

Journal Pre-proof

Online Barrier-Actor-Critic Learning for H_∞ Control with Full-State Constraints and Input Saturation

Yongliang Yang, Da-Wei Ding, Haoyi Xiong, Yixin Yin,
Donald C. Wunsch

PII: S0016-0032(19)30904-4
DOI: <https://doi.org/10.1016/j.jfranklin.2019.12.017>
Reference: FI 4334

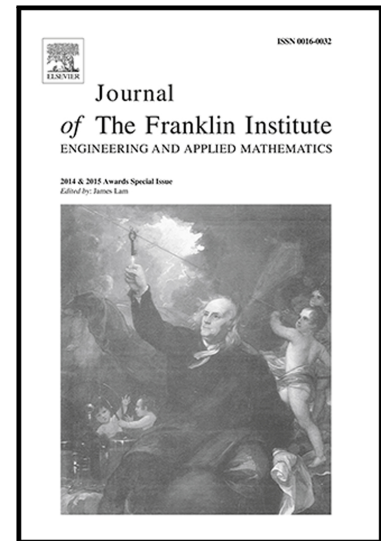
To appear in: *Journal of the Franklin Institute*

Received date: 17 June 2019
Revised date: 11 October 2019
Accepted date: 9 December 2019

Please cite this article as: Yongliang Yang, Da-Wei Ding, Haoyi Xiong, Yixin Yin, Donald C. Wunsch, Online Barrier-Actor-Critic Learning for H_∞ Control with Full-State Constraints and Input Saturation, *Journal of the Franklin Institute* (2019), doi: <https://doi.org/10.1016/j.jfranklin.2019.12.017>

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2019 Published by Elsevier Ltd on behalf of The Franklin Institute.



Online Barrier-Actor-Critic Learning for H_∞ Control with Full-State Constraints and Input Saturation[☆]

Yongliang Yang^{a,b,*}, Da-Wei Ding^{a,b}, Haoyi Xiong^{c,d}, Yixin Yin^{a,b}, Donald C. Wunsch^e

^a School of Automation & Electrical Engineering,
University of Science and Technology Beijing, Beijing 10083, China

^b Key Laboratory of Knowledge Automation for Industrial Processes, Ministry of Education,
University of Science and Technology Beijing, Beijing 10083, China

^c Big Data Laboratory, Baidu Research, Beijing 100193, China

^d National Engineering Laboratory of Deep Learning Technology and Application, Beijing 100193, China

^e Department of Electrical and Computer Engineering,
Missouri University of Science and Technology, Rolla, MO 65401, USA

Abstract

This paper develops a novel adaptive optimal control design method with full-state constraints and input saturation in the presence of external disturbance. First, to consider the full-state constraints, a barrier function is developed for system transformation. Moreover, it is shown that, with the barrier-function-based system transformation, the stabilization of the transformed system is equivalent to the original constrained control problem. Second, the disturbance attenuation problem is formulated within the zero-sum differential games framework. To determine the optimal control and the worst-case disturbance, a novel barrier-actor-critic algorithm is presented for adaptive optimal learning while guaranteeing the full-state constraints and input saturation. It is proven that the closed-loop signals remain bounded during the online learning phase. Finally, simulation studies are conducted to demonstrate the effectiveness of the presented barrier-actor-critic learning algorithm.

Keywords: full-state constraints, input saturation, disturbance attenuation, adaptive dynamic programming, barrier-actor-critic learning

[☆]This work was supported in part by the National Natural Science Foundation of China under Grant No. 61903028, No. 61873028 and No. 61333002, in part by the China Post-Doctoral Science Foundation under Grant 2018M641197, in part by the Fundamental Research Funds for the Central Universities under grant FRF-TP-18-031A1 and FRF-BD-19-002A, in part by the Lifelong Learning Machines program from DARPA/Microsystems Technology Office and in part by the Army Research Laboratory under Cooperative Agreement Number W911NF-18-2-0260. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the Army Research Laboratory or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation herein.

*Corresponding author: Da-Wei Ding (ddaweiauto@163.com)

Email addresses: yangyongliang@ieee.org (Yongliang Yang), xionghaoyi@baidu.com (Haoyi Xiong), yyx@ies.ustb.edu.cn (Yixin Yin), Wunsch@ieee.org (Donald C. Wunsch)

1. Introduction

Nonlinear dynamics commonly exists in engineering applications, such as input saturation [1, 2, 3] and dead-zone [4, 5], output constraints [6, 7], friction dynamics [8, 9], backlash-like hysteresis [10, 11, 12], unmodeled dynamics [13], etc. Modern control theory, such as the H_∞ control method [14, 15] and adaptive control method [16, 17], has received considerable attention to compensate for the system uncertainty and attenuate the effect of external disturbance for nonlinear systems. In addition to the closed-loop stability, practical constraints captured by user-defined performance is desired to be guaranteed. However, classical H_∞ control and adaptive control methods cannot guarantee the user-defined performance. In this paper, a novel adaptive optimal controller design is developed to stabilize the nonlinear systems while considering both the prescribed performance on full-state and input saturation simultaneously.

For the nonlinear systems with imperfect dynamical behavior, such as exogenous disturbance and system uncertainties, the adaptive control method is widely used for feedback design to compensate the system uncertainty and attenuate exogenous disturbances [16, 17]. However, classical adaptive control design methods only consider the closed-loop stability. In addition to the closed-loop stability, practical constraints are important for controller design. For example, in the control of Euler-Lagrange systems, the link and joint velocity cannot be arbitrarily large and has to be remained in the bounded region due to limitation imposed by mechanical characteristics. In many applications, the constraints are usually captured by the user-defined performance. Many efforts have been made to address this issue. Compared to classical quadratic Lyapunov function design, Lyapunov analysis is combined with barrier function design [18] to consider the constraints on output, which is essentially partial-state constraints [19, 20]. Since then, the barrier Lyapunov function design is extended to consider full-state constraints for stochastic nonlinear systems [21], pure-feedback systems [22], Euler-Lagrange systems [23], time-delay systems [24], to name a few. Another type of constrained controller design adopts a prescribed transient performance to develop a system transformation [25]. In the prescribed performance adaptive control, the prescribed transient performance is captured by a user-defined performance bound, which specifies the safety region for the tracking error. Recently, the prescribed performance adaptive control method is extended to deal with output feedback control problem [26], consensus problem of multi-agent systems [27], nonlinear systems with input dead-zone [28], controller design for flexible joint robots [29], synchronization problem of teleoperation robotics [30], and so on. To relax the requirement that both the reference signal along with its derivatives and every element of the state variable are available for feedback design, Arabia and Yucelen developed a set-theoretic model reference

34 adaptive control framework [31]. In the set-theoretic model reference adaptive control framework,
 35 the norm of the gap between the system state and the reference signal is guaranteed to be within
 36 a user-defined constant bound. However, in the existing adaptive controller design methods, only
 37 closed-loop stability and the prescribed user-defined performance constraints is considered with-
 38 out optimality discussion. In this paper, a novel adaptive constrained controller is presented with
 39 optimality discussions.

40 The centerpiece of optimal control theory is the Hamilton-Jacobi-Bellman/ Hamilton-Jacobi-
 41 Isaacs (HJB/HJI) equations for nonlinear systems, which is necessary and sufficient for the opti-
 42 mality condition [32]. However, the HJ equations are difficult to solve due to the inherent non-
 43 linearity. Therefore, adaptive dynamic programming (ADP) has been developed to approximate
 44 the nonlinear HJ equations in an online fashion, where an intelligent agent seeks optimal decisions
 45 to maximize the long-term cumulative reward [33]. Variants of ADP has been applied widely in
 46 control applications to solve the optimal control problems, including iterative ADP algorithms in
 47 discrete-time [34] and continuous-time [35] for optimal regulation problems, model-free learning
 48 algorithm for H_∞ control problem [36], online actor-critic learning algorithm [37] for optimal track-
 49 ing control problems [38, 39], robust stabilization problem [40], guaranteed cost control problem
 50 [41, 42], consensus control problem of multi-agent systems [43, 44], event-triggered control [45], to
 51 name a few. Besides, ADP has been successfully applied to differential games [46]. In addition,
 52 ADP extensions have been made to deal with constraints of input saturation in [47] and constraints
 53 on the state in [48]. However, these existing results do not consider the case with external distur-
 54 bance, input saturation, and full-state constraints. In this paper, all these issues are considered in
 55 a comprehensive framework.

56 The contributions of this paper are threefold. First, in this paper, both the full-state con-
 57 straints and input saturation are considered simultaneously for the controller design problem. This
 58 is achieved by introducing a barrier function based system transformation. It is also discussed the-
 59 oretically that the transform equivalence can be guaranteed in the sense that the stabilization of
 60 the transformed system ensures the constraints of the original system. Second, the disturbance
 61 attenuation is achieved within the framework of zero-sum differential games. A novel barrier-actor-
 62 critic algorithm is developed for adaptive optimal learning with the full-state constraints and input
 63 saturation. Finally, to obviate the requirement of persistent excitation condition, the experience
 64 replay technique is employed to utilize the history and current data concurrently.

65 The remainder of this paper is organized as follows. In Section 2, the problem of constrained
 66 control design with full-state constraints and input saturation is given. Section 3 presents the

67 barrier-function-based system transformation to deal with full-state constraints. In Section 4, a
 68 novel actor-critic-barrier algorithm is developed for the online learning of the adaptive optimal
 69 constrained controller.

70 2. Preliminaries

71 2.1. Notations and Definitions

72 The following standard notation will be adopted.

\mathbb{R}^+	\triangleq	the set of positive real numbers.
\mathbb{R}^n	\triangleq	n -dimensional vector space.
I	\triangleq	Identity matrix with proper dimension.
$\mathbf{1}$	\triangleq	vector with all entries being 1.
$\ \mathcal{M}\ $	\triangleq	$\sqrt{\text{tr}(\mathcal{M}\mathcal{M}^H)}$, the matrix Frobenius norm of matrix \mathcal{M} .
$\ v\ $	\triangleq	the euclidean norm of vector v .
\mathbb{Z}	\triangleq	the set of integers.
$\lambda_{\min}(A)$	\triangleq	the minimum eigenvalue of matrix A .

74 **Definition 1.** (Zero-State Observability) [15] The system (1) with the measured output $y = h(x)$
 75 is zero-state observable if $y(t) \equiv 0$ for $\forall t \geq 0$ implies that $x(t) \equiv 0$ for $\forall t \geq 0$.

76 **Definition 2.** (Persistent Excitation Condition) [16] The vector signal $z(\cdot) \in \mathbb{R}^n$ is said to be
 77 persistently excited (PE) on the interval $[T_1, T_2]$ if there exists positive constants $\gamma_1 > 0$ and
 78 $\gamma_2 > 0$ such that, for all $t \in [T_1, T_2]$,

$$\gamma_1 I \leq \int_t^{t+T_1} z(\tau) z^T(\tau) d\tau \leq \gamma_2 I$$

79 **Definition 3.** (Uniformly Ultimately Bounded Stability) [16] Consider the nonlinear system

$$\dot{x} = F(x, t), \quad \forall t \in \mathbb{R}^+, \quad x(t_0) = x_0 \quad (1)$$

80 with $x(t) \in \mathbb{R}^n$ is the system state and x_0 is the initial condition. The equilibrium point x_e of
 81 system (1) is said to be uniformly ultimately bounded (UUB) if there exists a compact set $\Omega \subset \mathbb{R}^n$
 82 so that for all $x_0 \in \Omega$, there exists a bound B and a time $T(B, x_0)$ such that $\|x(t) - x_0\| \leq B$ for
 83 all $t \geq t_0 + T$.

84 **Lemma 1.** [49] For $\forall w \in \mathbb{R}$, there exists a bounded \tilde{w} satisfying $\|\tilde{w}\| \leq \ln 4$, such that

$$-2 \ln(1 + e^{-2w}) = 2w - 2w \text{sgn}(w) + \tilde{w},$$

85 **Lemma 2.** [50] *The following inequality holds for any $a > 0$ and $y \in \mathbb{R}$*

$$0 \leq |y| - y \tanh\left(\frac{y}{a}\right) \leq \kappa a \quad (2)$$

86 where $\kappa = 0.2785$.

87 **2.2. Problem Statement**

88 In this paper, we consider the following continuous-time affine nonlinear dynamical systems

$$\begin{aligned} \dot{x}_1 &= x_2 \\ \dot{x}_2 &= x_3 \\ &\vdots \\ \dot{x}_{n-1} &= x_n \\ \dot{x}_n &= f(x) + g(x)u + k(x)d \end{aligned} \quad (3)$$

89 where $x = [x_1 \ \dots \ x_n]^T \in \mathbb{R}^n$ is the system state, $u(\cdot) : \mathbb{R}^n \rightarrow \mathbb{R}^{m_1}$ is the control policy,
90 $d(\cdot) : \mathbb{R}^n \rightarrow \mathbb{R}^{m_2}$ is the external disturbance, $f(\cdot), g(\cdot), k(\cdot) : \mathbb{R}^n \rightarrow \mathbb{R}$ are Lipschitz continuous
91 nonlinear functions. The constrained H_∞ control problem for system (3) with full-state constraints
92 and input saturation can be formulated as follows.

