

Data-Driven Reinforcement Learning Control for Quadrotor Systems

Ngoc Trung Dang¹ and Phuong Nam Dao^{2,*}

¹ Faculty of Electrical Engineering, Hanoi University of Technology, Thai Nguyen, Vietnam

² School of Electrical and Electronic Engineering, Hanoi University of Science and University, Hanoi, Vietnam
Email: trungcsktd@tnut.edu.vn (N.T.D.); nam.daophuong@hust.edu.vn (P.N.D.)

* Corresponding author

Abstract² This paper aims to solve the tracking problem and optimality effectiveness of an Unmanned Aerial Vehicle (UAV) by model-free data Reinforcement Learning (RL) algorithms in both sub-systems of attitude and position. First, a cascade UAV model structure is given to establish the control system diagram with two corresponding attitude and position control loops. Second, based on the computation of the time derivative of the Bellman function by two different methods, the combination of the Bellman function and the optimal control is adopted to maintain the control signal as time converges to infinity with the addition of a discount factor. Third, according to off policy technique, the two proposed model-free RL algorithms are designed for attitude and position sub-systems in UAV control structure with a discount factor, respectively. In particular, the designed algorithms not only solve the trajectory tracking problem but also guarantee the optimality performance. Finally an illustrative system is used to verify the performance of the proposed model-free data RL algorithms in the UAV control system.

Keywords² data Reinforcement Learning (RL), Unmanned Aerial Vehicles (UAVs), quadrotor, Approximate/Adaptive Dynamic Programming (ADP), model-free based control

I. INTRODUCTION

The tracking control problem of Unmanned Aerial Vehicles (UAVs) has been extensively studied in many real-world applications, such as cargo transportation, precise agriculture, military, etc. [1–8]. In practice, the position and attitude dynamics tend to be uncertain and perturbed, which necessitates the design of nonlinear controllers with strong robustness against dynamic uncertainties and external disturbances [4]. Several classical nonlinear control methods studied for UAVs are nowadays available in the literature, such as adaptive dynamic Sliding Mode Control (SMC) [4], finite-time SMC [9], etc. However, it should be noted that the backstepping technique laws have existed in almost all conventional nonlinear control designs [4, 9]. The majority of these methods is to handle the tracking performance of the entire UAV system while each subsystem guarantees stability. Therefore, the disadvantages of the

mentioned nonlinear control methods can be known not only in choosing the Lyapunov function candidate but also requiring complex computation. Moreover, the optimality problem has not been discussed in the unification with the tracking problem [4, 9].

Implementing the optimal control scheme for robotic systems, requires the approximate algorithms to solve Hamilton-Jacobi-Bellman (HJB) or Riccati equations, which are difficult to directly solve by analytical method. The development of Reinforcement Learning Control (RLC) and Approximate/Adaptive Dynamic Programming (ADP) theories have significant implications for developing optimal control problems in robotic control systems [10, 11]. In particular, RLC is approached by Actor/Critic consideration [10, 11] wherein through Neural Networks (NN) approximation and optimization solution, the learning algorithm was given to satisfy not only optimality effectiveness but also tracking problem. Nevertheless, the external disturbance and dynamic uncertainties were independently handled with the RL algorithm in [10, 11], which is considered the model-based method. Unlike the simultaneous learning in Actor/Critic in [10, 11], the work in [12–14] developed sequential training by Policy Iteration (PI) and Value Iteration (VI) for non-affine discrete-time systems and affine continuous-time systems, respectively. The convergence of optimal control and Bellman function is studied by considering the decreasing sequence to be limited. The data-driven PI in [14] was implemented with three successive phases for seeking Lipschitz admissible control and the optimal control policy. Moreover, the relation between PI and VI approaches was clarified in [12] and an extension of output PVI ADP strategy was given with state reconstruction in [18]. On the other hand, a data-driven learning algorithm was utilized to seek the output ADP for perturbed linear systems with system based on the sampling data, which was collected from two virtual systems [15]. It is different from the Policy technique considering the computation under the control signal being the control policy [10–15], the off-policy method was developed for model-free reinforcement learning with the addition of an observer [16]. An approach to handling model-free

requirements was introduced by using data-driven Actor/Critic Q learning for continuous linear systems [19]. On the other hand, to maintain the control signal as time comes to infinity a model-free Q-learning algorithm was applied for linear discrete-time systems with a discount factor in the cost function [20]. Furthermore, the constrained optimal problem was addressed by integrating the Barrier function into the inverse optimal control problem [21].

Recently, particularly relevant for this paper, a data RL control strategy for a UAV was proposed in [1] to achieve the unification of the optimality problem and trajectory tracking control. An appealing feature of data RL control is that the proposed control is developed for complete dynamic uncertainties [1]. Moreover, the fault-tolerant control problem was improved in the data RL algorithm [2]. However, the data RL technique was only developed for an attitude system of UAV and the cost function was mentioned but without a discount factor [1, 2, 6]. Additionally, a model-free RL control was also investigated for only an attitude system of UAV to address the actuator saturation by developing an additional saturation function and saturated sliding surface [3]. On the other hand, there have been RL strategies for multiple UAVs, such as for a position-subsystem [5], and an attitude subsystem [8]. Generally, the disadvantages of the recent references [10, 11] are to acquire the model parameters and to set the cost function being infinity, which means that the optimal control consideration is not satisfied. Thus, it is necessary to consider the addition of a discount factor in the cost function to develop the optimal control for UAV systems (see Table I).

