



## Supplementary materials for

Zhenyi ZHANG, Jie HUANG, Congjie PAN, 2024. Multi-agent reinforcement learning behavioral control for nonlinear second-order systems. *Front Inform Technol Electron Eng*, 25(6):869-886.

<https://doi.org/10.1631/FITEE.2300394>

### 1 Proof of Theorem 1

Consider the Lyapunov function candidate of Theorem 1 as

$$V_u = \frac{1}{2} \sum_{i=1}^N \mathbf{e}_{e,i}^T \mathbf{r}_{e,i} \mathbf{e}_{e,i} + \frac{1}{2} \sum_{i=1}^N \boldsymbol{\tau}_i^T \boldsymbol{\tau}_i + \frac{1}{2} \sum_{i=1}^N \text{tr} \left\{ \tilde{\boldsymbol{\omega}}_{f,i}^T \mathbf{r}_{f,i}^{-1} \tilde{\boldsymbol{\omega}}_{f,i} \right\} + \frac{1}{2} \sum_{i=1}^N \text{tr} \left\{ \tilde{\boldsymbol{\omega}}_{c,i}^T \tilde{\boldsymbol{\omega}}_{c,i} \right\} + \frac{1}{2} \sum_{i=1}^N \text{tr} \left\{ \tilde{\boldsymbol{\omega}}_{a,i}^T \tilde{\boldsymbol{\omega}}_{a,i} \right\}, \quad (\text{S1})$$

where  $\mathbf{r}_{e,i} = \begin{bmatrix} \xi_{e,i} \mathbf{I}_n & \mathbf{I}_n \\ \mathbf{I}_n & \mathbf{I}_n \end{bmatrix}$  is a positive-definite matrix if condition (42) of Theorem 1 is held;  $\xi_{e,i}$  is a designed parameter of  $\mathbf{r}_{e,i}$ ;  $\text{tr}\{\cdot\}$  is the trace of a matrix; and  $\tilde{\boldsymbol{\omega}}_{f,i} = \hat{\boldsymbol{\omega}}_{f,i} - \boldsymbol{\omega}_{f,i}^*$ ,  $\tilde{\boldsymbol{\omega}}_{c,i} = \hat{\boldsymbol{\omega}}_{c,i} - \boldsymbol{\omega}_{c,i}^*$ , and  $\tilde{\boldsymbol{\omega}}_{a,i} = \hat{\boldsymbol{\omega}}_{a,i} - \boldsymbol{\omega}_{a,i}^*$  are the weight errors of the identifier, critic, and actor, respectively.

The following result is yielded by differentiating Eq. (S1) along Eqs. (21), (30), (37), and (39)–(41):

$$\begin{aligned} \dot{V}_u = & \sum_{i=1}^N \left[ \xi_{e,i} \mathbf{e}_{p,i}^T \mathbf{e}_{v,i} + \mathbf{e}_{v,i}^T \mathbf{e}_{v,i} + (\mathbf{e}_{p,i}^T + \mathbf{e}_{v,i}^T) (-\xi_{p,i} \mathbf{e}_{p,i} - \xi_{v,i} \mathbf{e}_{v,i} - \hat{\boldsymbol{\omega}}_{f,i}^T \boldsymbol{\Psi}_{f,i}(\boldsymbol{\chi}_i) - \boldsymbol{\tau}_i - \frac{1}{2\beta_V} \hat{\boldsymbol{\omega}}_{a,i}^T \boldsymbol{\Psi}_{V,i}(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \right. \\ & \left. + \mathbf{f}_i(\boldsymbol{\chi}_i) - \dot{\mathbf{v}}_{i,r}) \right] + \sum_{i=1}^N \boldsymbol{\tau}_i^T (-\xi_{\tau,i} \boldsymbol{\tau}_i + \mathbf{u}_{\Delta,i}) + \sum_{i=1}^N \text{tr} \left\{ \tilde{\boldsymbol{\omega}}_{f,i}^T [\boldsymbol{\Psi}_{f,i}(\boldsymbol{\chi}_i) (\mathbf{e}_{p,i}^T + \mathbf{e}_{v,i}^T) - \boldsymbol{\nu}_{f,i} \hat{\boldsymbol{\omega}}_{f,i}] \right\} \\ & - \sum_{i=1}^N \text{tr} \left\{ \gamma_{c,i} \tilde{\boldsymbol{\omega}}_{c,i}^T \boldsymbol{\Psi}_{V,i}(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \boldsymbol{\Psi}_{V,i}^T(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \hat{\boldsymbol{\omega}}_{c,i} \right\} \\ & - \sum_{i=1}^N \text{tr} \left\{ \tilde{\boldsymbol{\omega}}_{a,i}^T \boldsymbol{\Psi}_{V,i}(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \boldsymbol{\Psi}_{V,i}^T(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) [\gamma_{a,i} (\hat{\boldsymbol{\omega}}_{a,i} - \hat{\boldsymbol{\omega}}_{c,i}) + \gamma_{c,i} \hat{\boldsymbol{\omega}}_{c,i}] \right\}. \end{aligned} \quad (\text{S2})$$

