(\theta_0) + \alpha L_g \sqrt{\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]} \sqrt{\sum_{t=0}^{T-1} \mathbb{E}[\|z_t\|^2]} + \sum_{t=0}^{T-1} \alpha \mathbb{E}[\zeta_G(\theta_t, O_t)] \\ &\quad + L_J \alpha^2 T C_g^2, \end{aligned} \tag{121}$$

686 where $\zeta_G(\theta_t, O_t) = \left\langle \nabla J(\theta_t), \frac{\nabla J(\theta_t)}{2} + G_{t+1}(\theta_t, \omega(\theta_t)) \right\rangle$. We then bound ζ_G in the following
687 lemma.

688 **Lemma 6.** For any $t \geq \tau_\beta$,

$$\mathbb{E}[\zeta_G(\theta_t, O_t)] \leq 2C_g^2 \beta + 2\alpha \tau_\beta L_\zeta C_g. \tag{122}$$

689 *Proof.* We only need to consider the case $t = \tau_\beta$, the proof for general case of $t \geq \tau_\beta$ is similar, and
690 thus is omitted here. We first have that

$$\begin{aligned} \zeta_G(\theta_{\tau_\beta}, O_{\tau_\beta}) &= \left\langle \nabla J(\theta_{\tau_\beta}), \frac{\nabla J(\theta_{\tau_\beta})}{2} + G_{\tau_\beta+1}(\theta_{\tau_\beta}, \omega(\theta_{\tau_\beta})) \right\rangle \\ &= \left\langle \nabla J(\theta_0), \frac{\nabla J(\theta_0)}{2} + G_{\tau_\beta+1}(\theta_0, \omega(\theta_0)) \right\rangle \\ &\quad + \left\langle \nabla J(\theta_{\tau_\beta}), \frac{\nabla J(\theta_{\tau_\beta})}{2} + G_{\tau_\beta+1}(\theta_{\tau_\beta}, \omega(\theta_{\tau_\beta})) \right\rangle \\ &\quad - \left\langle \nabla J(\theta_0), \frac{\nabla J(\theta_0)}{2} + G_{\tau_\beta+1}(\theta_0, \omega(\theta_0)) \right\rangle \\ &\leq \left\langle \nabla J(\theta_0), \frac{\nabla J(\theta_0)}{2} + G_{\tau_\beta+1}(\theta_0, \omega(\theta_0)) \right\rangle + 2L_\zeta \|\theta_{\tau_\beta} - \theta_0\| \\ &\leq \left\langle \nabla J(\theta_0), \frac{\nabla J(\theta_0)}{2} + G_{\tau_\beta+1}(\theta_0, \omega(\theta_0)) \right\rangle + 2\alpha \tau_\beta L_\zeta C_g, \end{aligned} \tag{123}$$

691 where $L_\zeta = 2C_g(L'_g + \frac{3L_J}{2})$ is the Lipschitz constant of $\zeta_G(\theta, O_t)$.

692 Then it follows that

$$\begin{aligned} &\mathbb{E}[\zeta_G(\theta_{\tau_\beta}, O_{\tau_\beta})] \\ &= \mathbb{E}\left[\left\langle \nabla J(\theta_0), \frac{\nabla J(\theta_0)}{2} + G_{\tau_\beta+1}(\theta_0, \omega(\theta_0)) \right\rangle\right] + 2\alpha L_\zeta C_g \tau_\beta \\ &\leq 2C_g^2 \beta + 2\alpha \tau_\beta L_\zeta C_g, \end{aligned} \tag{124}$$

693 where the last step follows from the uniform ergodicity of the MDP (Assumption 4). \square

694 Plugging the bound in (121), it follows that

$$\begin{aligned} &\frac{\alpha}{2} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2] \\ &\leq J(\theta_0) - J^* + \alpha L_g \sqrt{\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]} \sqrt{\sum_{t=0}^{T-1} \mathbb{E}[\|z_t\|^2]} \\ &\quad + \alpha^2 C_g^2 L_J T + \alpha (T(2C_g^2 \beta + 2\alpha \tau_\beta L_\zeta C_g) + 4\tau_\beta C_g^2), \end{aligned} \tag{125}$$

