

## Reinforcement learning based - sliding mode control for trajectory tracking of quadrotor unmanned aerial vehicles under disturbances

Tran Thai Duong<sup>1</sup>, Do Duc Manh<sup>1</sup>, Nguyen Chi Nhan<sup>1</sup>,  
Le Duc Thinh<sup>2</sup>, Nguyen Tung Lam<sup>1</sup>, Nguyen Danh Huy<sup>1\*</sup>

<sup>1</sup>School of Electrical and Electronic Engineering, Hanoi University of Science and Technology, 1 Dai Co Viet, Hai Ba Trung, Hanoi, Vietnam;

<sup>2</sup>Faculty of Electrical and Electronics Engineering, Thuyloi University, 175 Tay Son, Dong Da, Hanoi, Vietnam.

\*Corresponding author: huy.nguyendanh@hust.edu.vn

Received 4 Sep. 2024; Revised 15 Nov.2024; Accepted 5 Feb. 2025; Published 25 Feb. 2025.

DOI: <https://doi.org/10.54939/1859-1043.j.mst.101.2025.39-46>

### ABSTRACT

*In this article, a reinforcement learning (RL)-based sliding mode control (SMC) is proposed for trajectory tracking of a quadrotor unmanned aerial vehicle (QUAV) under external disturbances. First, an actor-critic RL framework sliding mode control is provided to tackle the optimal control problem without external disturbances. Secondly, the simulation in an environment with disturbances is carried out to show the robustness of the proposed controller. Theoretical analysis shows that the position and attitude tracking errors converge to a preset region, and the weight estimation errors of the actor-critic networks are uniformly ultimately bounded (UUB). Finally, a comparison of the numerical simulations between the proposed controller and traditional sliding mode controller and the Backstepping (BSP) technique is provided to indicate the advantages and improved performance of the RL-based SMC.*

**Keywords:** Reinforcement learning; Sliding mode control; Optimal control; Actor/critic structure; Quadrotor unmanned aerial vehicle (QUAV).

### 1. INTRODUCTION

Unmanned aerial vehicles (UAVs) have recently garnered significant interest and attention due to their extensive applications in logistics delivery, environmental monitoring, and various civil and military uses [1]. Quadrotor UAVs, in particular, are popular among users because of their simple structure, flexible flight capabilities, strong environmental adaptability, high maneuverability, and numerous other advantages [2]. Compared to other UAVs, quadrotor UAVs excel in vertical takeoff and landing.

However, designing a control scheme for a quadrotor UAV under external disturbances is challenging. The system's nonlinearity and strong coupling characteristics further complicate achieving high-precision control. Various control schemes, such as adaptive control, sliding mode control (SMC), and robust control, have been developed to address control issues involving disturbances and uncertainties [3]. Among these, SMC offers a faster response [4]. It is important to note that the energy resources UAVs consume are valuable and limited. Thus, designing a high-precision control scheme for a UAV that balances control performance and control cost is crucial. From a mathematical perspective, establishing the Hamilton–Jacobi–Bellman (HJB) function and its solution is essential for solving optimal control problems [5]. However, deriving an analytical solution from the HJB function is challenging due to the non-linear continuous-time systems [6]. To address this, Werbos [7] initially developed a reinforcement learning (RL) scheme with an actor-critic (AC) structure. In this scheme, a critic network approximates the value function, while an actor network determines the optimal control policy. Building on Werbos' work, Vamvoudakis et al. created an online AC algorithm to solve the continuous-time infinite horizon optimal control problem [8]. Ma et al. developed a learning-based adaptive sliding mode control scheme for a tethered space robot with limited inputs [9]. Given the unknown disturbances, accurately

determining the system dynamics for UAVs is practically difficult, limiting the methods' applicability. Thus, enhancing the robustness of these methods for practical systems with unknown disturbances remains a further challenge.

Motivated by these considerations, a novel RL-based tracking controller for quadrotors that accounts for external disturbances is being investigated. The main contributions and features of the proposed method can be summarized as follows:

- An online RL-based, nearly optimal controller is developed by combining the actor-critic (AC) structure with the hyperbolic tangent performance index. Additionally, the robustness of the learning-based controller is enhanced by integrating a super-twisting-like sliding mode control.
- Comparative numerical simulations demonstrate that the proposed method can enhance performance in terms of transient response and steady-state accuracy while requiring less control effort.