According to Property 1 and Eq. (32), Eq. (S2) can be further simplified as

$$\begin{aligned} \dot{V}_u = & - \sum_{i=1}^N [\xi_{p,i} \|\mathbf{e}_p\|^2 + (\xi_{v,i} - 1) \|\mathbf{e}_v\|^2 + \xi_{\tau,i} \|\boldsymbol{\tau}_i\|^2] - \sum_{i=1}^N (\xi_{p,i} + \xi_{v,i} - \xi_{e,i}) \mathbf{e}_{p,i}^T \mathbf{e}_{v,i} - \sum_{i=1}^N (\mathbf{e}_{p,i}^T + \mathbf{e}_{v,i}^T) \\ & \cdot \left[ \boldsymbol{\tau}_i + \frac{1}{2\beta_V} \hat{\boldsymbol{\omega}}_{a,i}^T \boldsymbol{\Psi}_{V,i}(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) + \dot{\mathbf{v}}_{i,r} - \boldsymbol{\sigma}_{f,i}(\boldsymbol{\chi}_i) \right] + \sum_{i=1}^N \boldsymbol{\tau}_i^T \mathbf{u}_{\Delta,i} - \sum_{i=1}^N \text{tr} \left\{ \tilde{\boldsymbol{\omega}}_{f,i}^T \boldsymbol{\nu}_{f,i} \hat{\boldsymbol{\omega}}_{f,i} \right\} \\ & - \sum_{i=1}^N \gamma_{c,i} \text{tr} \left\{ \tilde{\boldsymbol{\omega}}_{c,i}^T \boldsymbol{\Psi}_{V,i}(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \boldsymbol{\Psi}_{V,i}^T(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \hat{\boldsymbol{\omega}}_{c,i} \right\} - \sum_{i=1}^N \gamma_{a,i} \text{tr} \left\{ \tilde{\boldsymbol{\omega}}_{a,i}^T \boldsymbol{\Psi}_{V,i}(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \boldsymbol{\Psi}_{V,i}^T(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \hat{\boldsymbol{\omega}}_{a,i} \right\} \\ & + \sum_{i=1}^N (\gamma_{a,i} - \gamma_{c,i}) \text{tr} \left\{ \tilde{\boldsymbol{\omega}}_{a,i}^T \boldsymbol{\Psi}_{V,i}(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \boldsymbol{\Psi}_{V,i}^T(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \hat{\boldsymbol{\omega}}_{c,i} \right\}. \end{aligned} \quad (\text{S3})$$

The following results can be obtained by using Properties 1 and 2:

$$-(\mathbf{e}_{p,i}^T + \mathbf{e}_{v,i}^T) \boldsymbol{\tau}_i \leq \frac{1}{2} \|\mathbf{e}_{p,i}\|^2 + \frac{1}{2} \|\mathbf{e}_{v,i}\|^2 + \|\boldsymbol{\tau}_i\|^2, \quad (\text{S4})$$

$$-\frac{1}{2\beta_V}(\mathbf{e}_{p,i}^T + \mathbf{e}_{v,i}^T)\hat{\omega}_{a,i}^T \Psi_{V,i}(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \leq \frac{1}{2} \|\mathbf{e}_{p,i}\|^2 + \frac{1}{2} \|\mathbf{e}_{v,i}\|^2 + \frac{1}{4\beta_V^2} \text{tr} \{ \hat{\omega}_{a,i}^T \Psi_{V,i}(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \Psi_{V,i}^T(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \hat{\omega}_{a,i} \}, \quad (\text{S5})$$

$$-(\mathbf{e}_{p,i}^T + \mathbf{e}_{v,i}^T)\dot{\mathbf{v}}_{i,r} \leq \frac{1}{2} \|\mathbf{e}_{p,i}\|^2 + \frac{1}{2} \|\mathbf{e}_{v,i}\|^2 + \|\dot{\mathbf{v}}_{i,r}\|^2, \quad (\text{S6})$$

$$(\mathbf{e}_{p,i}^T + \mathbf{e}_{v,i}^T)\boldsymbol{\sigma}_{f,i}(\boldsymbol{\chi}_i) \leq \frac{1}{2} \|\mathbf{e}_{p,i}\|^2 + \frac{1}{2} \|\mathbf{e}_{v,i}\|^2 + \|\boldsymbol{\sigma}_{f,i}(\boldsymbol{\chi}_i)\|^2 \leq \frac{1}{2} \|\mathbf{e}_{p,i}\|^2 + \frac{1}{2} \|\mathbf{e}_{v,i}\|^2 + \varepsilon_f^2, \quad (\text{S7})$$

$$\boldsymbol{\tau}_i^T \mathbf{u}_{\Delta,i} \leq \frac{1}{2} \|\boldsymbol{\tau}_i\|^2 + \frac{1}{2} \|\mathbf{u}_{\Delta,i}\|^2 \leq \frac{1}{2} \|\boldsymbol{\tau}_i\|^2 + \frac{1}{2} \varepsilon_{\Delta}^2. \quad (\text{S8})$$

Substituting inequalities (S4)–(S8) into Eq. (S3), we obtain

$$\begin{aligned} \dot{V}_u \leq & - \sum_{i=1}^N \left[ (\xi_{p,i} - 2) \|\mathbf{e}_{p,i}\|^2 + (\xi_{v,i} - 3) \|\mathbf{e}_{v,i}\|^2 + \left( \xi_{\tau,i} - \frac{3}{2} \right) \|\boldsymbol{\tau}_i\|^2 \right] - \sum_{i=1}^N (\xi_{p,i} + \xi_{v,i} - \xi_{e,i}) \mathbf{e}_{p,i}^T \mathbf{e}_{v,i} \\ & - \sum_{i=1}^N \text{tr} \{ \tilde{\omega}_{f,i}^T \nu_{f,i} \hat{\omega}_{f,i} \} - \sum_{i=1}^N \gamma_{c,i} \text{tr} \{ \tilde{\omega}_{c,i}^T \Psi_{V,i}(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \Psi_{V,i}^T(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \hat{\omega}_{c,i} \} - \sum_{i=1}^N \gamma_{a,i} \text{tr} \{ \tilde{\omega}_{a,i}^T \Psi_{V,i}(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \\ & \cdot \Psi_{V,i}^T(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \hat{\omega}_{a,i} \} + \sum_{i=1}^N (\gamma_{a,i} - \gamma_{c,i}) \text{tr} \{ \tilde{\omega}_{a,i}^T \Psi_{V,i}(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \Psi_{V,i}^T(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \hat{\omega}_{c,i} \} \\ & + \frac{1}{4\beta_V^2} \sum_{i=1}^N \text{tr} \{ \hat{\omega}_{a,i}^T \Psi_{V,i}(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \Psi_{V,i}^T(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \hat{\omega}_{a,i} \} + \sum_{i=1}^N \|\dot{\mathbf{v}}_{i,r}\|^2 + N\varepsilon_f^2 + \frac{N}{2}\varepsilon_{\Delta}^2. \end{aligned} \quad (\text{S9})$$