695 and thus

$$\begin{aligned} &\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2] \\ &\leq \frac{2(J(\theta_0) - J^*)}{\alpha} + 2L_g \sqrt{\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]} \sqrt{\sum_{t=0}^{T-1} \mathbb{E}[\|z_t\|^2]} + 2\alpha C_g^2 L_J T \end{aligned}$$

$$+ 2 \left(T(2C_g^2\beta + 2\alpha\tau_\beta L_\zeta C_g) + 4\tau_\beta C_g^2 \right). \quad (126)$$

696 This further implies that

$$\begin{aligned} & \frac{\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]}{T} \\ & \leq \frac{2(J(\theta_0) - J^*)}{\alpha T} + 2L_g \sqrt{\frac{\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]}{T}} \sqrt{\frac{\sum_{t=0}^{T-1} \mathbb{E}[\|z_t\|^2]}{T}} + 2\alpha C_g^2 L_J \\ & \quad + 2 \left((2C_g^2\beta + 2\alpha\tau_\beta L_\zeta C_g) + 4C_g^2 \frac{\tau_\beta}{T} \right). \end{aligned} \quad (127)$$

697 We plug in the tracking error (94), and it follows that

$$\begin{aligned} & \frac{\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]}{T} \\ & \leq \frac{2(J(\theta_0) - J^*)}{\alpha T} + 2\alpha C_g^2 L_J + 2 \left((2C_g^2\beta + 2\alpha\tau_\beta L_\zeta C_g) + 4C_g^2 \frac{\tau_\beta}{T} \right) \\ & \quad + 2L_g \sqrt{\frac{\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]}{T}} \\ & \quad \cdot \sqrt{(2\|z_0\| + 2\beta\tau_\beta(b_{\max} + L_\omega C_g))^2 \left(\frac{1}{Tq} + \frac{\tau_\beta}{T} \right) + \frac{\alpha L_\omega \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]}{4q} + \frac{p}{q}} \\ & \leq \frac{2(J(\theta_0) - J^*)}{\alpha T} + 2\alpha C_g^2 L_J + 2 \left((2C_g^2\beta + 2\alpha\tau_\beta L_\zeta C_g) + 4C_g^2 \frac{\tau_\beta}{T} \right) \\ & \quad + 2L_g \sqrt{\frac{\alpha L_\omega \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]}{4q}} \\ & \quad + 2L_g \sqrt{\frac{\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]}{T}} \sqrt{(2\|z_0\| + 2\beta\tau_\beta(b_{\max} + L_\omega C_g))^2 \left(\frac{1}{Tq} + \frac{\tau_\beta}{T} \right) + \frac{p}{q}}. \end{aligned} \quad (128)$$

698 Note that $2L_g \sqrt{\frac{\alpha L_\omega}{4q}} = \mathcal{O}\left(\sqrt{\frac{\alpha}{\beta}}\right)$, hence we can choose α and β such that $2L_g \sqrt{\frac{\alpha L_\omega}{4q}} \leq \frac{1}{2}$. Hence
699 it follows that

$$\begin{aligned} & \frac{\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]}{T} \leq \frac{4(J(\theta_0) - J^*)}{\alpha T} + 4\alpha C_g^2 L_J + 4 \left((2C_g^2\beta + 2\alpha\tau_\beta L_\zeta C_g) + 4C_g^2 \frac{\tau_\beta}{T} \right) \\ & \quad + 4L_g \sqrt{\frac{\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]}{T}} \\ & \quad \cdot \sqrt{(2\|z_0\| + 2\beta\tau_\beta(b_{\max} + L_\omega C_g))^2 \left(\frac{1}{Tq} + \frac{\tau_\beta}{T} \right) + \frac{p}{q}} \\ & \triangleq U \sqrt{\frac{\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]}{T}} + V, \end{aligned} \quad (129)$$