## 2. PROBLEM

### 2.1. Dynamic models and problem statement

#### 2.1.1. Dynamic model

The Quad-rotor UAV is considered under two reference frames. The body-fixed reference frame is denoted by  $\mathbf{B} = \{O_b, x_b, y_b, z_b\}$ , and  $\mathbf{E} = \{O_e, x_e, y_e, z_e\}$  is the earth-fixed reference frame. In  $\mathbf{B}$  frame,  $\mathbf{P} = [x, y, z]^T$  is the position vector, and  $\Theta = [\phi, \theta, \psi]^T$  is Euler angles attitude vector (where  $\phi$  is roll angle,  $\theta$  is pitch angle, and  $\psi$  is yaw angles). The linear velocity vector of the QUAV in the earth-fixed frame and the angular velocity vector of QUAV in the body-fixed frame are denoted as  $\mathbf{V} = [v_x, v_y, v_z]^T$  and  $\boldsymbol{\omega} = [\omega_\phi, \omega_\theta, \omega_\psi]^T$ . The position and attitude dynamics of UAVs are considered as follows:

$$\begin{cases} \dot{\mathbf{P}} = \mathbf{V} \\ m\dot{\mathbf{V}} = \mathbf{R}\mathbf{F} - \mathbf{K}\mathbf{V} - mg\mathbf{E}_3 + \mathbf{d} \end{cases} \quad (1)$$

$$\begin{cases} \dot{\Theta} = \mathbf{T}\boldsymbol{\omega} \\ \mathbf{I}\dot{\boldsymbol{\omega}} = -\boldsymbol{\omega} \times (\mathbf{I}\boldsymbol{\omega}) + \mathbf{F}_\Theta - \mathbf{G} - \mathbf{K}_\Theta \dot{\Theta} + \mathbf{d}_\Theta \end{cases} \quad (2)$$

where  $\mathbf{R}$  is the kinematic rotation matrix,  $\mathbf{F} = [0, 0, u_m]^T$  is the force vector in  $\mathbf{B}$  frame.  $\mathbf{K} = \text{diag}(k_x, k_y, k_z)$  is the aerodynamic matrix,  $m$  is the mass of the QUAV,  $\mathbf{E}_3$  is the unit vector defined as  $\mathbf{E}_3 = [0, 0, 1]^T$ , and  $\mathbf{d} = [d_x, d_y, d_z]^T$  is the disturbance vector.  $\mathbf{T}$  is the attitude kinematic transformation matrix,  $\mathbf{I} = \text{diag}(I_x, I_y, I_z)$  is the inertia matrix,  $\mathbf{F}_\Theta = [u_\phi, u_\theta, u_\psi]^T$  is the attitude control torque vector.  $\mathbf{G} = \left( \sum_{i=1}^4 (-1)^i w_i \right) I_i (\boldsymbol{\omega} \times \mathbf{E}_3)$  is the gyroscopic torque vector,  $\mathbf{K}_\Theta = \text{diag}(k_\phi, k_\theta, k_\psi)$  is the attitude aerodynamic matrix.  $\mathbf{d}_\Theta = [d_\phi, d_\theta, d_\psi]^T$  is the attitude disturbance vector. The kinematic rotation matrix  $\mathbf{R}$  and the attitude kinematic transformation matrix are defined as:

$$\mathbf{R} = \begin{bmatrix} c_\psi c_\theta & c_\psi s_\theta s_\phi - s_\psi c_\phi & c_\psi s_\theta c_\phi + s_\psi s_\phi \\ s_\psi c_\theta & s_\psi s_\theta s_\phi + c_\psi c_\phi & s_\psi s_\theta c_\phi - c_\psi s_\phi \\ -s_\theta & c_\theta s_\phi & c_\theta c_\phi \end{bmatrix} \quad (3)$$

$$\mathbf{T} = \begin{bmatrix} 1 & s_\phi \tan \theta & c_\phi \tan \theta \\ 0 & c_\phi & -s_\phi \\ 0 & s_\phi \sec \theta & c_\phi \sec \theta \end{bmatrix} \quad (4)$$

Let  $\mathbf{f} = (-\mathbf{KV} - g\mathbf{E}_3) / m$ , (1) can be written as [10]:

$$\begin{cases} \dot{\mathbf{P}} = \mathbf{V} \\ \dot{\mathbf{V}} = \mathbf{RF} / m - \mathbf{f} + \mathbf{d} / m \end{cases} \quad (5)$$

Let  $\mathbf{f}_\Theta = \mathbf{T}^{-1}\mathbf{I}^{-1}(-\boldsymbol{\omega} \times (\mathbf{I}\boldsymbol{\omega}) - \mathbf{G} - \mathbf{K}_\Theta \dot{\boldsymbol{\Theta}})$ , equation (2) can be written as [10]:

$$\begin{cases} \dot{\boldsymbol{\Theta}} = \boldsymbol{\Omega} \\ \boldsymbol{\Omega} = \mathbf{f}_\Theta + \mathbf{T}^{-1}\mathbf{I}^{-1}\mathbf{F}_\Theta + \mathbf{T}^{-1}\mathbf{I}^{-1}\mathbf{d}_\Theta \end{cases} \quad (6)$$

where  $\boldsymbol{\Omega} = [\dot{\phi}, \dot{\theta}, \dot{\psi}]^T$  is the angular velocity vector of QUAV in the body-fixed frame.