The results obtained pursuant to the use of Properties 1–3, and from the fact that  $\tilde{\omega}_{f,i} = \hat{\omega}_{f,i} - \omega_{f,i}^*$ ,  $\tilde{\omega}_{c,i} = \hat{\omega}_{c,i} - \omega_{c,i}^*$ , and  $\tilde{\omega}_{a,i} = \hat{\omega}_{a,i} - \omega_{a,i}^*$ , are the following:

$$\text{tr} \{ \tilde{\omega}_{f,i}^T \hat{\omega}_{f,i} \} = \frac{1}{2} \text{tr} \{ \tilde{\omega}_{f,i}^T \tilde{\omega}_{f,i} \} + \frac{1}{2} \text{tr} \{ \hat{\omega}_{f,i}^T \hat{\omega}_{f,i} \} - \frac{1}{2} \text{tr} \{ (\omega_{f,i}^*)^T \omega_{f,i}^* \}, \quad (\text{S10})$$

$$\begin{aligned} \text{tr} \{ \tilde{\omega}_{c,i}^T \Psi_{V,i}(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \Psi_{V,i}^T(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \hat{\omega}_{c,i} \} &= \frac{1}{2} \text{tr} \{ \tilde{\omega}_{c,i}^T \Psi_{V,i}(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \Psi_{V,i}^T(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \tilde{\omega}_{c,i} \} \\ &+ \frac{1}{2} \text{tr} \{ \hat{\omega}_{c,i}^T \Psi_{V,i}(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \Psi_{V,i}^T(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \hat{\omega}_{c,i} \} - \frac{1}{2} \text{tr} \{ (\omega_{c,i}^*)^T \Psi_{V,i}(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \Psi_{V,i}^T(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \omega_{c,i}^* \}, \end{aligned} \quad (\text{S11})$$

$$\begin{aligned} \text{tr} \{ \tilde{\omega}_{a,i}^T \Psi_{V,i}(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \Psi_{V,i}^T(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \hat{\omega}_{a,i} \} &= \frac{1}{2} \text{tr} \{ \tilde{\omega}_{a,i}^T \Psi_{V,i}(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \Psi_{V,i}^T(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \tilde{\omega}_{a,i} \} \\ &+ \frac{1}{2} \text{tr} \{ \hat{\omega}_{a,i}^T \Psi_{V,i}(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \Psi_{V,i}^T(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \hat{\omega}_{a,i} \} - \frac{1}{2} \text{tr} \{ (\omega_{a,i}^*)^T \Psi_{V,i}(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \Psi_{V,i}^T(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \omega_{a,i}^* \}, \end{aligned} \quad (\text{S12})$$

$$\begin{aligned} & (\gamma_{a,i} - \gamma_{c,i}) \text{tr} \{ \tilde{\omega}_{a,i}^T \Psi_{V,i}(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \Psi_{V,i}^T(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \hat{\omega}_{c,i} \} \\ \leq & \frac{\gamma_{a,i} - \gamma_{c,i}}{2} \text{tr} \{ \tilde{\omega}_{a,i}^T \Psi_{V,i}(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \Psi_{V,i}^T(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \tilde{\omega}_{a,i} \} + \frac{\gamma_{a,i} - \gamma_{c,i}}{2} \text{tr} \{ \hat{\omega}_{c,i}^T \Psi_{V,i}(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \Psi_{V,i}^T(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \hat{\omega}_{c,i} \}. \end{aligned} \quad (\text{S13})$$

Substituting Eqs. (S10)–(S12) into inequality (S9), we obtain

$$\begin{aligned} \dot{V}_u \leq & - \sum_{i=1}^N \frac{1}{2} \mathbf{e}_i^T \boldsymbol{\Gamma}_{\xi,i} \mathbf{e}_i - \sum_{i=1}^N \left( \xi_{\tau,i} - \frac{3}{2} \right) \|\boldsymbol{\tau}_i\|^2 - \sum_{i=1}^N \frac{\nu_{f,i}}{2\lambda_f} \text{tr} \{ \tilde{\omega}_{f,i}^T \boldsymbol{\Gamma}_{f,i}^{-1} \tilde{\omega}_{f,i} \} \\ & - \sum_{i=1}^N \frac{\gamma_{c,i}}{2} \text{tr} \{ \tilde{\omega}_{c,i}^T \Psi_{V,i}(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \Psi_{V,i}^T(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \tilde{\omega}_{c,i} \} - \sum_{i=1}^N \frac{\gamma_{c,i}}{2} \text{tr} \{ \tilde{\omega}_{a,i}^T \Psi_{V,i}(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \Psi_{V,i}^T(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \tilde{\omega}_{a,i} \} \\ & - \sum_{i=1}^N \left( \gamma_{c,i} - \frac{\gamma_{a,i}}{2} \right) \text{tr} \{ \hat{\omega}_{c,i}^T \Psi_{V,i}(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \Psi_{V,i}^T(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \hat{\omega}_{c,i} \} \\ & - \sum_{i=1}^N \left( \frac{\gamma_{a,i}}{2} - \frac{1}{4\beta_V^2} \right) \text{tr} \{ \hat{\omega}_{a,i}^T \Psi_{V,i}(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \Psi_{V,i}^T(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \hat{\omega}_{a,i} \} + \Xi, \end{aligned} \quad (\text{S14})$$

where

$$\mathbf{r}_{\xi,i} = \begin{bmatrix} 2(\xi_{p,i} - 2)\mathbf{I}_n & (\xi_{p,i} + \xi_{v,i} - \xi_{e,i})\mathbf{I}_n \\ (\xi_{p,i} + \xi_{v,i} - \xi_{e,i})\mathbf{I}_n & 2(\xi_{v,i} - 3)\mathbf{I}_n \end{bmatrix},$$

$$\Xi = \sum_{i=1}^N \frac{\gamma_{a,i} + \gamma_{c,i}}{2} \text{tr} \left\{ (\boldsymbol{\omega}_{v,i}^*)^T \boldsymbol{\Psi}_{V,i}(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \boldsymbol{\Psi}_{V,i}^T(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \boldsymbol{\omega}_{v,i}^* \right\}$$

$$+ \frac{1}{2} \sum_{i=1}^N \nu_{f,i} \text{tr} \left\{ (\boldsymbol{\omega}_{f,i}^*)^T \boldsymbol{\omega}_{f,i}^* \right\} + \sum_{i=1}^N \|\dot{\mathbf{v}}_{i,r}\|^2 + N\varepsilon_f^2 + \frac{N}{2}\varepsilon_\Delta^2,$$

since all terms are bounded,  $\Xi$  is bounded as  $|\Xi| \leq \vartheta_2$ , and  $\bar{\lambda}_f$  is the maximal eigenvalue of  $\mathbf{r}_{f,i}^{-1}$ .