93 **Problem 1.** Design the proper performance output $L(x, u)$, where $L(\cdot, \cdot)$ is a positive definite
94 function of its argument, and the optimal policy u^* such that, with the saturation constraints on
95 the control input as

$$\|u_i\| \leq \lambda, \forall i = 1, \dots, m_1 \quad (4)$$

96 where $u = [u_1 \ \dots \ u_m]^T$, and the full-state constraints as

$$\begin{aligned} x_1 &\in (a_1, A_1) \\ &\vdots \\ x_n &\in (a_n, A_n) \end{aligned} \quad (5)$$

97 for $\forall d \in \mathcal{L}_2$, system (3) have L_2 -gain less than or equal to γ , i.e.,

$$\frac{\int_t^\infty L(x(\tau), u(\tau))d\tau}{\int_t^\infty \|d(\tau)\|^2 d\tau} \leq \gamma^2, \quad (6)$$

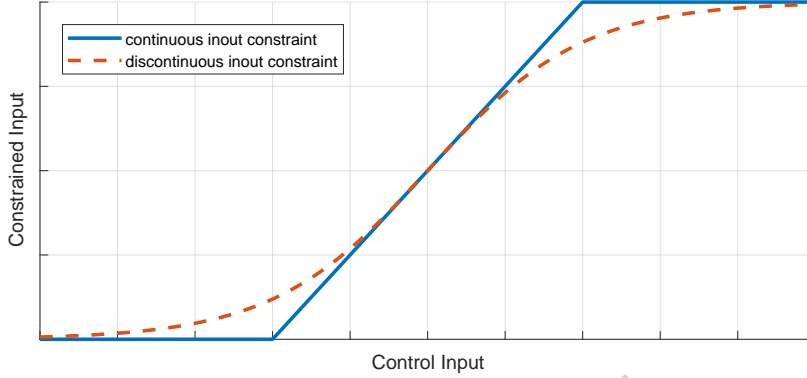


Figure 1: Evolution of the two-dimensional phase plot of the state trajectories $[x_1(t) \ x_2(t)]$. The black box denotes the safe region.

Remark 1. To deal with the input constraints (4), the saturation function can be applied, which is defined as [51, 52, 53, 54]

$$\Gamma(u_i) = \begin{cases} u_i, & \text{if } u_i \leq \lambda \\ \text{sign}(u_i), & \text{if } u_i > \lambda \end{cases}$$

Then, the system dynamics can be denoted as

$$\begin{aligned} \dot{x}_i &= x_{i+1}, \quad i = 1, 2, \dots, n-1 \\ \dot{x}_n &= f(x) + g(x)\Gamma(u) + k(x)d \end{aligned}$$

98 Note that the saturation function $\Gamma(\cdot)$ is a discontinuous function, which leads to discontinuity in
 99 the system dynamics. In this paper, we consider continuous constraints on the input signal, which
 100 is shown in Figure 1 and widely used in the literature, such as [33, 47, 55]. As shown later, the
 101 nonquadratic penalty function on the control input signal (17) is presented, which guarantees the
 102 boundedness of the optimal control input (23).

103 As shown by (4) – (6), the objective of Problem 1 can be divided into three parts, i.e., distur-
 104 bance attenuation, input saturation and full-state constraints. For the full-state constraints, we
 105 introduce the following barrier function.

106 **Definition 4.** (Barrier Function) The function $B(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$ defined on (a, A) is referred to as
 107 barrier function if

$$B(z; a, A) = \ln \left(\frac{A}{a} \frac{a-z}{A-z} \right), \quad \forall z \in (a, A) \quad (7)$$

108 where a and A are two constants satisfying $a < A$. Moreover, the barrier function is invertible on
 109 interval (a, A) , i.e.,

$$B^{-1}(y; a, A) = aA \frac{e^{\frac{y}{2}} - e^{-\frac{y}{2}}}{ae^{\frac{y}{2}} - Ae^{-\frac{y}{2}}}, \forall y \in \mathbb{R} \quad (8)$$

110 with the derivative

$$\frac{dB^{-1}(y; a, A)}{dy} = \frac{Aa^2 - aA^2}{a^2e^y - 2aA + A^2e^{-y}} \quad (9)$$

111 *Remark 2.* To guarantee that the full-state constraints is not violated for Problem 1, the barrier
 112 function in Definition 4 has the following desired properties

- 113 1) The barrier function $B(\cdot)$ takes finite value when the its arguments are within the user-defined
 114 region (a, A) .
 115 2) The barrier function $B(\cdot)$ approach to infinity as the state approach the boundary of the
 116 prescribed region (a, A) , i.e.,

$$\begin{aligned} \lim_{z \rightarrow a^+} B(z; a, A) &= -\infty \\ \lim_{z \rightarrow A^-} B(z; a, A) &= +\infty \end{aligned}$$

- 117 3) The barrier function $B(\cdot)$ vanishes at the equilibrium of the system (3), i.e.,

$$B(0; a, A) = 0, \forall a < A$$

118 3. Barrier-Function-Based Zero-Sum Game

119 In this section, the system (3) with full-state constraints is transformed into another system
 120 without state constraints by using the barrier function in Definition 4. Consider the barrier-
 121 function-based state transformation as

$$\begin{aligned} s_i &= B(x_i; a_i, A_i), \\ x_i &= B^{-1}(s_i; a_i, A_i), i = 1, \dots, n \end{aligned} \quad (10)$$

122 Then, by using the chain rule, one has

$$\frac{dx_i}{dt} = \frac{dx_i}{ds_i} \frac{ds_i}{dt} \quad (11)$$

123 From (11), the dynamics of the transformed state s can be written as

$$\begin{aligned}
 \dot{s}_i &= \frac{x_{i+1}(s_{i+1})}{\left. \frac{dB^{-1}(y; a_i, A_i)}{dy} \right|_{y=s_i}} \\
 &= \frac{a_{i+1}A_{i+1} \left(e^{\frac{s_{i+1}}{2}} - e^{-\frac{s_{i+1}}{2}} \right)}{a_{i+1}e^{\frac{s_{i+1}}{2}} - A_{i+1}e^{-\frac{s_{i+1}}{2}}} \frac{A_i^2 e^{-s_i} - 2a_i A_i + a_i^2 e^{s_i}}{A_i a_i^2 - a_i A_i^2} \\
 &= F_i(s_i, s_{i+1}), \quad i = 1, \dots, n-1 \\
 \dot{s}_n &= \frac{f(x) + g(x)u + k(x)d}{\left. \frac{dB^{-1}(y; a_n, A_n)}{dy} \right|_{y=s_n}} \\
 &= [f(x) + g(x)u + k(x)d] \frac{A_n^2 e^{-s_n} - 2a_n A_n + a_n^2 e^{s_n}}{A_n a_n^2 - a_n A_n^2} \\
 &= F_n(s) + g_n(s)u + k_n(s)d
 \end{aligned} \tag{12}$$

124 with

$$\begin{aligned}
 F_n(s) &= \frac{A_n^2 e^{-s_n} - 2a_n A_n + a_n^2 e^{s_n}}{A_n a_n^2 - a_n A_n^2} f \left(\left[B_1^{-1}(s_1) \quad \dots \quad B_n^{-1}(s_n) \right] \right) \\
 g_n(s) &= \frac{A_n^2 e^{-s_n} - 2a_n A_n + a_n^2 e^{s_n}}{A_n a_n^2 - a_n A_n^2} g \left(\left[B_1^{-1}(s_1) \quad \dots \quad B_n^{-1}(s_n) \right] \right) \\
 k_n(s) &= \frac{A_n^2 e^{-s_n} - 2a_n A_n + a_n^2 e^{s_n}}{A_n a_n^2 - a_n A_n^2} k \left(\left[B_1^{-1}(s_1) \quad \dots \quad B_n^{-1}(s_n) \right] \right)
 \end{aligned} \tag{13}$$

125 Note that system (12) with the state $s = \left[s_1 \quad \dots \quad s_n \right]^T$ can be expressed in a compact form
 126 as

$$\dot{s} = F(s) + G(s)u + K(s)d \tag{14}$$

127 with $F(s) = \begin{bmatrix} F_1(s_1, s_2) \\ \vdots \\ F_n(s) \end{bmatrix}$, $G(s) = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ g_n(s) \end{bmatrix}$, $K(s) = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ k_n(s) \end{bmatrix}$.

128 The following assumptions are imposed on system (14), which is commonly used for nonlinear
 129 systems controller design [47, 55].

130 **Assumption 1.** The system dynamics (14) is assumed to have the following properties.

- 131 1) $F(s)$ is Lipschitz with $F(0) = 0$, and there exists a constant b_f such that, for $s \in \Omega$, $\|F(s)\| \leq$
 132 $b_f \|s\|$ where Ω is a compact set containing the origin.
- 133 2) $G(s)$ and $K(s)$ are bounded on Ω , i.e., there exists a constant b_g and b_k such that $\|G(s)\| \leq b_g$
 134 and $\|K(s)\| \leq b_k$, respectively.

135 3) The system (3) is controllable over the compact set Ω .

136 In the following, to consider the input saturation and disturbance attenuation in Problem 1,
 137 the framework of the zero-sum differential game is introduced. For system (14) with the control
 138 input $u(t)$ and the disturbance policy $d(t)$, consider the following cost function

$$V(s_0; u, d) = \int_{t_0}^{\infty} U(s, u, d) dt \quad (15)$$

139 where $U(s, u, d)$ is the reward function with

$$\begin{aligned} U(s, u, d) &= L(x, u) - \gamma^2 \|d\|^2, \\ L(x, u) &= Q(s) + \Theta(u) - \gamma^2 \|d\|^2 \end{aligned} \quad (16)$$

140 where $Q(s)$ being a positive definite monotonically increasing function and $\Theta(u)$ being a positive
 141 definite integrand function. To deal with input saturation, the nonquadratic penalty function is
 142 used,

$$\begin{aligned} \Theta(u) &= 2 \int_0^u \left[\lambda \tanh^{-1} \left(\frac{v}{\lambda} \right) \right] R dv \\ &= 2\lambda (\tanh^{-1}(u/\lambda))^T Ru + \lambda^2 \bar{R} \ln(1 - u^2/\lambda^2) \end{aligned} \quad (17)$$

143 where $\lambda > 0$ is the saturation limit for the control input, $R = \text{diag}(r_1, \dots, r_m)$ and $\bar{R} = [r_1, \dots, r_m] \in$
 144 $\mathbb{R}^{1 \times m}$ with $r_i > 0$ for $i = 1, \dots, m$ is the weight on control effort for each input.

145 **Problem 2.** For system (14) with the control policy u and disturbance policy d , find the Nash
 146 equilibrium (u^*, d^*) of the zero-sum game with the constraints of input saturation (4).

147 Define the Hamiltonian for the cost (15) with the control policy u and disturbance policy d as

$$H(u, d, V) = \left(\frac{\partial V}{\partial s} \right)^T [F(s) + G(s)u + K(s)d] + U(s, u, d) \quad (18)$$

148 Then, differential equivalent of the cost (15) can be expressed in terms of the Hamiltonian (18) as

$$H \left(s, u, d, \frac{\partial V}{\partial s} \right) = 0 \quad (19)$$

149 which is referred to as the Bellman equation.

150 Based on the game theory [56], the disturbance attenuation problem is equivalent to solving
 151 the following two-player zero-sum game,

$$V^*(s) = \min_u \max_d V(s; u, d) \quad (20)$$

152 This two-player zero-sum game has a unique solution if the Nash condition holds

$$V^*(s) = \min_u \max_d V(s; u, d) = \max_d \min_u V(s; u, d) \quad (21)$$

153 According to [32], the stationary condition for optimality is

$$\frac{\partial H(u, d, V^*)}{\partial u} = 0, \quad \frac{\partial H(u, d, V^*)}{\partial d} = 0 \quad (22)$$

154 Then, one can obtain the optimal control input u^* and the worst-case disturbance d^* , respectively,

155 as

$$u^*(s) = -\lambda \tanh\left(\frac{1}{2\lambda} R^{-1} G^T(s) \frac{\partial V^*(s)}{\partial s}\right) \quad (23)$$

$$d^*(s) = \frac{1}{2\gamma^2} K^T(s) \frac{\partial V^*(s)}{\partial s} \quad (24)$$

156 where (u^*, d^*) is termed as Nash equilibrium for zero-sum game. Inserting the optimal control
157 policy and disturbance term (23) in (17) results in [55]

$$\Theta(u^*) = \lambda \left[\frac{\partial V^*(s)}{\partial s} \right]^T G(s) \tanh(D^*) + \lambda^2 \bar{R} \ln[1 - \tanh^2(D^*)] \quad (25)$$

158 where $D^* = (1/2\lambda) R^{-1} G(s)^T \frac{\partial V^*(s)}{\partial s}$. Inserting the Nash equilibrium (u^*, d^*) into (19) and using
159 (25), the Bellman equation becomes the Hamilton-Jacobi-Isaacs (HJI) equation

$$0 = Q(s) + \left[\frac{\partial V^*(s)}{\partial s} \right]^T F(s) + \lambda^2 \bar{R} \ln[1 - \tanh^2(D^*)] \\ + \frac{1}{4\gamma^2} \left[\frac{\partial V^*(s)}{\partial s} \right]^T K(s) K(s)^T \frac{\partial V^*(s)}{\partial s} \quad (26)$$

160 The following assumption on the cost function (15), which has been widely used in [14, 15], is
161 employed in this paper.

162 **Assumption 2.** The performance functional (15) satisfies zero-state observability.

163 The following lemma discusses the equivalence between Problems 1 and 2

164 **Lemma 3.** Suppose that the pair of policy $\{u^*(\cdot), d^*(\cdot)\}$ solve Problem 2 for system (14). Then,
165 the optimal control policy $\{u^*(\cdot)\}$ also solves Problem 1 provided that the initial state x_0 of system
166 (3) satisfies the constraints in (5).

167 Under Assumptions 1 and 2, suppose that $\mu^* = \{u^*, d^*\}$ solves Problem 2 for system (14) with
168 performance (15) and reward (16), then the following hold:

169 1) The closed-loop system satisfies the constraints (5) provided that the initial state x_0 of system
170 (3) is within the region (a_i, A_i) , $\forall i = 1, \dots, n$.