**Assumption 1:** UAV can be considered a rigid body in flight missions, and its rotation is limited as  $-\pi/2 < \phi < \pi/2$ ,  $-\pi/2 < \theta < \pi/2$ ,  $-\pi \leq \psi \leq \pi$ .

**Assumption 2:** The desired reference  $\mathbf{P}_d = [x_d, y_d, z_d]^T$ ,  $\psi_d$ , its first and second derivations  $\dot{\mathbf{P}}_d$ ,  $\dot{\psi}_d$ ,  $\ddot{\mathbf{P}}_d$ ,  $\ddot{\psi}_d$  are known and bounded. They are available signals in the UAV control.

**Remark 1:** In fact, the rotational velocities  $\omega_i, i = \overline{1,4}$  of four rotors are the direct control input of QUAV. The relationship between both total thrust  $u_m$  and control toque  $u_\phi, u_\theta, u_\psi$ :

$$\begin{aligned} u_m &= k_t (w_1^2 + w_2^2 + w_3^2 + w_4^2); u_\phi = k_l l (w_2^2 - w_4^2) \\ u_\theta &= k_l l (-w_1^2 + w_3^2); u_\psi = k_d (w_1^2 - w_2^2 + w_3^2 - w_4^2) \end{aligned} \quad (7)$$

where  $k_t$ ,  $k_d$ , and  $w_i$  are the thrust factor, drag factor, and rotational velocity of the  $i$ th rotor, respectively,  $l$  is the distance from the rotors to the mass center.

### 2.1.2. Preliminaries

**Lemma 1** [11]. Considering a nonlinear function  $f(x) = \ln[1 - \tanh^2(x)]$ , the equation:  $f(x) = \ln(4) - 2x \operatorname{sign}(x) + k_f$  always holds,  $k_f$  is bounded by a real positive constant.

**Lemma 2** [12]. For any  $\alpha > 0$ , the following inequality holds  $0 \leq |x| - x \tanh\left(\frac{x}{\alpha}\right) \leq k\alpha$ . Where

$k$  satisfies  $k = e^{-k+1}$ .

## 2.2. Optimized UAV position control

In this subsection, the position control problem will be tackled with the RL algorithm, the attitude loop will be carried out by the sliding mode controller. For simplicity, the (5) can be rewritten as:

$$\begin{cases} \dot{\mathbf{P}} = \mathbf{V} \\ \dot{\mathbf{V}} = -[0 \ 0 \ g]^T + \mathbf{U} \end{cases} \quad (8)$$

with  $\mathbf{U} = [U_1 \ U_2 \ U_3]^T = \mathbf{RF} / m$ . The relationship between  $\mathbf{U}$  and  $\mathbf{F}$  is:

$$\begin{aligned} U_1 &= (c_\phi s_\theta c_\psi + s_\phi s_\psi) u_m / m \\ U_2 &= (c_\phi s_\theta s_\psi - s_\phi c_\psi) u_m / m \\ U_3 &= (c_\phi c_\theta) u_m / m \end{aligned} \quad (9)$$

Therefore, the control signal  $u_m$  can be solved as:

$$u_m = m \sqrt{U_1^2 + U_2^2 + U_3^2} \quad (10)$$

with the desired position trajectory  $\mathbf{P}_{re}(t) = [x_{re}(t), y_{re}(t), z_{re}(t)]^T$ , the tracking error can be obtained as  $\mathbf{e}_{p1} = \mathbf{P}(t) - \mathbf{P}_{re}(t)$ , and  $\mathbf{e}_{p2} = \mathbf{V}(t) - \dot{\mathbf{P}}_{re}(t)$ , which yields:

$$\dot{\mathbf{e}}_{p1}(t) = \mathbf{e}_{p2}(t); \quad \dot{\mathbf{e}}_{p2}(t) = \mathbf{U} - [0 \ 0 \ g]^T - \ddot{\mathbf{P}}_{re}(t) + \mathbf{d} / m \quad (11)$$