Based on condition (42), inequality (S14) can be rewritten in the compact form as

$$\dot{V}_u \leq -\frac{1}{2} \sum_{i=1}^N \frac{\underline{\lambda}_\xi}{\bar{\lambda}_e} \mathbf{e}_i^T \mathbf{r}_{e,i} \mathbf{e}_i - \frac{1}{2} \sum_{i=1}^N \left[ 2 \left( \xi_{\tau,i} - \frac{3}{2} \right) \|\boldsymbol{\tau}_i\|^2 - \frac{1}{2} \sum_{i=1}^N \frac{\nu_{f,i}}{\bar{\lambda}_f} \text{tr} \left\{ \tilde{\boldsymbol{\omega}}_{f,i}^T \mathbf{r}_{f,i}^{-1} \tilde{\boldsymbol{\omega}}_{f,i} \right\} \right. \\ \left. - \frac{1}{2} \sum_{i=1}^N \underline{\lambda}_\psi \text{tr} \left\{ \tilde{\boldsymbol{\omega}}_{c,i}^T \tilde{\boldsymbol{\omega}}_{c,i} \right\} - \frac{1}{2} \sum_{i=1}^N \underline{\lambda}_\psi \text{tr} \left\{ \tilde{\boldsymbol{\omega}}_{a,i}^T \tilde{\boldsymbol{\omega}}_{a,i} \right\} \right] + \vartheta_2, \quad (\text{S15})$$

where  $\underline{\lambda}_\xi$  is the minimal eigenvalue of  $\mathbf{r}_{\xi,i}$ ,  $\bar{\lambda}_e$  is the maximal eigenvalue of  $\mathbf{r}_{e,i}$ , and  $\underline{\lambda}_\psi$  is the minimal eigenvalue of  $\boldsymbol{\Psi}_{V,i}(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \boldsymbol{\Psi}_{V,i}^T(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i)$ .

Let  $\vartheta_1 = \min \left\{ \frac{\underline{\lambda}_\xi}{\bar{\lambda}_e}, 2 \left( \xi_{\tau,i} - \frac{3}{2} \right), \frac{\nu_{f,i}}{\bar{\lambda}_f}, \underline{\lambda}_\psi \right\}$ . Then we have

$$\dot{V}_u \leq -\vartheta_1 V_u + \vartheta_2. \quad (\text{S16})$$

The final result is obtained by using Lemma 1 as

$$V_u \leq e^{-\vartheta_1 t} V_u(0) + \frac{\vartheta_2}{\vartheta_1} (1 - e^{-\vartheta_1 t}). \quad (\text{S17})$$

Therefore, Theorem 1 is proven completely.

## 2 Proof of mission stability

Let us consider the Lyapunov function candidate of mission stability as

$$V_\rho = \frac{1}{2} \tilde{\boldsymbol{\rho}}_{\text{OA}}^T \kappa_{\text{OA}} \tilde{\boldsymbol{\rho}}_{\text{OA}} + \frac{1}{2} \tilde{\boldsymbol{\rho}}_{\text{FM}}^T \kappa_{\text{FM}} \tilde{\boldsymbol{\rho}}_{\text{FM}} + \frac{1}{2} \tilde{\boldsymbol{\rho}}_{\text{FR}}^T \kappa_{\text{FR}} \tilde{\boldsymbol{\rho}}_{\text{FR}}, \quad (\text{S18})$$

where  $\kappa_{\text{OA}}$ ,  $\kappa_{\text{FM}}$ , and  $\kappa_{\text{FR}}$  are the designed positive constants.

We assume that the behavior priority satisfies OA>FM>FR; then the differential of Eq. (S18) is

$$\dot{V}_\rho = -\frac{1}{2} \tilde{\boldsymbol{\zeta}}^T \boldsymbol{\Gamma}_\rho \tilde{\boldsymbol{\zeta}}, \quad (\text{S19})$$

where

$$\boldsymbol{\Gamma}_\rho = \begin{bmatrix} \boldsymbol{\Gamma}_{\rho,11} & \boldsymbol{\Gamma}_{\rho,12} & \boldsymbol{\Gamma}_{\rho,13} \\ \boldsymbol{\Gamma}_{\rho,21} & \boldsymbol{\Gamma}_{\rho,22} & \boldsymbol{\Gamma}_{\rho,23} \\ \boldsymbol{\Gamma}_{\rho,31} & \boldsymbol{\Gamma}_{\rho,32} & \boldsymbol{\Gamma}_{\rho,33} \end{bmatrix}$$