171 2) The disturbance attenuation condition (6) can be guaranteed if the performance output $L(x, u)$
 172 is designed as

$$L(x, u) = Q(s) + \Theta(u).$$

Proof. 1) Based on Assumptions 1 and 2, the existence of a positive definite and continuously differentiable optimal value function $V^*(s)$ can be guaranteed. From (18), one can obtain that $\dot{V}^*(t) \leq 0$, i.e.,

$$V^*(s(t)) \leq V^*(s(0)), \forall t \geq 0.$$

Then, $V^*(s(t))$ remains bounded if $V^*(s(0))$ is bounded, which is guaranteed by the condition that the initial condition $x(0)$ of system (3) satisfies the constraints in (5). Finally, from the discussions in Remark 2, one can infer that

$$x_i(t) \in (a_i, A_i), \quad i = 1, 2, \dots, n.$$

173 Therefore, given $\mu^* = \{u_1^*, u_2^*\}$, the constraints of Problem 1 are satisfied.

174 2) Now consider the barrier-function-based state transformation described by (10). Then, each
 175 element of the state $s = \begin{bmatrix} b_1(x_1) & \dots & b_n(x_n) \end{bmatrix}^T$ is finite given that x satisfies the constraints
 176 given in (5). Note that the Nash equilibrium (u^*, d^*) and the optimal value function V^* satisfies the
 177 Bellman equation (19), i.e., $H\left(s, u^*, d^*, \frac{\partial V^*}{\partial s}\right) = 0$. Then, considering (16) and the performance
 178 output $L(x, u)$, one has,

$$H\left(s, u^*, d^*, \frac{\partial V^*}{\partial s}\right) = 0 \Rightarrow \frac{\int_t^\infty \|z(\tau)\|^2 d\tau}{\int_t^\infty \|d(\tau)\|^2 d\tau} \leq \gamma^2$$

179 provided that $L(x, u) = Q(s) + \Theta(u)$. This completes the proof. ■

180 *Remark 3.* As shown in (25), the optimal constrained control and disturbance solution $u^*(s)$ and
 181 $d^*(s)$ depend on solving the HJI equation (26) for the optimal value function $V^*(s)$. However,
 182 the HJI equation (26) is a nonlinear partial differential and extremely difficult to solve. In the
 183 following, an online algorithm is presented to find an approximate solution to the HJI equation
 184 (26).

185 4. Online Actor-Critic-Barrier Learning

186 As shown in Lemma 3, with the barrier-function-based system transformation (10), the equiv-
 187 alence between Problems 1 and 2 can be guaranteed. In this section, we present a novel barrier-
 188 actor-critic online algorithm to learn the optimal control policy and the worst disturbance with

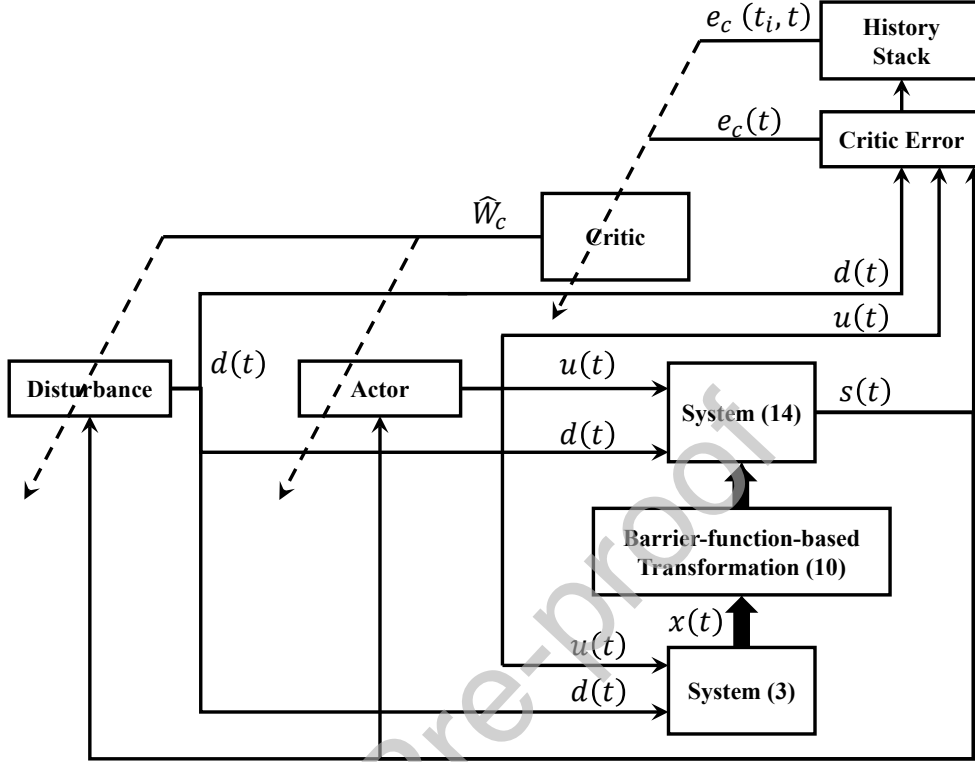


Figure 2: The overall barrier-actor-critic algorithm for disturbance attenuation with input saturation and full-state constraints. 1) Based on the barrier function defined in Definition 4, a novel system transformation is applied to original system (3) to obtain the transformed system (14). 2) The barrier-function-based system transformation is then combined with the actor-critic online algorithm to learn the optimal control policy u^* and worst-case disturbance d^* . 3) To obviate the requirement of PE condition for online critic learning, the experience replay technique is employed to concurrently utilize the online and history data.

189 respect to the performance of Problem 2. First, value function approximation for the critic learning
 190 is represented by using neural networks. Online critic learning is designed to approximate the HJI
 191 equation (26). In addition, two actor NNs are designed to learn the optimal control policy (23)
 192 and the worst-case disturbance (24), respectively.

193 4.1. Value Function Approximation

194 Using the NN approximation theorem, there exists a single-layer NN such that the value func-
 195 tion $V(s)$ and its gradient $\nabla V(s)$ can be uniformly approximated with a critic NN as the number

196 of basis sets increases, within a compact set $\Omega \subseteq \mathbb{R}^n$ that contains the origin, as

$$V^*(s) = (W^*)^T \phi(s) + \varepsilon(s) \quad (27)$$

$$\nabla V^*(s) = [\nabla \phi(s)]^T W^* + \nabla \varepsilon(s) \quad (28)$$

197 where $W^* \in \mathbb{R}^N$ is an ideal weight vector for the best N -dimensional value function approximation,
 198 $\phi(\cdot) : \mathbb{R}^n \rightarrow \mathbb{R}^N$ is the NN basis function, $\nabla = \partial/\partial s$, $\varepsilon(s)$ and $\nabla \varepsilon(s)$ are the NN approximation
 199 residual. For the value function approximation (27) and (28), the following standard assumption
 200 is adopted in this paper.

201 **Assumption 3.** The value function approximation as shown in (27) and (28) are assumed to have
 202 the following properties.

- 203 1) The ideal weight W is bounded by a constant, i.e., $\|W^*\| \leq b_*$;
 204 2) The value function approximation residual ε and $\nabla \varepsilon$ satisfies $\|\varepsilon(s)\| \leq b_\varepsilon$ and $\|\nabla \varepsilon(s)\| \leq b_{d\varepsilon}$;
 205 3) The NN basis function $\phi(s)$ and its gradient $\nabla \phi(s)$ satisfies $\|\phi(s)\| \leq b_\phi$ and $\|\nabla \phi(s)\| \leq b_{d\phi}$
 206 for $\forall s \in \Omega$.

207 For the optimal control policy $u^*(s)$ and the optimal disturbance inputs $d^*(s)$, the Bellman
 208 equation (19) approximation error using the value function approximation (27) can be expressed
 209 as

$$\xi = U(s, u^*, d^*) + (W^*)^T \sigma \quad (29)$$

210 where σ is a N -dimensional vector signal defined as

$$\sigma = \nabla \phi(s) [F(s) + G(s) u^* + K(s) d^*] \quad (30)$$

211 Considering the value gradient approximation (28), one can obtain that the Bellman residual
 212 results from the value gradient approximation error $\nabla \varepsilon(s)$, i.e.,

$$\xi = -[\nabla \varepsilon(s)]^T [F(s) + G(s) u^* + K(s) d^*] \quad (31)$$

213 Similarly, with the value function approximation (27), the HJI equation (26) can be approximated
 214 with a residual expressed as

$$\begin{aligned} \zeta &= Q(s) + (W^*)^T \sigma + \Theta(-\lambda \tanh(D_u)) - \frac{1}{4\gamma^2} (W^*)^T D_d W^* \\ &= Q(s) + (W^*)^T \nabla \phi(s) F(s) + \lambda^2 \bar{R} \ln(1 - \tanh^2(D_u^*)) + \frac{1}{4\gamma^2} (W^*)^T D_d W^* \end{aligned} \quad (32)$$

215 with

$$\begin{aligned} D_u^* &= \frac{1}{2\lambda} R^{-1} G^T(s) [\nabla\phi(s)]^T W^* \\ D_d &= \nabla\phi(s) K(s) K^T(s) [\nabla\phi(s)]^T \end{aligned} \quad (33)$$

216 *Remark 4.* From Assumptions 1 and 3, the policy representations in (23) and (24), the Bellman
217 equation approximation residual ξ is bounded in the sense that there exists a constant b_ξ such
218 that $\|\xi\| \leq b_\xi$. Similarly, the HJI approximation residual ζ using the ideal N -dimensional value
219 function approximation (27) and (28) is bounded as $\zeta \leq b_\zeta$.

220 4.2. Critic Learning

221 The ideal weight, W in (27), provides the best approximate to the optimal value function $V^*(s)$
222 on the compact set Ω and is unknown. Therefore, the estimation of W is implemented by the critic
223 network with the approximations of the value function and value gradient

$$\hat{V}(s) = \hat{W}_c^T \phi_c(s) \quad (34)$$

$$\nabla\hat{V}(s) = [\nabla\phi_c(s)]^T \hat{W}_c \quad (35)$$

224 Then, for a given policy $u(\cdot)$, the residual of Bellman equation approximation using the identifier
225 NN and the critic NN, can be determined as

$$\begin{aligned} e_c(t) &= U(s(t), u(t), d(t)) + \hat{W}_c^T \sigma(t) \\ &= -(\nabla\varepsilon)^T [F(s(t)) + G(s(t))u(t) + K(s(t))d(t)] \end{aligned} \quad (36)$$

226 Define the critic weight approximation error as

$$\tilde{W}_c = W^* - \hat{W}_c \quad (37)$$

227 Then, from (29), the relation between Bellman residual e_c and the Bellman equation approximation
228 error ζ can be written in terms of the critic weight error \tilde{W}_c as

$$e_c(t) = \xi(t) - \tilde{W}_c^T(t)\sigma(t) \quad (38)$$

$$e_c(t_i, t) = \xi(t_i) - \tilde{W}_c^T(t)\sigma(t_i) \quad (39)$$

229 Then $e_c \rightarrow \xi$ as $\hat{W}_c \rightarrow W^*$. The policy evaluation for an admissible control policy $u(\cdot)$ can be
230 formulated as adapting the critic weight \hat{W}_c to minimize the objective function

$$E_c = \frac{1}{2} \left(\frac{[e_c(t)]^2}{(1 + \sigma^T(t)\sigma(t))^2} + \sum_{i=1}^k \frac{[e_c^2(t_i, t)]^2}{(1 + \sigma^T(t_i)\sigma(t_i))^2} \right) \quad (40)$$

231 Using the chain rule yields adaptive critic online learning as

$$\begin{aligned}
 \dot{\tilde{W}}_c &= -\alpha_c \frac{\partial E_c}{\partial \tilde{W}_c} \\
 &= -\alpha_c \frac{\sigma(t) e_c(t)}{[1 + \sigma^T(t) \sigma(t)]^2} - \alpha_c \sum_{i=1}^k \frac{\sigma(t_i) e_c(t_i, t)}{[1 + \sigma^T(t_i) \sigma(t_i)]^2} \\
 &= -\alpha_c \frac{\sigma(t) [\xi(t) - \sigma(t)^T \tilde{W}_c(t)]}{[1 + \sigma^T(t) \sigma(t)]^2} - \alpha_c \sum_{i=1}^k \frac{\sigma(t_i) [\xi(t_i) - \sigma(t_i)^T \tilde{W}_c(t)]}{[1 + \sigma^T(t_i) \sigma(t_i)]^2}
 \end{aligned} \tag{41}$$

232 where $\alpha_c > 0$ is the critic learning rate.

233 **Condition 1.** The recorded data matrix $\begin{bmatrix} \sigma(t_1) & \dots & \sigma(t_k) \end{bmatrix}$ is full column rank.

234 **Theorem 1.** Let u be any given admissible control policy. Then, under Condition 1, the critic
 235 weight approximation error \tilde{W}_c in (37) is UUB with the adaptive critic learning (41).