The sliding variable  $\mathbf{s}$  is defined as:

$$\mathbf{s} = \mathbf{e}_{p2} + \mathbf{C} \mathbf{e}_{p1}$$

where  $\mathbf{C}$  is a positive diagonal matrix.

$$\dot{\mathbf{s}} = \mathbf{C} \mathbf{e}_{p2}(t) - [0 \ 0 \ g]^T - \ddot{\mathbf{P}}_{re}(t) + \mathbf{U} + \mathbf{d} / m \quad (12)$$

The conventional super-twisting controller is designed as follows:

$$\mathbf{U} = - \left( \mathbf{C} \mathbf{e}_{p2} - [0 \ 0 \ g]^T - \ddot{\mathbf{P}}_{re}(t) \right) + \lambda_1 \text{sig}(\mathbf{s}) - \mathbf{u}_d \quad (13)$$

where  $\text{sig}(\mathbf{s}) = C [ |s_i|^{1/2} \text{sign}(s_i) ]$ ;  $\mathbf{u}_d = -C [ \lambda_2 \text{sign}(s_i) ]$ . With the super-twisting sliding mode controller, the chattering phenomenon happening in traditional SMC can be attenuated thanks to the 2 continuous and smooth terms  $\text{sig}(\mathbf{s})$  and  $\mathbf{u}_d = \int -C [ \lambda_2 \text{sign}(s_i) ] dt$  instead of function  $\text{sign}(s_i)$  that can discontinuously change [17]. However, the cost index cannot be optimized, to tackle this problem, a RL-based controller is provided to improve the performance, assuming there is a family of admissible control inputs:

$$\mathbf{U} = \mathbf{u}_0 - \lambda_1 \text{sig}(\mathbf{s}) + \mathbf{u}_d; \quad (14)$$

Consider the following cost function:

$$V(s) = \int_t^\infty \{ V_s(\tau) + J[\mathbf{u}_0(\tau)] \} d\tau \quad (15)$$

where  $V_s = \mathbf{s}^T \mathbf{Q} \mathbf{s}$  with  $\mathbf{Q}$  is a positive diagonal matrix,  $J(\mathbf{u}_0)$  is specified as:

$$J(\mathbf{u}_0) = 2 \int_0^{\mathbf{u}_0} \left[ \lambda_1 \tanh^{-1} \left( \frac{\mathbf{v}}{\lambda_{\phi 1}} \right)^T \mathbf{\Omega} \right] d\mathbf{v} \quad (16)$$

where  $\mathbf{\Omega}$  is a positive diagonal matrix, the Lyapunov function of (15) is evaluated as:

$$V_s(t) + J[\mathbf{u}_0(t)] + \nabla V_s^T \dot{\mathbf{s}} = 0 \quad (17)$$

Therefore, the corresponding HJB equation can be formulated as:

$$\min_{\mathbf{u}_0} \{ V_s + J(\mathbf{u}_0) + \nabla V_s^{*T} \dot{\mathbf{s}} \} = 0 \quad (18)$$

the  $\mathbf{u}_0$  that minimizes the cost function is the solution of the following equation:

$$\frac{\partial}{\partial \mathbf{u}_0^*} [V_s + J(\mathbf{u}_0^*) + \nabla V_s^{*T} \dot{\mathbf{s}}] = 0 \quad (19)$$

Solving (21) yields:

$$\mathbf{u}_0^* = -\lambda_1 \tanh(\bar{\mathbf{u}}_0^*), \bar{\mathbf{u}}_0^* = \frac{1}{2\lambda_1} \boldsymbol{\omega}^{-1} \nabla V_s^* \quad (20)$$

Substituting  $\mathbf{u}_0^*$  into (16), the following equation can be obtained:

$$J(\mathbf{u}_0^*) = -\nabla V_s^{*T} \mathbf{I}_3 \mathbf{u}_0^* + \lambda_1^2 \bar{\boldsymbol{\omega}} \ln[\mathbf{1}_c - \tanh^2(\bar{\mathbf{u}}_0^*)] \quad (21)$$

From (19) and (21), it can be derived that:

$$V_s + \nabla V_s^{*T} \dot{\mathbf{s}} + \lambda_1^2 \bar{\boldsymbol{\omega}} \ln[\mathbf{1}_c - \tanh^2(\bar{\mathbf{u}}_0^*)] = 0 \quad (22)$$