700 where $U = 4L_g \sqrt{(2\|z_0\| + 2\beta\tau_\beta(b_{\max} + L_\omega C_g))^2 \left(\frac{1}{Tq} + \frac{\tau_\beta}{T} \right) + \frac{p}{q}} = \mathcal{O}\left(\sqrt{\beta\tau_\beta + \frac{1}{T\beta}}\right)$ and $V =$
701 $\frac{4(J(\theta_0) - J^*)}{\alpha T} + 4\alpha C_g^2 L_J + 4 \left((2C_g^2\beta + 2\alpha\tau_\beta L_\zeta C_g) + 4C_g^2 \frac{\tau_\beta}{T} \right) = \mathcal{O}\left(\frac{1}{T\alpha} + \alpha\tau_\beta + \beta\right)$. Thus it can
702 be shown that

$$\begin{aligned} & \frac{\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla J(\theta_t)\|^2]}{T} \leq \left(\frac{U + \sqrt{U^2 + 4V}}{2} \right)^2 \\ & \leq U^2 + 2V \\ & = 16L_g^2 \left((2\|z_0\| + 2\beta\tau_\beta(b_{\max} + L_\omega C_g))^2 \left(\frac{1}{Tq} + \frac{\tau_\beta}{T} \right) + \frac{p}{q} \right) \end{aligned}$$

$$\begin{aligned}
& + \frac{8(J(\theta_0) - J^*)}{\alpha T} + 8\alpha C_g^2 L_J + 8 \left((2C_g^2 \beta + 2\alpha \tau_\beta L_\zeta C_g) + 4C_g^2 \frac{\tau_\beta}{T} \right) \\
& = \mathcal{O} \left(\beta \tau_\beta + \frac{1}{T \beta} + \alpha \tau_\beta + \frac{1}{T \alpha} \right).
\end{aligned} \tag{130}$$

703 This completes the proof.

704 D Experiments

705 In this section, we provide some numerical experiments on two RL examples: the Garnet problem [1]
706 and the “spiral” counter example in [33].

707 D.1 Garnet Problem

708 The first experiment is on the Garnet problem [1], which can be characterized by $\mathcal{G}(|\mathcal{S}|, |\mathcal{A}|, b, N)$.
709 Here b is a branching parameter specifying how many next states are possible for each state-action
710 pair, and these b states are chosen uniformly at random. The transition probabilities are generated by
711 sampling uniformly and randomly between 0 and 1. The parameter N is the dimension of θ to be
712 updated. In our experiments, we generate a reward matrix uniformly and randomly between 0 and 1.
713 For every state s we randomly generate one feature function $k(s) \in [0, 1]$ using as the input. In both
714 experiments, we use a five-layer neural network with (1,2,2,3,1) neurons in each layer as the function
715 approximator. And for the activation function, we use the Sigmoid function, i.e., $f(x) = \frac{1}{1+e^{-x}}$. We
716 set all the weights and bias of the neurons as the parameter $\theta \in \mathbb{R}^{23}$.

717 We consider two sets of parameters: $\mathcal{G}(5, 2, 5, 23)$ and $\mathcal{G}(3, 2, 3, 23)$. We set the step-size $\alpha = 0.01$
718 and $\beta = 0.05$, and also the discount factor $\gamma = 0.95$. In Figures 1 and 2, we plot the squared gradient
719 norm v.s. the number of samples using 40 Garnet MDP trajectories, i.e., at each time t , we plot
720 $\|\nabla J(\theta_t)\|^2$. The upper and lower envelopes of the curves correspond to the 95 and 5 percentiles of
721 the 40 curves, respectively. We also plot the estimated variance of the stochastic update along the
722 iterations in Figures 1(b) and 2(b). Specifically, we first run the algorithm to get a sequence of θ_t and
723 ω_t . Then we generate 500 different trajectories $O^i = (O_1^i, O_2^i, \dots, O_t^i, \dots)$ where $i = 1, \dots, 500$, and
724 use them to estimate the variance $\|G_{t+1}^i(\theta_t, \omega_t) - \nabla J(\theta_t)\|^2$ and plot $\frac{\sum_{i=1}^{500} \|G_{t+1}^i(\theta_t, \omega_t) - \nabla J(\theta_t)\|^2}{500}$ at
725 each time t .

726 It can be seen from the figures that both gradient norm $\|\nabla J(\theta_t)\|$ and the estimated variance converge
727 to zero.