236 *Proof.* Based on (37) and (41), the dynamics of \tilde{W}_c can be expressed as

$$\dot{\tilde{W}}_c(t) = -N_1 \tilde{W}_c(t) + N_2 \tag{42}$$

237 where

$$N_1 = \alpha_c \left(\frac{\sigma(t) \sigma(t)^T}{[1 + \sigma^T(t) \sigma(t)]^2} + \sum_{i=1}^k \frac{\sigma(t_i) \sigma(t_i)^T}{[1 + \sigma^T(t_i) \sigma(t_i)]^2} \right) \tag{43}$$

$$N_2 = \alpha_c \left(\frac{\sigma(t) \xi(t)}{[1 + \sigma^T(t) \sigma(t)]^2} + \sum_{i=1}^k \frac{\sigma(t_i) \xi(t_i)}{[1 + \sigma^T(t_i) \sigma(t_i)]^2} \right) \tag{44}$$

238 Note the fact that $\left\| \frac{y}{1+y^T y} \right\| \leq \frac{1}{2}$ and $\left\| \frac{1}{1+y^T y} \right\| \leq 1$ for arbitrary vector signal y . Then, from Remark
 239 4, N_2 in (42) satisfies $\|N_2\| \leq \frac{\alpha_c}{2} (k+1) b_\xi$. Consider the following Lyapunov function:

$$V_c = \frac{1}{2\alpha_c} \tilde{W}_c^T(t) \tilde{W}_c(t) \tag{45}$$

240 By differentiating (45) along the critic weight error dynamics (42), one has

$$\dot{V}_c = -\tilde{W}_c^T(t) \left(\frac{\sigma(t) \sigma^T(t)}{[1 + \sigma^T(t) \sigma(t)]^2} + \Lambda \right) \tilde{W}_c(t) + \tilde{W}_c^T(t) N_2 \tag{46}$$

241 with

$$\Lambda = \sum_{i=1}^k \frac{\sigma(t_i) \sigma^T(t_i)}{(1 + \sigma^T(t_i) \sigma(t_i))^2} > 0 \tag{47}$$

which is guaranteed by Condition 1. Therefore, \dot{V}_c is negative if

$$\left\| \tilde{W}_c(t) \right\| > \frac{\alpha_c (k+1) b_\xi}{2\lambda_{\min}(\Lambda)}$$

Then, the critic weight error \tilde{W}_c converges to the residual set

$$\Omega_c = \left\{ \tilde{W}_c \left\| \tilde{W}_c(t) \right\| > \frac{\alpha_c(k+1)b_\xi}{2\lambda_{\min}(\Lambda)} \right\}$$

242 Therefore, the critic weight error \tilde{W}_c is UUB. This completes the proof. ■

243 *Remark 5.* In contrast to the stability analysis as in [37, 55] where the PE condition on the signal
 244 is required for the signal $\sigma(t)$, in this paper, only Condition 1 is required to be satisfied for the
 245 signal $\sigma(t_i)$ in the history stack. Note that Condition 1 is weaker than the traditional PE condition
 246 and is easier to be checked for online implementation.

247 4.3. Actor and Disturbance Learning

248 As shown in (23) and (24), the optimal control policy and disturbance depend on the optimal
 249 value gradient $\frac{\partial V^*(s)}{\partial s}$. Therefore, consider the value gradient approximation with the critic weight
 250 \hat{W}_c in (35), the control and disturbance policies can be determined using the critic weight as

$$u_c(s) = -\lambda \tanh(\hat{D}_c) \quad (48)$$

$$\hat{D}_c = \frac{1}{2\lambda} R^{-1} G^T (\nabla \phi)^T \hat{W}_c \quad (49)$$

$$d_c(s) = \frac{1}{2\gamma^2} K^T (\nabla \phi)^T \hat{W}_c \quad (50)$$

251 However, this policy improvement does not guarantee the stability of the closed-loop system [36,
 252 37, 47, 55]. Therefore, to ensure the closed-loop stability, the policy applied to the system is
 253 implemented by alternative approximators using actor and disturbance network as

$$u_a(s) = -\lambda \tanh(\hat{D}_u) \quad (51)$$

$$\hat{D}_u = \frac{1}{2\lambda} R^{-1} G^T (\nabla \phi)^T \hat{W}_u \quad (52)$$

$$d_a(s) = \frac{1}{2\gamma^2} K^T (\nabla \phi)^T \hat{W}_d \quad (53)$$

254 where \hat{W}_u is the actor network weight and \hat{W}_d is the disturbance network weight. Define the weight
 255 estimation errors for the actor and the disturbance as,

$$\tilde{W}_u = W^* - \hat{W}_u, \quad \tilde{W}_d = W^* - \hat{W}_d \quad (54)$$

256 The actor network is designed to minimize the objective function

$$E_u = \frac{1}{2} e_u^T R e_u \quad (55)$$

257 where

$$e_u = u_a - u_c = \lambda [\tanh(D_c) - \tanh(D_a)] \quad (56)$$

258 denotes the difference between the actor u_a (51) applied to the system and the control input
 259 u_c (48). Applying the actor (51) and disturbance (53) to the system (14) yields the closed-loop
 260 dynamics

$$\begin{aligned} \dot{s}(t) &= \sigma_a(t) \\ &= F(s) - G(s) \lambda \tanh(\hat{D}_u) + \frac{1}{2\gamma^2} K(s) K(s)^T [\nabla\phi(s)]^T \hat{W}_d \end{aligned} \quad (57)$$

261 Define

$$\begin{aligned} \xi_1 &= \left[\frac{\sigma_a \sigma_a^T}{(1 + \sigma_a^T \sigma_a)^2} + \sum_{i=1}^k \frac{\sigma_{ai} \sigma_{ai}^T}{(1 + \sigma_{ai}^T \sigma_{ai})^2} \right] \\ \xi_2 &= \left[\frac{\sigma_a}{(1 + \sigma_a^T \sigma_a)^2} + \sum_{i=1}^k \frac{\sigma_{ai}}{(1 + \sigma_{ai}^T \sigma_{ai})^2} \right] \\ \xi_3 &= \frac{\sigma_a \pi}{(1 + \sigma_a^T \sigma_a)^2} + \sum_{i=1}^k \frac{\sigma_{ai} \pi_i}{(1 + \sigma_{ai}^T \sigma_{ai})^2} \\ \xi_4 &= -\frac{\alpha_c}{4\gamma^2} \left[\frac{\sigma_a}{(1 + \sigma_a^T \sigma_a)^2} + \sum_{i=1}^k \frac{\sigma_{ai}}{(1 + \sigma_{ai}^T \sigma_{ai})^2} \right] \\ \psi(t) &= \nabla\phi G \lambda \left[\tanh\left(\frac{\hat{D}_u}{\rho}\right) - \tanh(\hat{D}_u) \right] \\ \pi(t) &= W^T \nabla\phi G \lambda \left[\tanh\left(\frac{D_u^*}{\rho}\right) - \tanh\left(\frac{\hat{D}_u}{\rho}\right) \right] + \varepsilon_J \end{aligned} \quad (58)$$

262 where $\sigma_{ai} = \sigma_a(t_i)$ and $\pi_i = \pi(t_i)$. Then, the stability and convergence of all the signals in
 263 the closed-loop system with the barrier-actor-disturbance learning algorithm is discussed in the
 264 following theorem.

265 **Theorem 2.** Consider the dynamical system (14) with the critic (34), the actor (51), the distur-
 266 bance input (53) with the design parameters in (58) and the following adaptive learning rules for
 267 the critic weight \hat{W}_c , actor weight \hat{W}_u and disturbance \hat{W}_d , respectively,

$$\begin{aligned} \dot{\hat{W}}_c &= -\alpha_c \frac{\sigma_a(t) \left[U(s(t), u_a, d_a) + \hat{W}_c^T \sigma_a(t) \right]}{(1 + \sigma_a^T(t) \sigma_a(t))^2} \\ &\quad - \alpha_c \sum_{i=1}^k \frac{\sigma_a(t_i) \left[U(s(t_i), u_a(t_i), d_a(t_i)) + \hat{W}_c^T(t) \sigma_a(t_i) \right]}{(1 + \sigma_a^T(t_i) \sigma_a(t_i))^2} \end{aligned} \quad (59)$$

$$\dot{\hat{W}}_u = -\alpha_u \left[Y_u \hat{W}_u + \nabla\phi G e_u + \nabla\phi G \tanh^2(\hat{D}_u) e_u \right], \quad (60)$$

$$\dot{\hat{W}}_d = -\alpha_d \left(Y_{d1} \hat{W}_d - Y_{d2} \hat{W}_c + D_d \hat{W}_d \xi_4^T \hat{W}_c \right) \quad (61)$$

268 where $\alpha_c \in \mathbb{R}$, $\alpha_u \in \mathbb{R}$ and $\alpha_d \in \mathbb{R}$ are the learning rate for the critic, actor and disturbance
 269 networks, $Y_u \in \mathbb{R}^{N \times N}$, $Y_{d1} \in \mathbb{R}^{N \times N}$ and $Y_{d2} \in \mathbb{R}^{N \times N}$ are the feedback gains for the actor and
 270 disturbance networks. Then, the augmented state $X = \left[s^T \quad \tilde{W}_c^T \quad \tilde{W}_u^T \quad \tilde{W}_d^T \right]^T$ is UUB provided
 271 that the design parameters are selected such that

$$\begin{aligned} q &> 0 \\ -\xi_1 + \frac{r_c}{2}\xi_2\xi_2^T + \frac{1}{2r_{d1}}I + \frac{r_{d2}}{2}\xi_4\xi_4^T &< 0 \\ \frac{1}{2r_c}\psi\psi^T - Y_u &< 0 \\ Y_{d1} + D_d\xi_4^T W^* + \frac{r_{d1}}{2}Y_{d2}Y_{d2}^T + \frac{1}{2r_{d2}}D_d W^* (W^*)^T D_d^T &< 0 \end{aligned} \quad (62)$$

272 where r_c , r_{d1} and r_{d2} are positive constants to be determined.

273 *Proof.* Consider the following Lyapunov candidate function:

$$J(X) = V^*(s) + V_c(\tilde{W}_c) + V_u(\tilde{W}_u) + V_d(\tilde{W}_d) \quad (63)$$

274 where $V^*(\cdot)$ is the optimal value function satisfying the HJI equation and

$$V_c(s) = \frac{1}{2}\tilde{W}_c^T \alpha_c^{-1} \tilde{W}_c, \quad V_u(s) = \frac{1}{2}\tilde{W}_u^T \alpha_u^{-1} \tilde{W}_u, \quad V_d(s) = \frac{1}{2}\tilde{W}_d^T \alpha_d^{-1} \tilde{W}_d$$

275 The derivative of the Lyapunov function (63) is given by

$$\dot{J} = \dot{V}^* + \dot{V}_c + \dot{V}_u + \dot{V}_d \quad (64)$$

276 For the first term of (64), one has

$$\begin{aligned} \dot{V}^* &= \left[(W^*)^T \nabla \phi + (\nabla \varepsilon)^T \right] [F(s) + G(s)u_a + K(s)d_a] \\ &= (W^*)^T \nabla \phi F - (W^*)^T \nabla \phi G \lambda \tanh(\hat{D}_u) + \frac{1}{2\gamma^2} (W^*)^T D_d \hat{W}_d + \varepsilon_0 \end{aligned} \quad (65)$$

277 with D_d is defined in (33) and

$$\begin{aligned} \varepsilon_0 &= (\nabla \varepsilon)^T \sigma_a \\ \sigma_a &= F(s) - G(s) \lambda \tanh(\hat{D}_u) + \frac{1}{2\gamma^2} K(s) K(s)^T [\nabla \phi(s)]^T \hat{W}_d \end{aligned} \quad (66)$$

278 Based on Assumptions 1 and 3 and Remark 4, ε_0 can be upper bounded as

$$\varepsilon_0 \leq b_{d\varepsilon} b_f \|s\| + b_{d\varepsilon} b_g \lambda + \frac{1}{2\gamma^2} b_{d\varepsilon} b_k^2 b_{d\phi} b_* - \frac{1}{2\gamma^2} (\nabla \varepsilon)^T K(s) K(s)^T [\nabla \phi(s)]^T \tilde{W}_d \quad (67)$$

279 From (25) and (32), one has

$$\begin{aligned} (W^*)^T \nabla \phi F &= -Q(s) - \Theta(-\lambda \tanh(D_u^*)) + (W^*)^T \nabla \phi G \lambda \tanh(D_u^*) \\ &\quad - \frac{1}{4\gamma^2} (W^*)^T D_d W^* + \zeta \end{aligned}$$

280 with

$$\Theta(-\lambda \tanh(D_u^*)) = (W^*)^T \nabla \phi G \lambda \tanh(D_u^*) + \lambda^2 \bar{R} \ln(1 - \tanh^2(D_u^*)) \quad (68)$$

281 where D_u has been defined as in (33). Inserting $(W^*)^T \nabla \phi F$ and (68) into (65) yields

$$\begin{aligned} \dot{V}^* &= -Q(s) - \Theta(-\lambda \tanh(D_u^*)) + (W^*)^T \nabla \phi G \lambda \tanh(D_u^*) \\ &\quad - \frac{1}{4\gamma^2} (W^*)^T D_d W^* - (W^*)^T \nabla \phi G \lambda \tanh(\hat{D}_u) + \frac{1}{2\gamma^2} (W^*)^T D_d W^* \\ &\quad - \frac{1}{2\gamma^2} (W^*)^T D_d \tilde{W}_d + \zeta + \varepsilon_0 \end{aligned} \quad (69)$$

282 Since $Q(\cdot)$ and $\Theta(\cdot)$ are positive definite functions, then, there exists a positive constant $q > 0$ such
283 that

$$s^T q s \leq Q(s) \leq Q(s) + \Theta(-\lambda \tanh(D_u^*)) \quad (70)$$

284 The third term in (69) can be upper bounded by

$$(W^*)^T \nabla \phi(x) G \lambda \tanh(D_u^*) \leq \lambda b_g b_{d\phi} b_* \quad (71)$$

285 Considering $W^* = W_u + \tilde{W}_u$, then, the fourth term in (69) can be rewritten as

$$\begin{aligned} &-(W^*)^T \nabla \phi G \lambda \tanh(\hat{D}_u) \\ &= -\tilde{W}_u^T \nabla \phi G \lambda \tanh(\hat{D}_u) - \hat{W}_u^T \nabla \phi G \lambda \tanh(\hat{D}_u) \\ &\leq -\tilde{W}_u^T \nabla \phi G \lambda \tanh(\hat{D}_u) \end{aligned} \quad (72)$$