with  $\boldsymbol{\Gamma}_{\rho,11} = \kappa_{\text{OA}} \boldsymbol{\Lambda}_{\text{OA}}$ ,  $\boldsymbol{\Gamma}_{\rho,21} = \boldsymbol{\Gamma}_{\rho,12}^T = \frac{1}{2} \kappa_{\text{FM}} \mathbf{J}_{\text{FM}} \mathbf{J}_{\text{OA}}^\dagger \boldsymbol{\Lambda}_{\text{OA}}$ ,  $\boldsymbol{\Gamma}_{\rho,22} = \kappa_{\text{FM}} \mathbf{J}_{\text{FM}} (\mathbf{I}_n - \mathbf{J}_{\text{OA}}^\dagger \mathbf{J}_{\text{OA}}) \mathbf{J}_{\text{FM}}^\dagger \boldsymbol{\Lambda}_{\text{FM}}$ ,  $\boldsymbol{\Gamma}_{\rho,31} = \boldsymbol{\Gamma}_{\rho,13}^T = \frac{1}{2} \kappa_{\text{FR}} \mathbf{J}_{\text{FR}} \mathbf{J}_{\text{OA}}^\dagger \boldsymbol{\Lambda}_{\text{OA}}$ ,  $\boldsymbol{\Gamma}_{\rho,23} = \kappa_{\text{FR}} \mathbf{J}_{\text{FR}} (\mathbf{I}_n - \mathbf{J}_{\text{OA}}^\dagger \mathbf{J}_{\text{OA}}) (\mathbf{I}_n - \mathbf{J}_{\text{FM}}^\dagger \mathbf{J}_{\text{FM}}) \mathbf{J}_{\text{FR}}^\dagger \boldsymbol{\Lambda}_{\text{FR}}$ ,  $\boldsymbol{\Gamma}_{\rho,32} =$

$\kappa_{\text{FM}} \mathbf{J}_{\text{FR}} (\mathbf{I}_n - \mathbf{J}_{\text{OA}}^\dagger \mathbf{J}_{\text{OA}}) \mathbf{J}_{\text{FM}}^\dagger \mathbf{A}_{\text{FM}}$ ,  $\Gamma_{\rho,33} = \kappa_{\text{FR}} \mathbf{J}_{\text{FR}} (\mathbf{I}_n - \mathbf{J}_{\text{OA}}^\dagger \mathbf{J}_{\text{OA}}) (\mathbf{I}_n - \mathbf{J}_{\text{FM}}^\dagger \mathbf{J}_{\text{FM}}) \mathbf{J}_{\text{FR}}^\dagger \mathbf{A}_{\text{FR}}$ ,  $\tilde{\zeta} = [\tilde{\rho}_{\text{OA}}^\text{T}, \tilde{\rho}_{\text{FM}}^\text{T}, \tilde{\rho}_{\text{FR}}^\text{T}]^\text{T}$ . Eq. (S19) satisfies the following inequality:

$$\begin{aligned} \dot{V}_\rho \leq & -\underline{\Gamma}_{\rho,11} \|\tilde{\rho}_{\text{OA}}\|^2 - \underline{\Gamma}_{\rho,22} \|\tilde{\rho}_{\text{FM}}\|^2 - \underline{\Gamma}_{\rho,33} \|\tilde{\rho}_{\text{FR}}\|^2 + 2\overline{\Gamma}_{\rho,21} \|\tilde{\rho}_{\text{OA}}\| \|\tilde{\rho}_{\text{FM}}\| \\ & + 2\overline{\Gamma}_{\rho,31} \|\tilde{\rho}_{\text{OA}}\| \|\tilde{\rho}_{\text{FR}}\| + 2\overline{\Gamma}_{\rho,32} \|\tilde{\rho}_{\text{FM}}\| \|\tilde{\rho}_{\text{FR}}\| = -\hat{\zeta}^\text{T} \hat{\Gamma}_\rho \hat{\zeta}, \end{aligned} \quad (\text{S20})$$

where

$$\hat{\Gamma}_\rho = \begin{bmatrix} \kappa_{\text{OA}} \underline{\Delta}_{\text{OA}} & -\frac{\kappa_{\text{FR}} \overline{\Delta}_{\text{OA}}}{2} & -\frac{\kappa_{\text{FR}} \overline{\Delta}_{\text{OA}}}{2} \\ -\frac{\kappa_{\text{FM}} \overline{\Delta}_{\text{OA}}}{2} & \kappa_{\text{FM}} \underline{\Delta}_{\text{FM}} & -\kappa_{\text{FM}} \overline{\Delta}_{\text{FM,FR}} \\ -\frac{\kappa_{\text{FR}} \overline{\Delta}_{\text{OA}}}{2} & -\kappa_{\text{FM}} \overline{\Delta}_{\text{FM,FR}} & \kappa_{\text{FR}} \underline{\Delta}_{\text{FR}} \end{bmatrix}.$$

$\underline{\Gamma}_{\rho,11} = \kappa_{\text{OA}} \underline{\Delta}_{\text{OA}}$ ,  $\underline{\Gamma}_{\rho,22} = \kappa_{\text{FM}} \underline{\Delta}_{\text{FM}}$ , and  $\underline{\Gamma}_{\rho,33} = \kappa_{\text{FR}} \underline{\Delta}_{\text{FR}}$  are the lower bounds on the induced norms of  $\Gamma_{\rho,11}$ ,  $\Gamma_{\rho,22}$ , and  $\Gamma_{\rho,33}$ , respectively;  $\underline{\Delta}_{\text{OA}} = \min_i \{\Delta_{\text{OA},i}\}$ ,  $\underline{\Delta}_{\text{FM}} = \min_i \{\Delta_{\text{FM},i}\}$ ,  $\underline{\Delta}_{\text{FR}} = \min_i \{\Delta_{\text{FR},i}\}$ ;  $\overline{\Gamma}_{\rho,21} = \kappa_{\text{FM}} \overline{\Delta}_{\text{OA}}$ ,  $\overline{\Gamma}_{\rho,31} = \kappa_{\text{FR}} \overline{\Delta}_{\text{OA}}$ , and  $\overline{\Gamma}_{\rho,32} = 2\kappa_{\text{FM}} \overline{\Delta}_{\text{FM,FR}}$  are the upper bounds on the induced norms of  $\Gamma_{\rho,21}$ ,  $\Gamma_{\rho,31}$ , and  $\Gamma_{\rho,32}$ , respectively;  $\overline{\Delta}_{\text{OA}} = \max_i \{\Delta_{\text{OA},i}\}$ ,  $\overline{\Delta}_{\text{FM,FR}} = \max_i \{\Delta_{\text{FM},i}, \Delta_{\text{FR},i}\}$ ,  $\hat{\zeta} = [\|\tilde{\rho}_{\text{OA}}\|, \|\tilde{\rho}_{\text{FM}}\|, \|\tilde{\rho}_{\text{FR}}\|]^\text{T}$ ,  $\|\mathbf{I}_n - \mathbf{J}_{\text{OA}}^\dagger \mathbf{J}_{\text{OA}}\| \leq 1$ ,  $\|\mathbf{I}_n - \mathbf{J}_{\text{FM}}^\dagger \mathbf{J}_{\text{FM}}\| \leq 1$ . When