Nevertheless, it is impossible to obtain the analytical solutions of the HJB equations in (22), which are nonlinear partial differential equations. Based on the fine approximation ability of NN, two RL-based control policies will be proposed to deal with this intractable problem. The following approximation of  $V_s^*$  can be established

$$V^*(\mathbf{s}) = \mathbf{W}^{*T} \boldsymbol{\sigma}[\mathbf{s}(t)] + \varepsilon[\mathbf{s}(t)] \quad (23)$$

where  $\boldsymbol{\sigma}[\mathbf{s}(t)]$  are the basis function vectors and  $\mathbf{W}^*$  denotes the optimal weight vector,  $\varepsilon[\mathbf{s}(t)]$  represent the approximation errors. Subsequently, the gradient of  $V^*(\mathbf{s})$  is

$$\nabla V_s^* = (\nabla \boldsymbol{\sigma}_s[\mathbf{s}(t)])^T \mathbf{W}^* + \nabla \varepsilon_s[\mathbf{s}(t)] \quad (24)$$

where  $\nabla \boldsymbol{\sigma}_s = \partial \boldsymbol{\sigma} / \partial \mathbf{s}$ ,  $\nabla \varepsilon_s = \partial \varepsilon / \partial \mathbf{s}$ . In light of the universal approximation property of NN for smooth functions on prescribed compact sets, the approximation errors  $\varepsilon[\mathbf{s}(t)]$  and  $\nabla \varepsilon[\mathbf{s}(t)]$  are bounded with a finite dimension of  $\boldsymbol{\sigma}[\mathbf{s}(t)]$ . Moreover, it is assumed that  $\|\mathbf{W}^*\|$ ,  $\boldsymbol{\sigma}[\mathbf{s}(t)]$  and  $\nabla \boldsymbol{\sigma}[\mathbf{s}(t)]$  are bounded. Recalling (20), the NN-based nearly optimal control law can be formulated as:

$$\begin{cases} \mathbf{U} = \hat{\mathbf{U}}_0^* - \mathbf{I}_3^{-1} \{ \lambda_1 \text{sig}(\mathbf{s}) - \mathbf{u}_d \} \\ \hat{\mathbf{U}}_0^* = -\lambda_1 \tanh \left[ \frac{1}{2\lambda_1} \boldsymbol{\omega}^{-1} \mathbf{I}_3^T (\nabla \boldsymbol{\sigma}_s^T \mathbf{W}^* + \nabla \varepsilon_s) \right] \end{cases} \quad (25)$$

Defining the Bellman error as:

$$B_e = \lambda_1^2 \bar{\boldsymbol{\omega}} \left\{ \ln[\mathbf{1}_c - \tanh^2(\bar{\mathbf{U}}_0^*)] - \ln[\mathbf{1}_c - \tanh^2(\hat{\mathbf{U}}_0^*)] \right\} \quad (26)$$

with  $\hat{\mathbf{U}}_0^* = \frac{1}{2\lambda_1} \boldsymbol{\omega}^{-1} \mathbf{I}_3^T \nabla \boldsymbol{\sigma}_s^T \mathbf{W}^*$ , (19) can be rewritten as:

$$V_s + (\nabla \boldsymbol{\sigma}_s)^T \mathbf{W}^* \boldsymbol{\Gamma} + \lambda_1^2 \bar{\boldsymbol{\omega}} \ln[\mathbf{1}_c - \tanh^2(\hat{\mathbf{U}}_0^*)] + \varepsilon_H = 0 \quad (27)$$

where  $\varepsilon_H = \nabla \varepsilon_s^T \boldsymbol{\Gamma} + B_e + \mathbf{W}^* \nabla \boldsymbol{\sigma}_s D_2$  is the bounded HJB error. In this paper, the optimal weight  $W^*$  is generated by the online RL scheme with the AC structure. In this context, the nearly optimal control policy is formulated as:

$$\begin{cases} \mathbf{U} = \hat{\mathbf{U}}_a^* - \mathbf{I}_3^{-1} \{ \lambda_1 \text{sig}(\mathbf{s}) - \mathbf{u}_d \} \\ \hat{\mathbf{U}}_a^* = -\lambda_1 \tanh(\hat{\mathbf{U}}_a^*), \hat{\mathbf{U}}_a^* = \frac{1}{2\lambda_1} \boldsymbol{\omega}^{-1} \mathbf{I}_3^T \nabla \boldsymbol{\sigma}_s^T \hat{\mathbf{W}}_a \end{cases} \quad (28)$$

where  $\hat{\mathbf{W}}_a$  is the weight of the actor network. Moreover, the performance index (15) can be