286 where the above inequality results from the fact that $\hat{W}_u^T \nabla \phi G \lambda \tanh(\hat{D}_u) = 2\lambda^2 \bar{R} [\hat{D}_u \tanh(\hat{D}_u)]$
287 and $x^T \tanh(x) \geq 0$, for arbitrary vector signal x . Considering now the facts (67), (69), (70), (71)
288 and (72), \dot{V}^* further satisfies

$$\begin{aligned} \dot{V}^* &\leq -\tilde{W}_u^T \nabla \phi G \lambda \tanh(\hat{D}_u) - s^T q s + b_{d\varepsilon} b_f \|s\| \\ &\quad + \lambda b_g b_{d\phi} b_* + \frac{1}{4\gamma^2} b_*^2 b_{d\phi}^2 b_k^2 + b_\zeta + \lambda b_{d\varepsilon} b_g + \frac{1}{2\gamma^2} b_{d\varepsilon} b_k^2 b_{d\phi} b_* \\ &\quad - \frac{1}{2\gamma^2} (\nabla \varepsilon)^T K(s) K(s)^T [\nabla \phi(s)]^T \tilde{W}_d - \frac{1}{2\gamma^2} (W^*)^T D_d \tilde{W}_d \\ &= -\tilde{W}_u^T \nabla \phi G \lambda \tanh(\hat{D}_u) - s^T q s + M_s \|s\| + N_s + M_{d1} \tilde{W}_d \end{aligned} \quad (73)$$

289 where

$$\begin{aligned} M_s &= b_{d\varepsilon} b_f \\ N_s &= \lambda b_g b_{d\phi} b_* + \frac{1}{4\gamma^2} b_*^2 b_{d\phi}^2 b_k^2 + b_\zeta + \lambda b_{d\varepsilon} b_g + \frac{1}{2\gamma^2} b_{d\varepsilon} b_k^2 b_{d\phi} b_* \\ M_{d1} &= -\frac{1}{2\gamma^2} \left\{ (\nabla \varepsilon)^T K(s) K(s)^T [\nabla \phi(s)]^T + (W^*)^T D_d \right\} \end{aligned} \quad (74)$$

290 Second, for the critic weight error \tilde{W}_c , from (41) one has

$$\dot{\tilde{W}}_c(t) = \alpha_c \frac{\sigma_a}{(1 + \sigma_a^T \sigma_a)^2} e_c(t) + \alpha_c \sum_{i=1}^k \frac{\sigma_{ai}}{(1 + \sigma_{ai}^T \sigma_{ai})^2} e_c(t_i, t) \quad (75)$$

291 where σ_a has been defined in (66) and $\sigma_{ai} = \sigma_a(t_i)$. Differentiating V_c along with (75), one has

$$\begin{aligned} \dot{V}_c &= \tilde{W}_c^T \alpha_c^{-1} \dot{\tilde{W}}_c \\ &= \tilde{W}_c^T \left[\frac{\sigma_a}{(1 + \sigma_a^T \sigma_a)^2} e_c(t) + \sum_{i=1}^k \frac{\sigma_{ai}}{(1 + \sigma_{ai}^T \sigma_{ai})^2} e_c(t_i, t) \right] \end{aligned} \quad (76)$$

292 From (32), one has

$$-Q(s) - \Theta(-\lambda \tanh(D_u^*)) - (W^*)^T \sigma + \frac{1}{4\gamma^2} (W^*)^T D_d W^* + \zeta = 0. \quad (77)$$

293 Therefore, one can obtain

$$\begin{aligned} e_c &= Q(s) + \Theta(-\lambda \tanh(\hat{D}_u)) + \hat{W}_c^T \sigma_a - \frac{1}{4\gamma^2} \hat{W}_d^T D_d \hat{W}_d \\ &= Q(s) + \Theta(-\lambda \tanh(\hat{D}_u)) + \hat{W}_c^T \sigma_a - \frac{1}{4\gamma^2} \hat{W}_d^T D_d \hat{W}_d \\ &\quad - Q(s) - \Theta(-\lambda \tanh(D_u^*)) - (W^*)^T \sigma + \frac{1}{4\gamma^2} (W^*)^T D_d W^* + \zeta \end{aligned}$$

294 Adding and subtracting $(W^*)^T \sigma_a$ to e_c yields

$$\begin{aligned} e_c &= \Theta(-\lambda \tanh(\hat{D}_u)) - \Theta(-\lambda \tanh(D_u^*)) - \tilde{W}_c^T \sigma_a + (W^*)^T (\sigma_a - \sigma) \\ &\quad - \frac{1}{4\gamma^2} \hat{W}_d^T D_d \hat{W}_d + \frac{1}{4\gamma^2} (W^*)^T D_d W^* + \zeta \end{aligned} \quad (78)$$

295 Moreover, note that

$$\begin{aligned} &\Theta(-\lambda \tanh(\hat{D}_u)) - \Theta(-\lambda \tanh(D_u^*)) \\ &= \lambda \hat{W}_a^T \nabla \phi G \tanh(\hat{D}_u) + \lambda^2 \bar{R} \ln(1 - \tanh^2(\hat{D}_u)) \\ &\quad - \lambda W^T \nabla \phi G \tanh(D_u^*) - \lambda^2 \bar{R} \ln(1 - \tanh^2(D_u^*)) \end{aligned} \quad (79)$$

296 Note that the term $\lambda^2 \bar{R} \ln(1 - \tanh^2(D_u^*))$ in (79) can be rewritten as

$$\lambda^2 \bar{R} \ln(1 - \tanh^2(D_u^*)) = \lambda^2 \bar{R} \left[\ln 4 - 2D_u^* - 2 \ln(1 + e^{-2D_u^*}) \right], \quad (80)$$

297 where $-2 \ln(1 + e^{-2D_u^*})$ can be approximated using Lemma 1 as

$$-2 \ln(1 + e^{-2D_u^*}) = 2D_u^* - 2D_u^* \operatorname{sgn}(D_u^*) + \varepsilon_{D_u^*}, \quad (81)$$

298 where $\|\varepsilon_{D_u^*}\| \leq \ln 4$. Then, inserting (81) into (80) yields

$$\lambda^2 R \ln(1 - \tanh^2(D_u^*)) = \lambda^2 \bar{R} [\ln 4 - 2D_u^* \operatorname{sgn}(D_u^*) + \varepsilon_{D_u^*}]. \quad (82)$$

299 Similarly,

$$\lambda^2 \bar{R} \ln(1 - \tanh^2(\hat{D}_u)) = \lambda^2 R [\ln 4 - 2\hat{D}_u \operatorname{sgn}(\hat{D}_u) + \varepsilon_{\hat{D}_u}], \quad (83)$$

300 where $\|\varepsilon_{\hat{D}_u}\| \leq \ln 4$. Consider (79), (82) and (83), one has

$$\begin{aligned} & \Theta(-\lambda \tanh(\hat{D}_u)) - \Theta(-\lambda \tanh(D_u^*)) \\ &= \lambda \hat{W}_a^T \nabla \phi G \tanh(\hat{D}_u) - \lambda W^T \nabla \phi G \tanh(D_u^*) \\ & \quad + \lambda^2 \bar{R} [2D_u^* \operatorname{sgn}(D_u^*) - 2\hat{D}_u \operatorname{sgn}(\hat{D}_u) + \varepsilon_{\hat{D}_u} - \varepsilon_{D_u^*}] \end{aligned} \quad (84)$$

301 The nonsmooth function $\operatorname{sgn}(\cdot)$ in (84) can be approximated by the function $\tanh(\cdot)$ by using
302 Lemma 2. Then, based on (84), one has

$$\begin{aligned} & \lambda^2 \bar{R} (2D_u^* \operatorname{sgn}(D_u^*) - 2\hat{D}_u \operatorname{sgn}(\hat{D}_u)) \\ &= (W^*)^T \nabla \phi G \lambda \tanh\left(\frac{D_u^*}{\rho}\right) - \hat{W}_u^T \nabla \phi G \lambda \tanh\left(\frac{\hat{D}_u}{\rho}\right) + \lambda^2 \bar{R} \varepsilon_\rho \end{aligned} \quad (85)$$

303 with approximation error satisfying $0 \leq \varepsilon_\rho \leq 2\kappa\rho$ where $\kappa = 0.2785$ is defined in Lemma 2. Based
304 on (84) and (85), adding and subtracting $(W^*)^T \nabla \phi G \lambda \tanh(\hat{D}_u)$, one has

$$\begin{aligned} e_c &= -\tilde{W}_c^T \sigma_a + \tilde{W}_u^T \nabla \phi G \lambda \left[\tanh\left(\frac{\hat{D}_u}{\rho}\right) - \tanh(\hat{D}_u) \right] \\ & \quad + (W^*)^T \nabla \phi G \lambda \left[\tanh\left(\frac{D_u^*}{\rho}\right) - \tanh\left(\frac{\hat{D}_u}{\rho}\right) \right] + \zeta + \lambda^2 \bar{R} (\varepsilon_{\hat{D}_u} - \varepsilon_{D_u^*} + \varepsilon_\rho) + \epsilon_c \end{aligned} \quad (86)$$

305 where $\epsilon_c = -\frac{1}{4\gamma^2} \hat{W}_d^T D_d \hat{W}_d + \frac{1}{2\gamma^2} (W^*)^T D_d \hat{W}_d - \frac{1}{4\gamma^2} (W^*)^T D_d W^*$, which can be further rewritten
306 as

$$\begin{aligned} \epsilon_c &= -\frac{1}{4\gamma^2} \hat{W}_d^T D_d \hat{W}_d - \frac{1}{4\gamma^2} (W^*)^T D_d W^* + \frac{1}{4\gamma^2} (W^*)^T D_d \hat{W}_d + \frac{1}{4\gamma^2} (W^*)^T D_d \hat{W}_d \\ &= \frac{1}{4\gamma^2} \tilde{W}_d^T D_d \hat{W}_d - \frac{1}{4\gamma^2} (W^*)^T D_d \tilde{W}_d \\ &= -\frac{1}{4\gamma^2} \tilde{W}_d^T D_d \tilde{W}_d \end{aligned} \quad (87)$$

307 Denote $\varepsilon_J = \lambda^2 R (\varepsilon_{\hat{D}_u} - \varepsilon_{D_u^*} + \varepsilon_\rho) + \zeta$, one has

$$e_c(t) = -\tilde{W}_c^T(t) \sigma_a(t) + \tilde{W}_a^T(t) \psi(t) + \pi(t) - \frac{1}{4\gamma^2} \tilde{W}_d^T(t) D_d \tilde{W}_d(t) \quad (88)$$

308 where ψ and π is defined in (58). Similarly,

$$e_c(t_i, t) = -\tilde{W}_c^T(t) \sigma_a(t_i) + \tilde{W}_a^T(t) \psi(t) + \pi(t_i) - \frac{1}{4\gamma^2} \tilde{W}_d^T(t) D_d \tilde{W}_d(t) \quad (89)$$

309 where Based on Assumptions 1 and 3, both ψ and π are bounded. Substituting (88) and (89) into
310 (75) yields,

$$\begin{aligned} \dot{\tilde{W}}_c &= -\alpha_c \left[\frac{\sigma_a \sigma_a^T}{(1 + \sigma_a^T \sigma_a)^2} + \sum_{i=1}^k \frac{\sigma_{ai} \sigma_{ai}^T}{(1 + \sigma_{ai}^T \sigma_{ai})^2} \right] \tilde{W}_c \\ &+ \alpha_c \left[\frac{\sigma_a}{(1 + \sigma_a^T \sigma_a)^2} + \sum_{i=1}^k \frac{\sigma_{ai}}{(1 + \sigma_{ai}^T \sigma_{ai})^2} \right] \psi^T \tilde{W}_a \\ &+ \alpha_c \left[\frac{\sigma_a \pi}{(1 + \sigma_a^T \sigma_a)^2} + \sum_{i=1}^k \frac{\sigma_{ai} \pi_i}{(1 + \sigma_{ai}^T \sigma_{ai})^2} \right] \\ &- \frac{\alpha_c}{4\gamma^2} \left[\frac{\sigma_a}{(1 + \sigma_a^T \sigma_a)^2} + \sum_{i=1}^k \frac{\sigma_{ai}}{(1 + \sigma_{ai}^T \sigma_{ai})^2} \right] \tilde{W}_d^T D_d \tilde{W}_d \end{aligned} \quad (90)$$

311 Substituting (90) into (76), one has

$$\begin{aligned} \dot{V}_c &= \tilde{W}_c^T \alpha_c^{-1} \dot{\tilde{W}}_c \\ &= -\tilde{W}_c^T \xi_1 \tilde{W}_c + \tilde{W}_c^T \xi_2 \psi^T \tilde{W}_a + \tilde{W}_c^T \xi_3 + \tilde{W}_c^T \xi_4 \tilde{W}_d^T D_d \tilde{W}_d \\ &\leq -\tilde{W}_c^T \xi_1 \tilde{W}_c + \frac{r_c}{2} \tilde{W}_c^T \xi_2 \xi_2^T \tilde{W}_c + \frac{1}{2r_c} \tilde{W}_a^T \psi \psi^T \tilde{W}_a + \tilde{W}_c^T \xi_3 + \tilde{W}_c^T \xi_4 \tilde{W}_d^T D_d \tilde{W}_d \\ &= \tilde{W}_c^T \left[-\xi_1 + \frac{r_c}{2} \xi_2 \xi_2^T \right] \tilde{W}_c + \frac{1}{2r_c} \tilde{W}_a^T \psi \psi^T \tilde{W}_a + \tilde{W}_c^T \xi_3 + \tilde{W}_c^T \xi_4 \tilde{W}_d^T D_d \tilde{W}_d \end{aligned} \quad (91)$$

312 where ξ_i for $i = 1, 2, 3, 4$ has been defined in (58).