$$\kappa_{\text{OA}} > \max \left\{ \frac{\kappa_{\text{FM}} \overline{\Delta}_{\text{OA}}^2}{4 \underline{\Delta}_{\text{OA}} \underline{\Delta}_{\text{FM}}}, \frac{\kappa_{\text{FR}} \overline{\Delta}_{\text{OA}} (2\kappa_{\text{FM}} \overline{\Delta}_{\text{FM,FR}} + \kappa_{\text{FR}} \underline{\Delta}_{\text{FM}} + \kappa_{\text{FM}} \underline{\Delta}_{\text{FR}})}{4 (\kappa_{\text{FR}} \underline{\Delta}_{\text{FM}} \underline{\Delta}_{\text{FR}} - \kappa_{\text{FM}} \overline{\Delta}_{\text{FM,FR}}^2)} \right\},$$

we may infer that  $\hat{\Gamma}_\rho$  is a positive-definite symmetric matrix,

$$\dot{V}_\rho \leq -\underline{\lambda}_{\hat{\Gamma}} \hat{\zeta}^\text{T} \hat{\zeta} \leq 0, \quad (\text{S21})$$

where  $\underline{\lambda}_{\hat{\Gamma}}$  is the minimum eigenvalue of  $\hat{\Gamma}_\rho$ . Other behavior priority cases are proved by similar ways. Since Assumption 3 is held, control inputs are not always saturated, and mission errors are converged eventually.

We consider a special case in which the FM and FR behaviors are conflicting, and the OA behavior is conflict-free. We also assume that the priority is FM>FR. Then, the differential of Eq. (S18) is expressed as

$$\dot{V}_\rho = -\kappa_{\text{OA}} \tilde{\rho}_{\text{OA}}^\text{T} \mathbf{J}_{\text{OA}} \mathbf{J}_{\text{OA}}^\dagger \mathbf{A}_{\text{OA}} \tilde{\rho}_{\text{OA}} - \tilde{\zeta}^\text{T} \Gamma_\rho \tilde{\zeta}, \quad (\text{S22})$$

where  $\tilde{\zeta} = [\tilde{\rho}_{\text{FM}}^\text{T}, \tilde{\rho}_{\text{FR}}^\text{T}]^\text{T}$ ,  $\Gamma_\rho = \begin{bmatrix} \Gamma_{\rho,11} & \Gamma_{\rho,12} \\ \Gamma_{\rho,21} & \Gamma_{\rho,22} \end{bmatrix}$ ,  $\Gamma_{\rho,11} = \kappa_{\text{FM}} \mathbf{A}_{\text{FM}}$ ,  $\Gamma_{\rho,21} = \Gamma_{\rho,12}^\text{T} = \frac{1}{2} \kappa_{\text{FR}} \mathbf{J}_{\text{FR}} \mathbf{J}_{\text{FM}}^\dagger \mathbf{A}_{\text{FM}}$ , and  $\Gamma_{\rho,22} = \kappa_{\text{FR}} \mathbf{J}_{\text{FM}} (\mathbf{I} - \mathbf{J}_{\text{FM}} \mathbf{J}_{\text{FM}}^\dagger) \mathbf{J}_{\text{FR}}^\dagger \mathbf{A}_{\text{FR}}$ . Eq. (S22) satisfies the following inequality:

$$\dot{V}_\rho \leq -\underline{\Delta}_{\text{OA}} \|\tilde{\rho}_{\text{OA}}\|^2 - \underline{\Gamma}_{\rho,11} \|\tilde{\rho}_{\text{FM}}\|^2 - \underline{\Gamma}_{\rho,22} \|\tilde{\rho}_{\text{FR}}\|^2 + 2\overline{\Gamma}_{\rho,21} \|\tilde{\rho}_{\text{FM}}\| \|\tilde{\rho}_{\text{FR}}\| = -\hat{\zeta}^\text{T} \hat{\Gamma}_\rho \hat{\zeta}, \quad (\text{S23})$$

where  $\underline{\Gamma}_{\rho,11}$  and  $\underline{\Gamma}_{\rho,22}$  are the lower bounds on the induced norms of  $\Gamma_{\rho,11}$  and  $\Gamma_{\rho,22}$  respectively,  $\underline{\Gamma}_{\text{FM}} = \min_i \{\Gamma_{\text{FM},i}\}$ ,  $\underline{\Gamma}_{\text{FR}} = \min_i \{\Gamma_{\text{FR},i}\}$ ,  $\overline{\Gamma}_{\rho,21}$  is the upper bound on the induced norm of  $\Gamma_{\rho,12}$ ,  $\overline{\Gamma}_{\text{FM}} = \max_i \{\Gamma_{\text{FM},i}\}$ ,  $\hat{\zeta} = [\|\tilde{\rho}_{\text{FM}}\|^\text{T}, \|\tilde{\rho}_{\text{FR}}\|^\text{T}]^\text{T}$ ,  $\hat{\Gamma}_\rho = \begin{bmatrix} \hat{\Gamma}_{\rho,11} & \hat{\Gamma}_{\rho,12} \\ \hat{\Gamma}_{\rho,21} & \hat{\Gamma}_{\rho,22} \end{bmatrix}$ ,  $\hat{\Gamma}_{\rho,11} = \kappa_{\text{FM}} \underline{\Delta}_{\text{FM}}$ ,  $\hat{\Gamma}_{\rho,12} = \hat{\Gamma}_{\rho,21} = -\kappa_{\text{FR}} \overline{\Delta}_{\text{FR}}/2$ , and  $\hat{\Gamma}_{\rho,22} = \kappa_{\text{FR}} \underline{\Delta}_{\text{FR}}$ . When  $\kappa_{\text{FM}} > \frac{\kappa_{\text{FR}} \overline{\Delta}_{\text{FM}}^2}{4 \overline{\Delta}_{\text{FM}} \overline{\Delta}_{\text{FR}}}$ ,  $\hat{\Gamma}_\rho$  is a positive-definite symmetric matrix,