estimated as:  $\hat{V}_s = \hat{\mathbf{W}}_c^T \boldsymbol{\sigma}_s[\mathbf{s}(t)]$ . Where  $\hat{\mathbf{W}}_c$  is the weight of the critic network. The adaptation laws for  $\hat{\mathbf{W}}_c$  and  $\hat{\mathbf{W}}_a$  are designed as:

$$\begin{aligned} \dot{\hat{\mathbf{W}}}_c &= -A_1 \left\{ \bar{\boldsymbol{\theta}} \left[ \boldsymbol{\theta}^T \hat{\mathbf{W}}_c + V_s + J(\hat{\mathbf{U}}_a^*) \right] + p_c (\hat{\mathbf{W}}_c - \hat{\mathbf{W}}_a) \right\} \\ \dot{\hat{\mathbf{W}}}_a &= -A_2 \left[ p_a (\hat{\mathbf{W}}_a - \hat{\mathbf{W}}_c) - \frac{\lambda_1 p_a}{p_c} \nabla \boldsymbol{\sigma}_s \mathbf{I}_3 \boldsymbol{\Psi} \bar{\boldsymbol{\theta}}^T \hat{\mathbf{W}}_c \right] \end{aligned} \quad (29)$$

where  $\boldsymbol{\theta} = \nabla \boldsymbol{\sigma}_s(\boldsymbol{\Gamma} + \mathbf{I}_3 \hat{\mathbf{U}}_a^*)$ ,  $\bar{\boldsymbol{\theta}} = \boldsymbol{\theta} / (\boldsymbol{\theta}^T \boldsymbol{\theta} + 1)^2$ ,  $\boldsymbol{\Psi} = \tanh\left(\frac{\hat{\mathbf{U}}_a^*}{\kappa}\right) - \tanh\left(\frac{\hat{\mathbf{U}}_a}{\kappa}\right)$ ,  $p_c, p_a, k > 0, A_1,$

$A_2$  are positive diagonal matrices. Assuming  $\boldsymbol{\theta}_1 = \boldsymbol{\theta} / (\boldsymbol{\theta}^T \boldsymbol{\theta} + 1)$  is persistently excited.

**Theorem 1.** Considering the systems in (8), if the control policies are designed as (28), with the weight update rules in (29), then the weight errors and the tracking errors will be guaranteed to be ultimately uniformly bounded under the proposed RL-based control policies. Based on SMC, attitude control for the angle can be calculated as:

$$\mathbf{F}_\Theta = \mathbf{B}^{-1} \left( -\varepsilon \text{sat}(\mathbf{s}) - k\mathbf{s} - c\dot{\mathbf{s}} - \mathbf{f}_\Theta - \ddot{\mathbf{n}}_{re} \right) \quad (30)$$

If the control policies are designed as (30), the angular errors  $e_\phi, e_\theta, e_\psi$  will converge to the origin. Therefore, the angles converge to the reference values.

### 3. NUMERICAL SIMULATION RESULTS

#### 3.1. Simulation setting

To demonstrate the performance of the control system, Matlab/Simulink environment is employed, parameters of the drone is given as follows:  $m = 1.776 \text{ kg}$ ,  $g = 9.81 \text{ m/s}^2$ ,  $\mathbf{I} = \text{diag}(I_x, I_y, I_z) = \text{diag}(35, 35, 55) \times 10^{-4} \text{ kg.m}^2$ ,  $l = 0.225 \text{ m}$ ,  $I = 2.8 \times 10^{-6} \text{ kg.m}^2$ , Drag coefficient  $k_c = 5.567 \times 10^{-4} \text{ s}^{-1}$  [14]. The following disturbances affecting the system is provided:  $D = -5(0.8 \sin(0.7t) + 0.3 \sin(t) + 0.15 \sin(2t) + 0.055 \sin(\pi t / 2))$

#### 3.2. Simulation results

In this subsection, a comparison between the proposed controller and the other 2 controllers, which are the sliding mode controller (SMC) and the Backstepping-SMC (BSP) is provided. As can be seen in figures 1-3, the trajectory error controlled by RL-SMC is the smallest among the 3 controllers, the error along the x-axis and y-axis being approximately  $10^{-4}$  and along the z-axis being  $10^{-5}$ .