313 Next, we give the upper bound of \dot{V}_u . Based on (60), differentiating V_u yields

$$\begin{aligned} \dot{V}_u &= \tilde{W}_u^T \alpha_u^{-1} \dot{\tilde{W}}_u \\ &= -\tilde{W}_u^T \left[\nabla \phi G e_u + \nabla \phi G \tanh^2(\hat{D}_u) e_u + Y_u \hat{W}_u \right] \\ &= -\tilde{W}_u^T Y_u \tilde{W}_u + \tilde{W}_u^T \nabla \phi G \lambda \tanh(\hat{D}_u) + \tilde{W}_u^T M_u \\ &\leq -\tilde{W}_u^T Y_u \tilde{W}_u + \tilde{W}_u^T \nabla \phi G \lambda \tanh(\hat{D}_u) + M_u^T \tilde{W}_u \end{aligned} \quad (92)$$

314 where $M_u = \left[-\nabla \phi G \lambda \tanh(\hat{D}_c) + \nabla \phi G \tanh^2(\hat{D}_u) e_u + Y_u W^* \right]$

315 Based on Assumption 1 - 3 and the definition of the actor learning error e_u in (56), M_u is also
316 bounded.

317 For the derivative of V_d , according to (61) one has

$$\begin{aligned}
 \dot{V}_d &= -\tilde{W}_d^T Y_{d1} W^* + \tilde{W}_d^T Y_{d1} \tilde{W}_d + \tilde{W}_d^T Y_{d2} W^* - \tilde{W}_d^T Y_{d2} \tilde{W}_c \\
 &\quad - \tilde{W}_d^T D_d W^* \xi_4^T W^* + \tilde{W}_d^T D_d \tilde{W}_d \xi_4^T W^* + \tilde{W}_d^T D_d W^* \xi_4^T \tilde{W}_c - \tilde{W}_d^T D_d \tilde{W}_d \xi_4^T \tilde{W}_c \\
 &= \tilde{W}_d^T Y_{d1} \tilde{W}_d + \tilde{W}_d^T D_d \tilde{W}_d \xi_4^T W^* - \tilde{W}_d^T Y_{d1} W^* - \tilde{W}_d^T D_d W^* \xi_4^T W^* + \tilde{W}_d^T Y_{d2} W^* \\
 &\quad - \tilde{W}_d^T Y_{d2} \tilde{W}_c + \tilde{W}_d^T D_d W^* \xi_4^T \tilde{W}_c - \tilde{W}_d^T D_d \tilde{W}_d \xi_4^T \tilde{W}_c
 \end{aligned} \tag{93}$$

318 Using Young's inequality to (93) yields

$$\begin{aligned}
 \dot{V}_d &\leq \tilde{W}_d^T (Y_{d1} + D_d \xi_4^T W^*) \tilde{W}_d + \tilde{W}_d^T [Y_{d2} W^* - Y_{d1} W^* - D_d W^* \xi_4^T W^*] \\
 &\quad + \frac{r_{d1}}{2} \tilde{W}_d^T Y_{d2} Y_{d2}^T \tilde{W}_d + \frac{1}{2r_{d1}} \|\tilde{W}_c\|^2 + \frac{1}{2r_{d2}} \tilde{W}_d^T D_d W^* (W^*)^T D_d^T \tilde{W}_d + \frac{r_{d2}}{2} \tilde{W}_c^T \xi_4 \xi_4^T \tilde{W}_c \\
 &\quad - \tilde{W}_d^T D_d \tilde{W}_d \xi_4^T \tilde{W}_c \\
 &= \frac{1}{2r_{d1}} \|\tilde{W}_c\|^2 + \frac{r_{d2}}{2} \tilde{W}_c^T \xi_4 \xi_4^T \tilde{W}_c - \tilde{W}_d^T Q_d \tilde{W}_d + M_{d2}^T \tilde{W}_d - \tilde{W}_d^T D_d \tilde{W}_d \xi_4^T \tilde{W}_c
 \end{aligned} \tag{94}$$

319 where

$$\begin{aligned}
 Q_d &= - \left[Y_{d1} + D_d \xi_4^T W^* + \frac{r_{d1}}{2} Y_{d2} Y_{d2}^T + \frac{1}{2r_{d2}} D_d W^* (W^*)^T D_d^T \right] \\
 M_{d2} &= Y_{d2} W^* - Y_{d1} W^* - D_d W^* \xi_4^T W^*
 \end{aligned} \tag{95}$$

320 Finally, collecting the results in (73), (91), (92) and (94), one has

$$\begin{aligned}
 \dot{J} &\leq -s^T Q_s s + M_s \|s\| + N_s - \tilde{W}_c^T Q_c \tilde{W}_c + \tilde{W}_c^T \xi_3 \\
 &\quad - \tilde{W}_u^T Q_u \tilde{W}_u + M_u^T \tilde{W}_u - \tilde{W}_d^T Q_d \tilde{W}_d + (M_{d1} + M_{d2}^T) \tilde{W}_d \\
 &\leq -\lambda_{\min}(Q_s) \|s\|^2 + \|M_s\| \|s\| + \|N_s\| - \lambda_{\min}(Q_c) \|\tilde{W}_c\|^2 + \|\xi_3\| \|\tilde{W}_c\| \\
 &\quad - \lambda_{\min}(Q_u) \|\tilde{W}_u\|^2 + \|M_u\| \|\tilde{W}_u\| - \lambda_{\min}(Q_d) \|\tilde{W}_d\|^2 + \|M_{d1} + M_{d2}^T\| \|\tilde{W}_d\|
 \end{aligned} \tag{96}$$

321 where

$$\begin{aligned}
 Q_s &= qI \\
 Q_c &= \xi_1 - \frac{r_c}{2} \xi_2 \xi_2^T - \frac{1}{2r_{d1}} I - \frac{r_{d2}}{2} \xi_4 \xi_4^T \\
 Q_u &= -\frac{1}{2r_c} \psi \psi^T + Y_u
 \end{aligned}$$

322 Based on Assumptions 1 and 3, M_s , M_u , M_{d1} , M_{d2} and N_s are bounded. Note that the parameters

323 design in (62) guarantees that $Q_s > 0$, $Q_c > 0$, $Q_u > 0$ and $Q_d > 0$. Therefore, $\dot{J} < 0$ if

$$\begin{aligned}
 \|s\| &> \frac{\|s\|}{2\lambda_{\min}(Q_s)} + \sqrt{\frac{\|M_s\|^2}{4\lambda_{\min}^2(Q_s)} + \frac{\|N_s\|}{\lambda_{\min}(Q_s)}} \\
 \|\tilde{W}_c\| &> \frac{\|\tilde{W}_c\|}{2\lambda_{\min}(Q_c)} + \sqrt{\frac{\|M_c\|^2}{4\lambda_{\min}^2(Q_c)}} \\
 \|\tilde{W}_u\| &> \frac{\|\tilde{W}_u\|}{2\lambda_{\min}(Q_u)} + \sqrt{\frac{\|M_u\|^2}{4\lambda_{\min}^2(Q_u)}} \\
 \|\tilde{W}_d\| &> \frac{\|\tilde{W}_d\|}{2\lambda_{\min}(Q_d)} + \sqrt{\frac{\|M_d\|^2}{4\lambda_{\min}^2(Q_d)}}
 \end{aligned} \tag{97}$$

324 Then, the augmented state $X = [s^T \ \tilde{W}_c^T \ \tilde{W}_u^T \ \tilde{W}_d^T]^T$ converges to the residual set Ω_X
 325 defined as

$$\begin{aligned}
 \Omega_X = \{X \mid \|s\| < \frac{\|s\|}{2\lambda_{\min}(Q_s)} + \sqrt{\frac{\|M_s\|^2}{4\lambda_{\min}^2(Q_s)} + \frac{\|N_s\|}{\lambda_{\min}(Q_s)}}, \\
 \|\tilde{W}_c\| < \frac{\|\tilde{W}_c\|}{2\lambda_{\min}(Q_c)} + \sqrt{\frac{\|M_c\|^2}{4\lambda_{\min}^2(Q_c)}}, \|\tilde{W}_u\| < \frac{\|\tilde{W}_u\|}{2\lambda_{\min}(Q_u)} + \sqrt{\frac{\|M_u\|^2}{4\lambda_{\min}^2(Q_u)}}, \\
 \|\tilde{W}_d\| < \frac{\|\tilde{W}_d\|}{2\lambda_{\min}(Q_d)} + \sqrt{\frac{\|M_d\|^2}{4\lambda_{\min}^2(Q_d)}} \}
 \end{aligned} \tag{98}$$

326 This completes the proof. ■

327 5. Simulation Study

328 To verify the effectiveness of the presented online safe RL algorithm with the actor-critic-barrier
 329 structure, we consider the following nonlinear systems of a single link robot arm

$$\ddot{\theta}(t) = -\frac{Mgl}{\tilde{G}} \sin(\theta(t)) - \frac{\tilde{D}}{\tilde{G}} \dot{\theta}(t) + \frac{1}{\tilde{G}} u(t) + kd(t) \tag{99}$$

330 where θ is the angle position, $\dot{\theta}$ is the angle velocity, M is the mass of the payload, g is the
 331 acceleration of gravity, l is the length of the arm, \tilde{D} is the viscous friction and \tilde{G} is the moment of
 332 inertia. In this experiment, $M = 10\text{kg}$, $g = 9.81\text{m/s}^2$, $l = 0.5\text{m}$, $\tilde{D} = 2\text{N}$ and $\tilde{G} = 10\text{kgm}^2$. Let
 333 $x_1 = \theta$, $x_2 = \dot{\theta}$ and $x = [x_1 \ x_2]^T$, then the dynamics of x can be written as

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} x_2 \\ f(x) \end{bmatrix} + \begin{bmatrix} 0 \\ g(x) \end{bmatrix} u + \begin{bmatrix} 0 \\ k(x) \end{bmatrix} d \tag{100}$$

334 where

$$\begin{aligned} f(x) &= -\frac{Mgl}{\tilde{G}} \sin(x_1) - \frac{\tilde{D}}{\tilde{G}} x_2 \\ g(x) &= \frac{1}{\tilde{G}}, k(x) = k \end{aligned}$$

335 For Problem 1, the performance output is selected as $L(x, u) = x^T H x + u^T R u$ with $H = 50I$,
336 $R = 10I$. In addition, the following safety constraints are considered

$$x_i \in (a_i, A_i), \forall i \in \{1, 2\} \quad (101)$$

337 where $a_1 = -1.6$, $A_1 = 3$, $a_2 = -4$ and $A_2 = 3$. By using the classical actor-critic reinforcement
338 learning algorithm, the state evolution with respect to time can be shown in Figure 3. The phase
339 portrait of the state evolution in the state space is shown in Figure 5. As can be seen from Figure 3,
340 the full-state constraints cannot be guaranteed by the classical actor-critic reinforcement learning
341 algorithm. The evolution of the actor-critic-disturbance is shown in Figure 4.

342 To deal with the full-state constraints, the barrier-function-based system transformation (10)
343 is employed. With the barrier function, one can obtain the transformed system as $\dot{s} = F(s) +$
344 $G(s)u + K(s)d$ with

$$\begin{aligned} F(s) &= \begin{bmatrix} \frac{a_2 A_2 (e^{\frac{s_2^2}{2}} - e^{-\frac{s_2^2}{2}})}{a_2 e^{\frac{s_2^2}{2}} - A_2 e^{-\frac{s_2^2}{2}}} \frac{A_1^2 e^{-s_1} - 2a_1 A_1 + a_1^2 e^{s_1}}{A_1 a_1^2 - a_1 A_1^2} \\ f(B^{-1}(s)) \frac{A_2^2 e^{-s_2} - 2a_2 A_2 + a_2^2 e^{s_2}}{A_2 a_2^2 - a_2 A_2^2} \end{bmatrix} \\ G(s) &= \begin{bmatrix} 0 \\ \frac{1}{\tilde{G}} \frac{A_2^2 e^{-s_2} - 2a_2 A_2 + a_2^2 e^{s_2}}{A_2 a_2^2 - a_2 A_2^2} \end{bmatrix} \\ K(s) &= \begin{bmatrix} 0 \\ k \frac{A_2^2 e^{-s_2} - 2a_2 A_2 + a_2^2 e^{s_2}}{A_2 a_2^2 - a_2 A_2^2} \end{bmatrix} \end{aligned} \quad (102)$$

345 with the initial condition

$$\begin{aligned} s_0 &= \begin{bmatrix} s_0(1) & s_0(2) \end{bmatrix}^T \\ s_0(1) &= b(x_0(1); a_1, A_1), s_0(2) = b(x_0(2); a_2, A_2) \end{aligned}$$

Based on the actor-critic-barrier online learning algorithm, the state evolution of state $s(t)$ in
system (102) is given in Figure 6. One can observe that the state $s(t)$ of system (102) converges
to the origin asymptotically. Based on the state evolution of $s(t)$, by using the barrier function
inverse mapping (10), one can obtain the state $x(t)$ as

$$x_1(t) = b^{-1}(s_1(t); a_1, A_1), x_2(t) = b^{-1}(s_2(t); a_2, A_2)$$

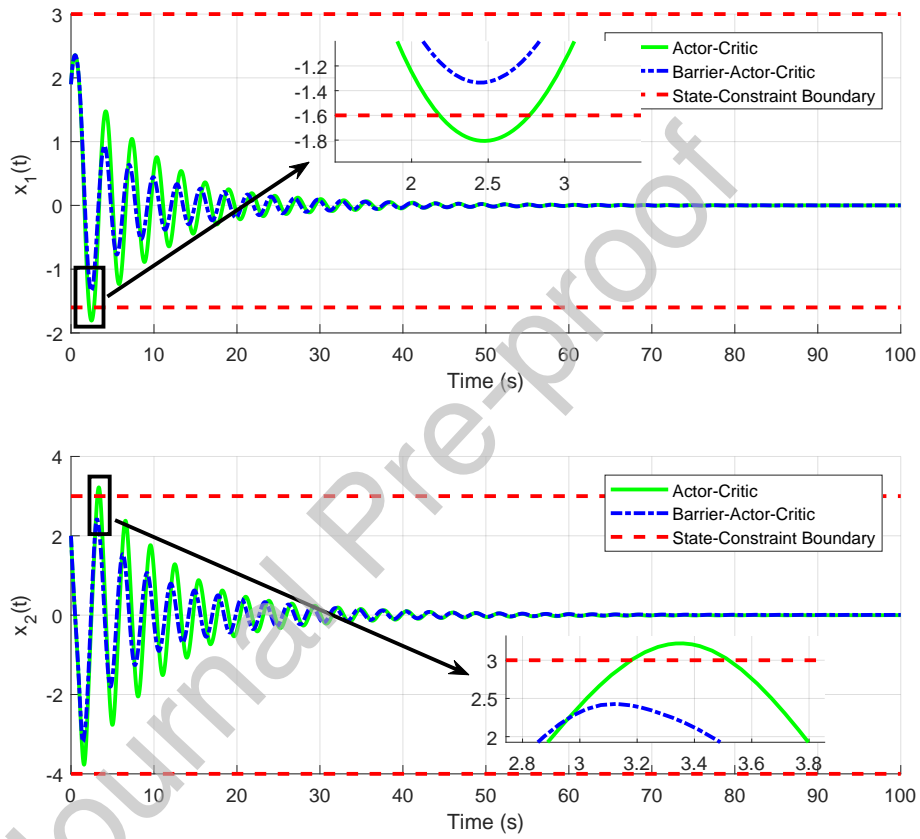


Figure 3: Evolution of the state $x(t)$ by using the presented actor-critic-barrier learning and classical actor-critic learning. The dashed line represents the boundary of the safe region.