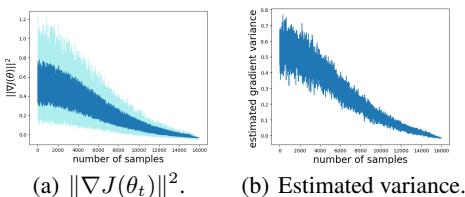


Figure 1: Garnet problem 1: $\mathcal{G}(5, 2, 5, 23)$.

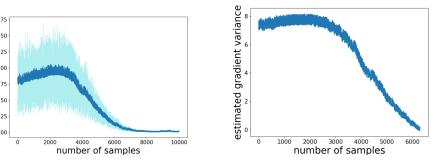


Figure 2: Garnet problem 2: $\mathcal{G}(3, 2, 3, 23)$.

728 D.2 Spiral Counter Example

729 In our second experiment, we consider the spiral counter example proposed in [33], which is often
730 used to show the TD algorithm may diverge with nonlinear function approximation. The problem
731 setting is given in Figure 3. There are three states and each state can transit to the next one with
732 probability $\frac{1}{2}$ or stay at the current state with probability $\frac{1}{2}$. The reward is always zero with the
733 discount factor $\gamma = 0.9$. Similar to [5], we consider the value function approximation:

$$V_\theta(s) = (a(s) \cos(k\theta) + b(s) \sin(k\theta)) e^{\epsilon\theta}, \tag{131}$$

734 where in Figure 4, $a = [0.94, -0.43, 0.18]$ and $b = [0.21, -0.52, 0.76]$; and in Figure 5, $a =$
735 $[0.21, -0.33, 0.29]$ and $b = [0.68, 0.41, 0.82]$. We let $k = 0.866$ and $\epsilon = 0.1$. The step-size are

736 chosen as $\alpha = 0.01$ and $\beta = 0.05$. In Figures 4(a) and 5(a), we plot the squared gradient norm v.s.
 737 the number of samples using 40 MDP trajectories. The upper and lower envelopes of the curves
 738 correspond to the 95 and 5 percentiles of the 40 curves. Similarly, we also plot the estimated variance
 739 $\|G_{t+1}(\theta_t, \omega_t) - \nabla J(\theta_t)\|^2$ of the stochastic update along the iterations using 50 samples at each time
 740 step. More specifically, we first run the algorithm to get a sequence of θ_t and ω_t . Then we generate
 741 50 different trajectories $O^i = (O_1^i, O_2^i, \dots, O_t^i, \dots)$ where $i = 1, \dots, 50$, and use them to estimate the
 742 variance $\|G_{t+1}^i(\theta_t, \omega_t) - \nabla J(\theta_t)\|^2$ and plot $\frac{\sum_{i=1}^{50} \|G_{t+1}^i(\theta_t, \omega_t) - \nabla J(\theta_t)\|^2}{50}$ at each time t .
 743 It can be seen that in both experiments, the gradient norm $\|\nabla J(\theta_t)\|$ converges to 0, i.e., the algorithm
 744 converges to a stationary point. The estimated variance also decreases to zero.

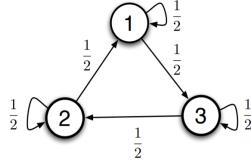
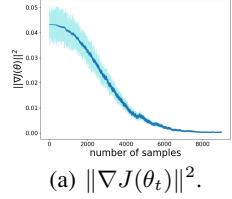
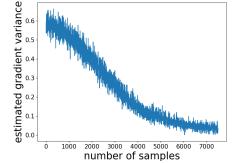


Figure 3: Spiral counter example.

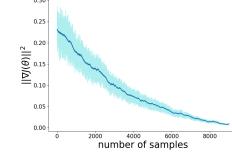


(a) $\|\nabla J(\theta_t)\|^2$.

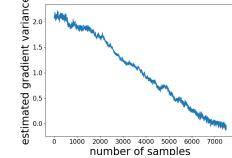


(b) Estimated variance.

Figure 4: Spiral counter example 1:
 $a = [0.94, -0.43, 0.18], b = [0.21, -0.52, 0.76]$.



(a) $\|\nabla J(\theta_t)\|^2$.



(b) Estimated variance.

Figure 5: Spiral counter example 2:
 $a = [0.21, -0.33, 0.29], b = [0.68, 0.41, 0.82]$.

745