In comparison in figure 5, the other 2 controllers are affected by the disturbances and, therefore, can not converge to the origin. Regarding figure 4 yaw angles, the proposed controller also takes a shorter time to converge than the other 2, approximately 0.45 seconds. This study employs three performance indices to assess tracking performance in quantitation: ISE, IAE, and ITAE, as described in [15]. The ISEp and ISEa indices represent the integral squared errors for position and attitude, respectively. IAEp and IAEa are the integral absolute errors for position and attitude, respectively. The ITAEp and ITAEa indices, which measure the integral time-weighted absolute errors for position and attitude. Therefore, the following index table is obtained as in table 1. According to these performance indices shown in table 1, the position of the designed controller in this paper provides the smallest statistics, which show that the proposed controller can maintain its robustness and effectiveness.

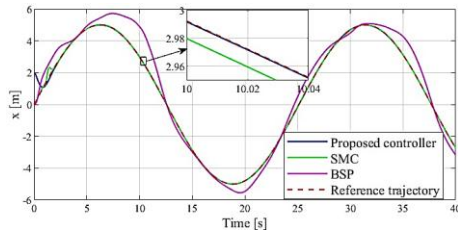


Figure 1. Response of x position.

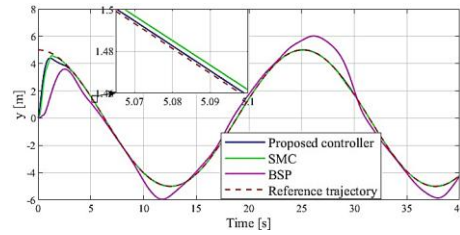


Figure 2. Response of y position.

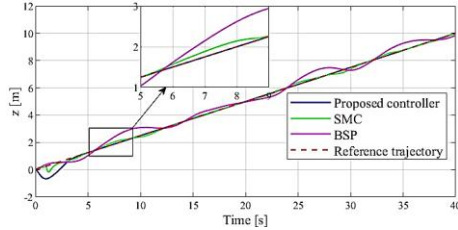


Figure 3. Response of z position.

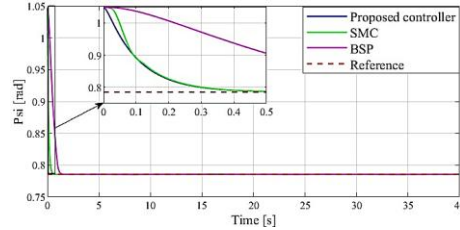


Figure 4. Response of yaw angle.

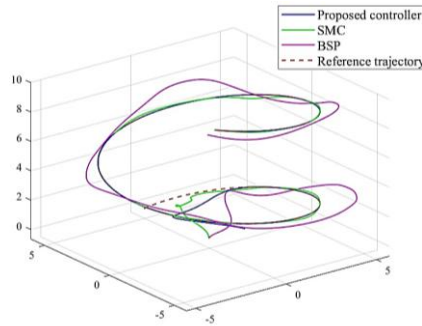


Figure 5. Trajectory in 3D space.

Table 1. The performance index table of the 3 controllers.

INDEX	$ISE_p$	$ISE_a$	$IAE_p$	$IAE_a$	$ITAE_p$	$ITAE_a$
SMC	0.081	14.923	0.559	8.795	7.535	71.995
BSP	2.179	59.801	7.558	55.624	137.48	1027.12
RL-SMC	1.923	<b>10.768</b>	0.62	<b>6.029</b>	0.206	<b>8.306</b>

#### 4. CONCLUSIONS

In this article, the model of the UAV with disturbances is demonstrated, based on the combination of the RL and SMC for the position loop, the simulated results on Matlab-Simulink proved that the proposed controller can effectively control the UAV tracking the reference trajectory in the appearance of disturbances.

Given the constraints of time and financial resources, a real-life quadrotor system has not been built. As a result, future work could focus on developing a practical model to validate the proposed controllers and improve the control algorithm for the inner loop as well.

#### REFERENCES

- [1]. Mohsan, Syed Agha Hassnain, et al. "Unmanned aerial vehicles (UAVs): Practical aspects, applications, open challenges, security issues, and future trends." Intelligent Service Robotics, 16.1, 109-137, (2023).
- [2]. Liu, Hui, et al. "Reinforcement learning-based tracking control for a quadrotor unmanned aerial vehicle under external disturbances." International Journal of Robust and Nonlinear Control, 33.17, 10360-10377, (2023).