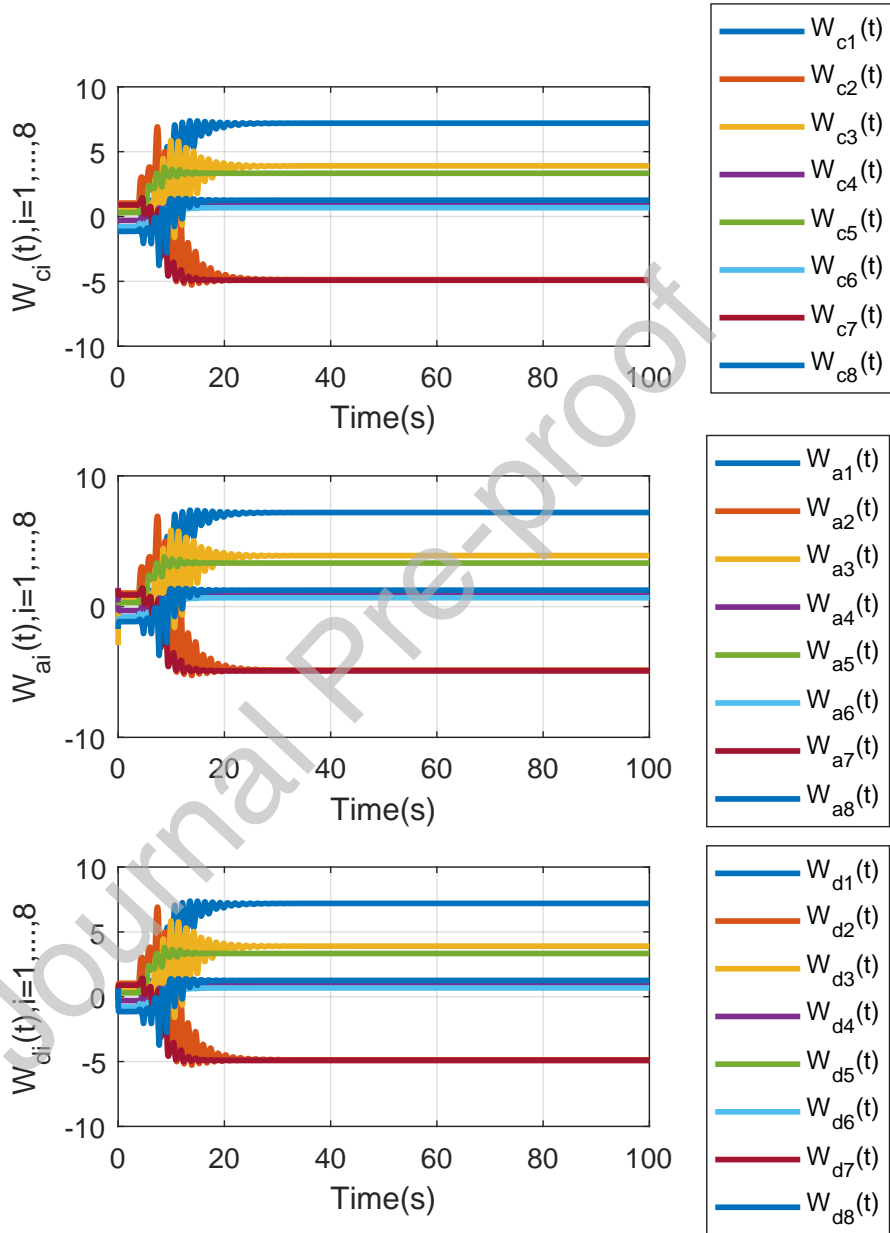


Figure 4: Evolution of the actor and critic weights using classical actor-critic learning.

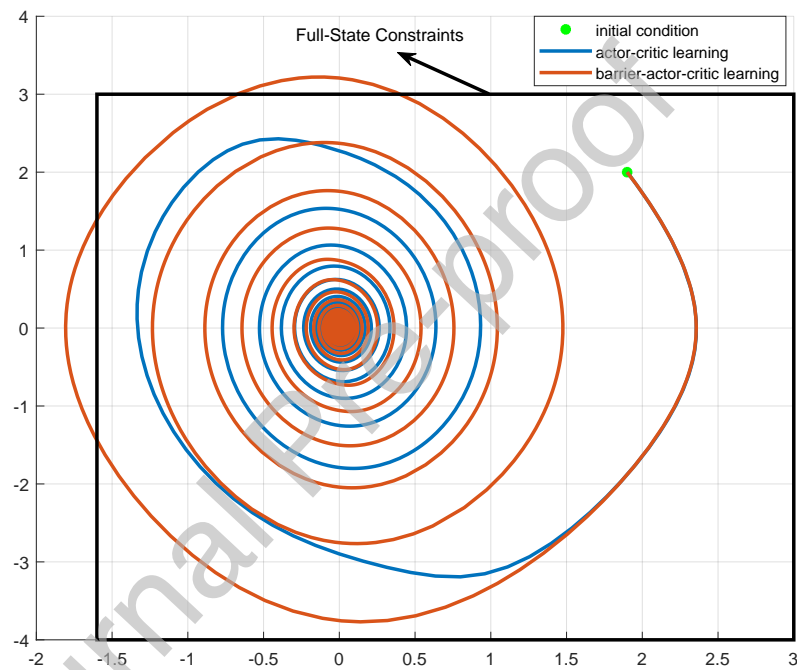


Figure 5: Evolution of the two-dimensional phase plot of the state trajectories $[x_1(t) \ x_2(t)]$. The black box denotes the safe region.

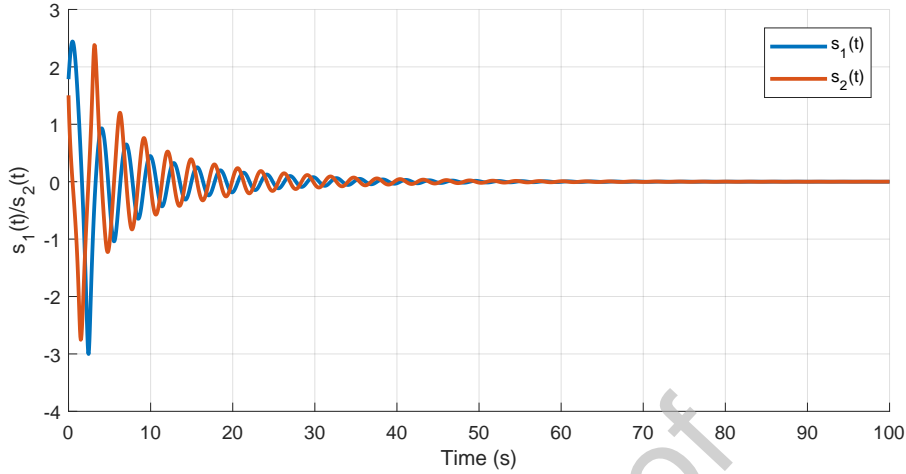


Figure 6: Evolution of the state $s(t)$ by using the presented actor-critic-barrier learning and classical actor-critic learning.

346 Then, the evolution of the state $x(t)$ is shown in Figure 3. The phase portrait of the state evolution
 347 $[x_1(t) \ x_2(t)]$ is provided in Figure 5. The black box represents the full-state constraints. One can
 348 observe that with the barrier-actor-critic learning algorithm, the state evolution does not exceed the
 349 boundary of the prescribed region and full-state constraints can be guaranteed. That is, the state
 350 $x(t)$ converges to the origin asymptotically while satisfying the safety constraints (101). Finally,
 351 the learning process of the barrier-actor-critic networks is shown in Figure 7.

352 6. Conclusions

353 In this paper, the disturbance attenuation problem with both full-state constraints and input
 354 saturation is considered. An adaptive optimal controller design with the barrier-actor-critic al-
 355 gorithm is developed. First, a novel barrier function is defined to deal with full-state saturation.
 356 Based on this barrier function, a novel system transformation is applied to the original system
 357 to obtain the transformed system. Second, the barrier-function-based system transformation is
 358 then combined with the actor-critic online algorithm to learn the optimal control policy and the
 359 worst-case disturbance. To obviate the requirement of PE condition for online critic learning, the
 360 experience replay technique is employed to utilize the online and history data concurrently. The
 361 stability of the closed-loop system and the convergence of the actor-critic parameters to the op-
 362 timal condition are discussed in the framework of Lyapunov analysis. The input saturation and

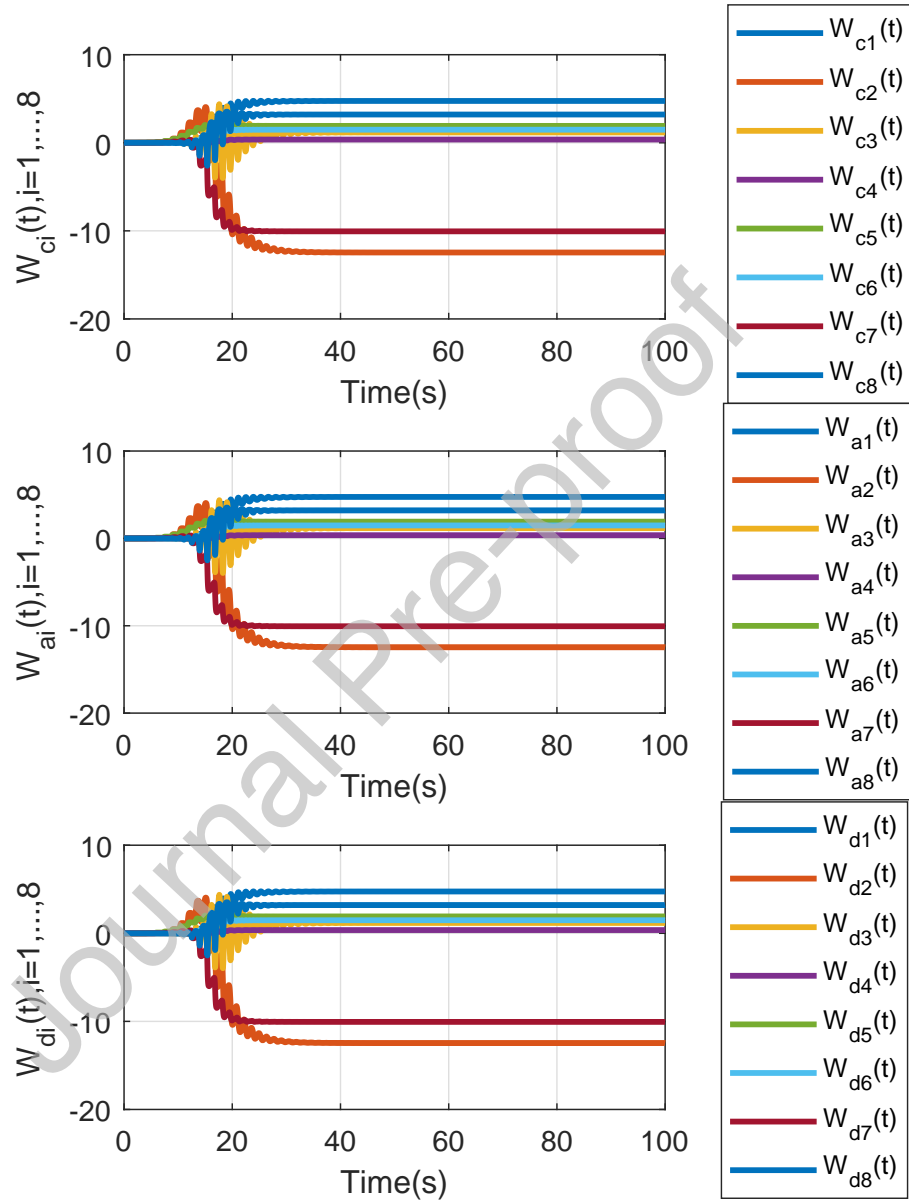


Figure 7: Evolution of the actor and critic weights using barrier-actor-critic learning.

363 full-state constraints are guaranteed to be satisfied during the learning phase. Finally, simulation
364 studies are conducted to verify the efficacy of the presented barrier-actor-critic online learning.