$$\dot{V}_\rho \leq -\underline{\Delta}_{\text{OA}} \|\tilde{\rho}_{\text{OA}}\|^2 - \underline{\lambda}_{\hat{\Gamma}} \hat{\zeta}^\text{T} \hat{\zeta} \leq 0, \quad (\text{S24})$$

where  $\underline{\lambda}_{\hat{\Gamma}}$  is the minimum eigenvalue of  $\hat{\Gamma}_\rho$ . The proof is completed.

### 3 Proof of boundedness

Since inequality (S17) only gives the total bound of all error signals, the bounds of each error signal should be further analyzed. Let us define the Lyapunov function candidates as  $V_\tau = \frac{1}{2} \sum_{i=1}^N \boldsymbol{\tau}_i^T \boldsymbol{\tau}_i$ ,  $V_e = \frac{1}{2} \sum_{i=1}^N \mathbf{e}_i^T \boldsymbol{\Upsilon}_{e,i} \mathbf{e}_i$ ,  $V_f = \frac{1}{2} \sum_{i=1}^N \text{tr} \left\{ \tilde{\boldsymbol{\omega}}_{f,i}^T \boldsymbol{\Upsilon}_{f,i}^{-1} \tilde{\boldsymbol{\omega}}_{f,i} \right\}$ ,  $V_c = \frac{1}{2} \sum_{i=1}^N \text{tr} \left\{ \tilde{\boldsymbol{\omega}}_{c,i}^T \tilde{\boldsymbol{\omega}}_{c,i} \right\}$ , and  $V_a = \frac{1}{2} \sum_{i=1}^N \text{tr} \left\{ \tilde{\boldsymbol{\omega}}_{a,i}^T \tilde{\boldsymbol{\omega}}_{a,i} \right\}$ . Following the above proof, it is easy to yield the results as

$$\dot{V}_\tau \leq - \sum_{i=1}^N \left( \xi_{\tau,i} - \frac{1}{2} \right) \|\boldsymbol{\tau}_i\| + \frac{N}{2} \epsilon_\Delta. \quad (\text{S25})$$

Let  $\vartheta_\tau = 2\xi_{\tau,i} - 1$  and  $\vartheta_\Delta = \frac{N}{2} \epsilon_\Delta$ ; we then have

$$\dot{V}_\tau \leq -\vartheta_\tau V_\tau + \vartheta_\Delta. \quad (\text{S26})$$

The inequality is obtained by using Lemma 1 as

$$V_\tau \leq e^{-\vartheta_\tau t} V_\tau(0) + \frac{\vartheta_\Delta}{\vartheta_\tau} (1 - e^{-\vartheta_\tau t}), \quad (\text{S27})$$

where  $\boldsymbol{\tau}_i$  is bounded with  $\|\boldsymbol{\tau}_i\| \leq \epsilon_\tau = \left[ 2e^{-\vartheta_\tau t} V_\tau(0) + \frac{2\vartheta_\Delta}{\vartheta_\tau} (1 - e^{-\vartheta_\tau t}) \right]^{1/2}$ .

$$\dot{V}_e + \dot{V}_f \leq -\frac{1}{2} \sum_{i=1}^N \frac{\lambda_\xi}{\lambda_e} \mathbf{e}_i^T \boldsymbol{\Upsilon}_{e,i} \mathbf{e}_i - \frac{1}{2} \sum_{i=1}^N \frac{\nu_{f,i}}{\lambda_f} \text{tr} \left\{ \tilde{\boldsymbol{\omega}}_{f,i}^T \boldsymbol{\Upsilon}_{f,i}^{-1} \tilde{\boldsymbol{\omega}}_{f,i} \right\} + \Xi_\xi, \quad (\text{S28})$$

where  $\Xi_\xi = \sum_{i=1}^N \|\dot{\boldsymbol{v}}_{i,r}\|^2 + N\epsilon_f^2 + \epsilon_\tau + \frac{1}{4\beta_V^2} \sum_{i=1}^N \text{tr} \left\{ \hat{\boldsymbol{\omega}}_{a,i}^T \boldsymbol{\Psi}_{V,i}(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \boldsymbol{\Psi}_{V,i}^T(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \hat{\boldsymbol{\omega}}_{a,i} \right\}$ . Let  $\vartheta_e = \frac{\lambda_\xi}{\lambda_e}$ ,  $\vartheta_f = \frac{\nu_{f,i}}{\lambda_{f,i}}$ , and  $|\Xi_\xi| \leq \vartheta_\xi$ ; we then have

$$\dot{V}_e \leq -\vartheta_e V_e + \vartheta_\xi, \quad (\text{S29})$$

$$\dot{V}_f \leq -\vartheta_f V_f + \vartheta_\xi. \quad (\text{S30})$$

The inequalities are obtained by using Lemma 1 as

$$V_e \leq e^{-\vartheta_e t} V_e(0) + \frac{\vartheta_\xi}{\vartheta_e} (1 - e^{-\vartheta_e t}), \quad (\text{S31})$$

$$V_f \leq e^{-\vartheta_{f,1} t} V_f(0) + \frac{\vartheta_\xi}{\vartheta_f} (1 - e^{-\vartheta_{f,1} t}), \quad (\text{S32})$$

where  $\mathbf{e}_i$  is bounded with  $\|\mathbf{e}_i\| \leq \epsilon_e = \left[ 2e^{-\vartheta_e t} V_e(0) + \frac{2\vartheta_\xi}{\vartheta_e} (1 - e^{-\vartheta_e t}) \right]^{1/2}$ ,  $\tilde{\boldsymbol{\omega}}_{f,i}$  is bounded with  $\|\tilde{\boldsymbol{\omega}}_{f,i}\| \leq \epsilon_f = \left\{ \left[ 2e^{-\vartheta_{f,1} t} V_f(0) + \frac{2\vartheta_\xi}{\vartheta_f} (1 - e^{-\vartheta_{f,1} t}) \right] / \Delta_f \right\}^{1/2}$ ,  $\Delta_f$  is the minimal eigenvalue of  $\boldsymbol{\Upsilon}_{f,i}^{-1}$ .