- [3]. Li, Bo, et al. "Fixed-time integral sliding mode control of a high-order nonlinear system." *Nonlinear Dynamics*, 107, 909-920, (2022).
- [4]. Qi, Wenhai, Guangdeng Zong, and Wei Xing Zheng. "Adaptive event-triggered SMC for stochastic switching systems with semi-Markov process and application to boost converter circuit model." *IEEE Transactions on Circuits and Systems I: Regular Papers*, 68.2, 786-796, (2020).
- [5]. Yong, Jiongmin, et al. "Dynamic programming and HJB equations." *Stochastic controls: Hamiltonian systems and HJB equations*, 157-215, (1999).
- [6]. Evans, L. C., and M. R. James. "The Hamiltonian–Jacobi–Bellman equation for time-optimal control." *SIAM journal on control and optimization*, 27.6, 1477-1489, (1989).
- [7]. Werbos, Paul. "Approximate dynamic programming for real-time control and neural modeling." *Handbook of intelligent control*, (1992).
- [8]. Vamvoudakis, Kyriakos G., and Frank L. Lewis. "Online actor-critic algorithm to solve the continuous-time infinite horizon optimal control problem." *Automatica*, 46.5, 878-888, (2010).
- [9]. Ma, Zhiqiang, Panfeng Huang, and Yuxin Lin. "Learning-based sliding-mode control for underactuated deployment of tethered space robot with limited input." *IEEE Transactions on Aerospace and Electronic Systems*, 58.3, 2026-2038, (2021).
- [10]. Modares, Hamidreza, Mohammad-Bagher Naghibi Sistani, and Frank L. Lewis. "A policy iteration approach to online optimal control of continuous-time constrained-input systems." *ISA transactions*, 52.5, 611-621, (2013).
- [11]. Tan, Jian, and Shijun Guo. "Backstepping control with fixed-time prescribed performance for fixed wing UAV under model uncertainties and external disturbances." *International Journal of Control*, 95.4, 934-951, (2022).
- [12]. Wen, Guoxing, et al. "Optimized backstepping tracking control using reinforcement learning for quadrotor unmanned aerial vehicle system." *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, 52.8, 5004-5015, (2021).
- [13]. Xu, Shihao, et al. "Reinforcement-learning-based tracking control with fixed-time prescribed performance for reusable launch vehicle under input constraints." *Applied Sciences*, 12.15, 7436, (2022).
- [14]. Nguyen, H. S., et al. "Advanced Motion Control of a Quadrotor Unmanned Aerial Vehicle based on Extended State Observer." *International Conference on System Science and Engineering*. IEEE, (2023).
- [15]. Liu, K., Wang, R., Zheng, S., Dong, S., & Sun, G. "Fixed-time disturbance observer-based robust fault-tolerant tracking control for uncertain quadrotor UAV subject to input delay". *Nonlinear Dynamics*, 107(3), 2363-2390, (2022).
- [16]. Chen, Fuyang, et al. "Robust backstepping sliding-mode control and observer-based fault estimation for a quadrotor UAV." *IEEE Transactions on Industrial Electronics*, 63.8, 5044-5056, (2016).
- [17]. Tan, Lingwei, et al. "Super-twisting sliding mode control with defined boundary layer for chattering reduction of permanent magnet linear synchronous motor." *Journal of Mechanical Science and Technology*, 35: 1829-1840, (2021).

## TÓM TẮT

### Điều khiển trượt dựa trên học tăng cường trong điều khiển bám quỹ đạo của máy bay không người lái dưới các nhiễu bất định

Bài báo đề xuất phương pháp điều khiển trượt (SMC) dựa trên học tăng cường (RL) để điều khiển bám quỹ đạo của quadrotor UAV (QUAV) dưới tác động của nhiễu bên ngoài. Đầu tiên, một bộ điều khiển trượt dựa trên học tăng cường actor-critic (actor-critic RL) được giới thiệu để giải quyết bài toán điều khiển tối ưu trong điều kiện không có nhiễu. Tiếp theo, mô phỏng trong môi trường có nhiễu được thực hiện nhằm chứng minh tính bền vững của bộ điều khiển được đề xuất. Phân tích lý thuyết cho thấy sai số vị trí và góc của UAV hội tụ về một miền đặt trước, trong khi sai số ước lượng của mạng actor-critic được giới hạn cuối cùng một cách thống nhất (uniformly ultimately bounded - UUB). Cuối cùng, một phân tích so sánh mô phỏng số giữa bộ điều khiển đề xuất, bộ điều khiển trượt truyền thống và bộ điều khiển trượt kết hợp kỹ thuật Backstepping (BSP) được thực hiện để làm rõ các ưu điểm và hiệu suất cải thiện của SMC dựa trên RL.

**Từ khóa:** Học tăng cường (RL); Cấu trúc Actor/Critic; Điều khiển trượt (SMC); Điều khiển tối ưu; Máy bay không người lái 4 cánh (QUAV).