365 References

- 366 [1] M. Rehan, C. K. Ahn, M. Chadli, Consensus of one-sided lipschitz multi-agents under in-
367 put saturation, *IEEE Transactions on Circuits and Systems II: Express Briefs* doi:10.1109/
368 TCSII.2019.2923721.
- 369 [2] Z. Liu, Z. Zhao, C. K. Ahn, Boundary constrained control of flexible string systems subject
370 to disturbances, *IEEE Transactions on Circuits and Systems II: Express Briefs* doi:10.1109/
371 TCSII.2019.2901283.
- 372 [3] Z. Zhao, Z. Liu, Z. Li, N. Wang, J. Yang, Control design for a vibrating flexible marine riser
373 system, *Journal of the Franklin Institute* 354 (18) (2017) 8117 – 8133.
- 374 [4] R. R. Selmic, F. L. Lewis, Deadzone compensation in motion control systems using neural
375 networks, *IEEE Transactions on Automatic Control* 45 (4) (2000) 602–613.
- 376 [5] W. He, B. Huang, Y. Dong, Z. Li, C. Su, Adaptive neural network control for robotic manip-
377 ulators with unknown deadzone, *IEEE Transactions on Cybernetics* 48 (9) (2018) 2670–2682.
- 378 [6] Y. Liu, S. Lu, S. Tong, Neural network controller design for an uncertain robot with time-
379 varying output constraint, *IEEE Transactions on Systems, Man, and Cybernetics: Systems*
380 47 (8) (2017) 2060–2068.
- 381 [7] Q. Zhou, L. Wang, C. Wu, H. Li, H. Du, Adaptive fuzzy control for nonstrict-feedback sys-
382 tems with input saturation and output constraint, *IEEE Transactions on Systems, Man, and*
383 *Cybernetics: Systems* 47 (1) (2017) 1–12.
- 384 [8] R. R. Selmic, F. L. Lewis, Neural-network approximation of piecewise continuous functions:
385 application to friction compensation, *IEEE Transactions on Neural Networks* 13 (3) (2002)
386 745–751.
- 387 [9] J. Na, Q. Chen, X. Ren, Y. Guo, Adaptive prescribed performance motion control of servo
388 mechanisms with friction compensation, *IEEE Transactions on Industrial Electronics* 61 (1)
389 (2014) 486–494.

- 390 [10] G. Tao, P. V. Kokotovic, Adaptive control of plants with unknown hystereses, *IEEE Trans-*
391 *actions on Automatic Control* 40 (2) (1995) 200–212.
- 392 [11] M. Chen, S. S. Ge, Adaptive neural output feedback control of uncertain nonlinear systems
393 with unknown hysteresis using disturbance observer, *IEEE Transactions on Industrial Elec-*
394 *tronics* 62 (12) (2015) 7706–7716.
- 395 [12] Z. Zhao, S. Lin, D. Zhu, G. Wen, Vibration control of a riser-vessel system subject to input
396 backlash and extraneous disturbances, *IEEE Transactions on Circuits and Systems II: Express*
397 *Briefs* doi:10.1109/TCSII.2019.2914061.
- 398 [13] G. A. Rovithakis, Robust redesign of a neural network controller in the presence of unmodeled
399 dynamics, *IEEE Transactions on Neural Networks* 15 (6) (2004) 1482–1490.
- 400 [14] J. C. Doyle, K. Glover, P. P. Khargonekar, B. A. Francis, State-space solutions to standard H_2
401 and H_∞ control problems, *IEEE Transactions on Automatic Control* 34 (8) (1989) 831–847.
- 402 [15] A. J. van der Schaft, l_2 -gain analysis of nonlinear systems and nonlinear state-feedback h_∞
403 control, *IEEE Transactions on Automatic Control* 37 (6) (1992) 770–784.
- 404 [16] P. A. Ioannou, J. Sun, *Robust Adaptive Control*, Prentice-Hall, Inc., Upper Saddle River, NJ,
405 USA, 1995.
- 406 [17] M. Krstic, P. V. Kokotovic, I. Kanellakopoulos, *Nonlinear and Adaptive Control Design*, 1st
407 Edition, John Wiley & Sons, Inc., New York, NY, USA, 1995.
- 408 [18] S. Boyd, L. Vandenberghe, *Convex Optimization*, Cambridge University Press, New York,
409 NY, USA, 2004.
- 410 [19] K. P. Tee, S. S. Ge, E. H. Tay, Barrier lyapunov functions for the control of output-constrained
411 nonlinear systems, *Automatica* 45 (4) (2009) 918 – 927.
- 412 [20] B. Ren, S. S. Ge, K. P. Tee, T. H. Lee, Adaptive neural control for output feedback nonlinear
413 systems using a barrier lyapunov function, *IEEE Transactions on Neural Networks* 21 (8)
414 (2010) 1339–1345.
- 415 [21] Y.-J. Liu, S. Lu, S. Tong, X. Chen, C. P. Chen, D.-J. Li, Adaptive control-based barrier
416 lyapunov functions for a class of stochastic nonlinear systems with full state constraints,
417 *Automatica* 87 (2018) 83 – 93.

- 418 [22] Y.-J. Liu, S. Tong, Barrier Lyapunov functions-based adaptive control for a class of nonlinear
419 pure-feedback systems with full state constraints, *Automatica* 64 (2016) 70 – 75.
- 420 [23] W. He, Y. Chen, Z. Yin, Adaptive neural network control of an uncertain robot with full-state
421 constraints, *IEEE Transactions on Cybernetics* 46 (3) (2016) 620–629.
- 422 [24] D. Li, D. Li, Adaptive tracking control for nonlinear time-varying delay systems with full state
423 constraints and unknown control coefficients, *Automatica* 93 (2018) 444 – 453.
- 424 [25] C. P. Bechlioulis, G. A. Rovithakis, Robust adaptive control of feedback linearizable mimo
425 nonlinear systems with prescribed performance, *IEEE Transactions on Automatic Control*
426 53 (9) (2008) 2090–2099.
- 427 [26] A. K. Kostarigka, G. A. Rovithakis, Adaptive dynamic output feedback neural network control
428 of uncertain mimo nonlinear systems with prescribed performance, *IEEE Transactions on*
429 *Neural Networks and Learning Systems* 23 (1) (2012) 138–149.
- 430 [27] C. P. Bechlioulis, G. A. Rovithakis, Decentralized robust synchronization of unknown high
431 order nonlinear multi-agent systems with prescribed transient and steady state performance,
432 *IEEE Transactions on Automatic Control* 62 (1) (2017) 123–134.
- 433 [28] A. Theodorakopoulos, G. A. Rovithakis, Guaranteeing preselected tracking quality for un-
434 certain strict-feedback systems with deadzone input nonlinearity and disturbances via low-
435 complexity control, *Automatica* 54 (2015) 135 – 145.
- 436 [29] A. K. Kostarigka, Z. Doulgeri, G. A. Rovithakis, Prescribed performance tracking for flexible
437 joint robots with unknown dynamics and variable elasticity, *Automatica* 49 (5) (2013) 1137 –
438 1147.
- 439 [30] Y. Yang, C. Ge, H. Wang, X. Li, C. Hua, Adaptive neural network based prescribed per-
440 formance control for teleoperation system under input saturation, *Journal of the Franklin*
441 *Institute* 352 (5) (2015) 1850 – 1866.
- 442 [31] E. Arabi, T. Yucelen, B. C. Gruenwald, M. Fravolini, S. Balakrishnan, N. T. Nguyen, A
443 neuroadaptive architecture for model reference control of uncertain dynamical systems with
444 performance guarantees, *Systems & Control Letters* 125 (2019) 37 – 44.
- 445 [32] F. L. Lewis, D. Vrabie, V. L. Syrmos, *Optimal control*, John Wiley & Sons, Hoboken, NJ,
446 USA, 2012.

- 447 [33] B. Kiumarsi, K. G. Vamvoudakis, H. Modares, F. L. Lewis, Optimal and autonomous control
448 using reinforcement learning: A survey, *IEEE Transactions on Neural Networks and Learning*
449 *Systems* 29 (6) (2018) 2042–2062.
- 450 [34] D. Liu, Q. Wei, Policy iteration adaptive dynamic programming algorithm for discrete-time
451 nonlinear systems, *IEEE Transactions on Neural Networks and Learning Systems* 25 (3) (2014)
452 621–634.
- 453 [35] Y. Yang, D. Wunsch, Y. Yin, Hamiltonian-driven adaptive dynamic programming for con-
454 tinuous nonlinear dynamical systems, *IEEE Transactions on Neural Networks and Learning*
455 *Systems* 28 (8) (2017) 1929–1940.
- 456 [36] Y. Yang, K. G. Vamvoudakis, H. Ferraz, H. Modares, Dynamic intermittent Q-learning-based
457 model-free suboptimal co-design of L_2 -stabilization, *International Journal of Robust and Non-*
458 *linear Control* 29 (9) (2019) 2673–2694.
- 459 [37] K. G. Vamvoudakis, F. L. Lewis, Online actor-critic algorithm to solve the continuous-time
460 infinite horizon optimal control problem, *Automatica* 46 (5) (2010) 878 – 888.
- 461 [38] H. Modares, F. L. Lewis, Linear quadratic tracking control of partially-unknown continuous-
462 time systems using reinforcement learning, *IEEE Transactions on Automatic Control* 59 (11)
463 (2014) 3051–3056.
- 464 [39] H. Modares, F. L. Lewis, Z. Jiang, H_∞ tracking control of completely unknown continuous-
465 time systems via off-policy reinforcement learning, *IEEE Transactions on Neural Networks*
466 *and Learning Systems* 26 (10) (2015) 2550–2562.
- 467 [40] Y. Yang, Z. Guo, H. Xiong, D. Ding, Y. Yin, D. C. Wunsch, Data-driven robust control of
468 discrete-time uncertain linear systems via off-policy reinforcement learning, *IEEE Transactions*
469 *on Neural Networks and Learning Systems* 30 (12) (2019) 3735 – 3747.
- 470 [41] D. Wang, D. Liu, Learning and guaranteed cost control with event-based adaptive critic
471 implementation, *IEEE Transactions on Neural Networks and Learning Systems* 29 (12) (2018)
472 6004–6014.
- 473 [42] D. Wang, Intelligent critic control with robustness guarantee of disturbed nonlinear plants,
474 *IEEE Transactions on Cybernetics*.

- 475 [43] Y. Yang, H. Modares, D. C. Wunsch, Y. Yin, Optimal containment control of unknown
476 heterogeneous systems with active leaders, *IEEE Transactions on Control Systems Technology*
477 27 (3) (2019) 1228–1236.
- 478 [44] Y. Yang, H. Modares, D. C. Wunsch, Y. Yin, Leaderfollower output synchronization of linear
479 heterogeneous systems with active leader using reinforcement learning, *IEEE Transactions on*
480 *Neural Networks and Learning Systems* 29 (6) (2018) 2139–2153.
- 481 [45] D. Wang, D. Liu, Neural robust stabilization via event-triggering mechanism and adaptive
482 learning technique, *Neural Networks* 102 (2018) 27 – 35.
- 483 [46] D. Zhao, Q. Zhang, D. Wang, Y. Zhu, Experience replay for optimal control of nonzero-
484 sum game systems with unknown dynamics, *IEEE Transactions on Cybernetics* 46 (3) (2016)
485 854–865.
- 486 [47] H. Modares, F. L. Lewis, M. Naghibi-Sistani, Adaptive optimal control of unknown
487 constrained-input systems using policy iteration and neural networks, *IEEE Transactions on*
488 *Neural Networks and Learning Systems* 24 (10) (2013) 1513–1525.
- 489 [48] J. Sun, C. Liu, Disturbance observer-based robust missile autopilot design with full-state
490 constraints via adaptive dynamic programming, *Journal of the Franklin Institute* 355 (5)
491 (2018) 2344 – 2368.
- 492 [49] H. Modares, F. L. Lewis, M.-B. N. Sistani, Online solution of nonquadratic two-player zero-
493 sum games arising in the h_∞ control of constrained input systems, *International Journal of*
494 *Adaptive Control and Signal Processing* 28 (3-5) (2014) 232–254.
- 495 [50] M. Polycarpou, P. Ioannou, A robust adaptive nonlinear control design, *Automatica* 32 (3)
496 (1996) 423 – 427.
- 497 [51] N. us Saqib, M. Rehan, M. Hussain, N. Iqbal, H. ur Rashid, Delay-range-dependent static anti-
498 windup compensator design for nonlinear systems subjected to input-delay and saturation,
499 *Journal of the Franklin Institute* 354 (14) (2017) 5919 – 5948.
- 500 [52] M. Hussain, M. Rehan, C. Ki Ahn, M. Tufail, Robust antiwindup for one-sided lipschitz sys-
501 tems subject to input saturation and applications, *IEEE Transactions on Industrial Electronics*
502 65 (12) (2018) 9706–9716.

- 503 [53] Z. Zhao, X. He, G. Wen, Boundary robust adaptive anti-saturation control of vibrating flexible
504 riser systems, *Ocean Engineering* 179 (2019) 298 – 306.
- 505 [54] Z. Zhao, X. He, Z. Ren, G. Wen, Boundary adaptive robust control of a flexible riser system
506 with input nonlinearities, *IEEE Transactions on Systems, Man, and Cybernetics: Systems*
507 49 (10) (2019) 1971–1980.
- 508 [55] H. Modares, F. L. Lewis, Optimal tracking control of nonlinear partially-unknown constrained-
509 input systems using integral reinforcement learning, *Automatica* 50 (7) (2014) 1780 – 1792.
- 510 [56] T. Başar, P. Bernhard, H_∞ Optimal Control and Related Minimax Design Problems: A
511 Dynamic Game Approach, Springer, Berlin, Germany, 2008.

Conflict of interest statement

The authors declared that they have no conflicts of interest to this work.

We declare that we do not have any commercial or associative interest that represents a conflict of interest in connection with the work submitted.

Journal Pre-proof