$$\dot{V}_c + \dot{V}_a \leq -\frac{1}{2} \sum_{i=1}^N \lambda_\psi \text{tr} \left\{ \tilde{\boldsymbol{\omega}}_{c,i}^T \tilde{\boldsymbol{\omega}}_{c,i} \right\} - \frac{1}{2} \sum_{i=1}^N \lambda_\psi \text{tr} \left\{ \tilde{\boldsymbol{\omega}}_{a,i}^T \tilde{\boldsymbol{\omega}}_{a,i} \right\} + \Xi_\psi, \quad (\text{S33})$$

where  $\Xi_\psi = \frac{1}{4\beta_V^2} \sum_{i=1}^N \text{tr} \left\{ \hat{\boldsymbol{\omega}}_{a,i}^T \boldsymbol{\Psi}_{V,i}(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \boldsymbol{\Psi}_{V,i}^T(\boldsymbol{\chi}_i, \boldsymbol{\tau}_i) \hat{\boldsymbol{\omega}}_{a,i} \right\}$  with  $|\Xi_\psi| \leq \vartheta_\psi$ . Let  $\vartheta_{c,1} = \lambda_\psi$  and  $\vartheta_{a,1} = \lambda_\psi$ ; we then have

$$\dot{V}_c \leq -\vartheta_c V_c + \vartheta_\psi, \quad (\text{S34})$$

$$\dot{V}_a \leq -\vartheta_a V_a + \vartheta_\psi. \quad (\text{S35})$$

The inequalities are obtained by using Lemma 1 as

$$V_c \leq e^{-\vartheta_c t} V_c(0) + \frac{\vartheta_\Psi}{\vartheta_c} (1 - e^{-\vartheta_c t}), \quad (\text{S36})$$

$$V_a \leq e^{-\vartheta_a t} V_a(0) + \frac{\vartheta_\Psi}{\vartheta_a} (1 - e^{-\vartheta_a t}), \quad (\text{S37})$$

where  $\tilde{\omega}_{c,i}$  is bounded with  $\|\tilde{\omega}_{c,i}\| \leq \epsilon_c = \left[ 2e^{-\vartheta_c t} V_c(0) + \frac{2\vartheta_\Psi}{\vartheta_c} (1 - e^{-\vartheta_c t}) \right]^{1/2}$  and  $\tilde{\omega}_{a,i}$  is bounded with  $\|\tilde{\omega}_{a,i}\| \leq \epsilon_a = \left[ 2e^{-\vartheta_a t} V_a(0) + \frac{2\vartheta_\Psi}{\vartheta_a} (1 - e^{-\vartheta_a t}) \right]^{1/2}$ .

---

**Algorithm S1** MARLMS

---

**Input:** total number of training episodes  $T_e$ , total number of time steps  $T_s$

- 1: Initialize  $Q(\mathbf{s}_t, \mathbf{b}_t; \omega_Q, \omega_V, \omega_B) = V(\mathbf{s}_t; \omega_Q, \omega_V) + B(\mathbf{s}_t, \mathbf{b}_t; \omega_Q, \omega_B)$  with random weights
  - 2: Initialize experience replay buffer  $\mathfrak{D}$
  - 3: Initialize  $\bar{T}(\phi(\mathbf{s}_t))$ -greedy exploration and leniency  $\mathfrak{L}(\mathbf{s}_t, \mathbf{b}_t)$
  - 4: **For** episode=1 :  $T_e$  **do**
  - 5: Reset the joint state  $\mathbf{s}_t$  to initial value  $\mathbf{s}_0$
  - 6:   **For**  $t = 1 : T_s$  **do**
  - 7:      $Q(\mathbf{s}_{t+1}, \mathbf{b}_{t+1}; \omega_Q^-, \omega_V^-, \omega_B^-) = \frac{1}{\lambda} \sum_{l=1}^{\lambda} Q(\mathbf{s}_{t+1}, \mathbf{b}_{t+1}; \omega_{Q_{t-l}}^-, \omega_{V_{t-l}}^-, \omega_{B_{t-l}}^-)$
  - 8:      $y_{\mathbf{s}_t, \mathbf{b}_t} = \mathbb{E}_D[r + \gamma \max_{\mathbf{b}_{t+1}} Q(\mathbf{s}_{t+1}, \mathbf{b}_{t+1}; \omega_Q^-, \omega_V^-, \omega_B^-) | \mathbf{s}_t, \mathbf{b}_t]$
  - 9:      $\begin{cases} \omega_{Q_t}, \omega_{V_t}, \omega_{B_t} \approx \arg \min_{\omega_Q, \omega_V, \omega_B} \mathbb{E}_D[(y_{\mathbf{s}_t, \mathbf{b}_t} - Q(\mathbf{s}_t, \mathbf{b}_t; \omega_Q, \omega_V, \omega_B))^2], & \delta_t > 0 \text{ or } \vartheta > \mathfrak{L}_t \\ \omega_{Q_t}, \omega_{V_t}, \omega_{B_t} = \omega_{Q_t}, \omega_{V_t}, \omega_{B_t}, & \delta_t \leq 0 \text{ and } \vartheta \leq \mathfrak{L}_t \end{cases}$
  - 10:     Update  $\mathfrak{D}$ ,  $\bar{T}(\phi(\mathbf{s}_t))$ , and  $\mathfrak{L}_t$
  - 11:   **End for**
  - 12: **End for**
- Output:**  $Q(\mathbf{s}_t, \mathbf{b}_t; \omega_Q, \omega_V, \omega_B)$
-