TABLE I. THE PHYSICAL MEANINGS OF NOTATIONS

Notation	Meaning
\mathcal{F}^i	The inertial frame
\mathcal{F}^b	The Body-fixed frame
X	a vector X expressed in \mathcal{F}^i , where $j = i, b$
e_n^i	a vector in \mathcal{F}^i with the i^{th} element equal to 1 and 0 at the others
$0_{m,n}$	a $m \times n$ matrix with all 0 elements
R_b^i	The rotation matrix denotes a rotation from \mathcal{F}^b to \mathcal{F}^i
c, s	$\cos(*)$, $\sin(*)$
$\begin{matrix} 1, 2, \\ 3, 4 \end{matrix} \dot{z}$	The velocities of the four rotors in Fig. 1
l	Distance between center of mass and rotor (Fig. 1)
m	Mass of the quadrotor (Fig. 1)

The rest of the paper is structured as follows. In Section II, we briefly introduce the position and attitude models of UAVs and present a control design objective. In Section III, we further present the RLC for nonlinear systems with a discount factor and two model-free RL algorithms are developed for attitude and position-subsystems. In Section IV, the proposed algorithms are illustrated with a UAV control system to demonstrate their

performance. This article concludes with a summary of the findings in Section V.

II. PRELIMINARIES AND PROBLEM STATEMENTS

Fig. 1 depicts the kinematics and dynamics of a quadrotor with two coordinate frames to be considered as a Body-fixed frame \mathcal{F}^b and a NED inertial frame \mathcal{F}^i .

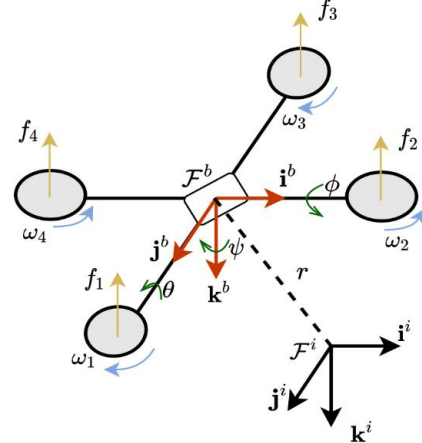


Fig. 1. Quadrotor model in NED coordinate

According to Ref. [1], the position vector of the quadrotor and the orientation angles are defined under the NED inertial frame, the Euler angles Roll-Pitch-Yaw, respectively. It should be noted that the Euler angles are bounded such that $|\psi| \leq \psi_{max}$, $|\theta| \leq \theta_{max}$ and $|\phi| \leq \phi_{max}$. On the other hand, the Euler angles vector is utilized to represent the following rotation matrix:

$$R = \begin{bmatrix} c_\phi c_\psi & c_\phi s_\psi & s_\phi \\ -s_\phi c_\psi & -s_\phi s_\psi & c_\phi \\ s_\psi & c_\psi & 0 \end{bmatrix} \quad (1)$$

In Fig. 1, the angular velocities vector of the quadrotor concerning coordinate \mathcal{F}^i is given by $\bar{\omega} = (p, q, r)^T$, which is expressed about Euler angles rate: $(\dot{\psi}, \dot{\theta}, \dot{\phi})^T$ as follows:

$$\bar{\omega} = \begin{bmatrix} 1 & 0 & s_\psi \\ 0 & c_\psi & c_\psi \\ 0 & s_\psi & c_\psi \end{bmatrix} \begin{bmatrix} \dot{\psi} \\ \dot{\theta} \\ \dot{\phi} \end{bmatrix} \quad (2)$$

According to [1], the dynamics model of a quadrotor can be depicted as:

$$m \ddot{z} = mg - R \bar{f} \quad (3)$$

where $\bar{f} \in \mathbb{R}^3$ indicates the total lift forces in the frame \mathcal{F}^i by $\bar{f} = [f_1, f_2, f_3, f_4]^T$, and is the torque in a frame \mathcal{F}^b (see Fig 1). Their relationship is described in Eq. (4):

$$\begin{aligned}
 & f \quad k_f \quad k_f \quad k_f \quad k_f \quad k_f \\
 T_x & \quad 0 \quad k_z \quad 0 \quad k_z \quad k_z \\
 T_y & \quad k_z \quad 0 \quad k_z \quad 0 \quad k_z \\
 T_z & \quad k_z \quad k_z \quad k_z \quad k_z \quad k_z
 \end{aligned} \quad (4)$$

with l to be the distance between the center of mass and the rotor and k_f, k_z, k_z are constant dynamics coefficients of force and torque respectively. Similarly, we define $u_x = \frac{1}{4} Z^2, u_y = \frac{1}{4} Z^2, u_z = \frac{1}{4} Z^2$. Consequently, from Eqs. (2) and (3), the dynamics equations can be written as:

$$m \ddot{r} = f R e_3 - mg e_3 \quad (5)$$

$$J \ddot{C} = T C \quad (6)$$

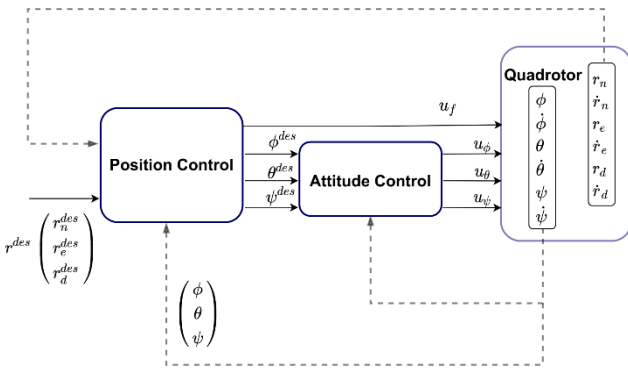


Fig. 2. Quadrotor control schematic

The objective is to seek an optimal controller for perturbed quadrotors to satisfy not only the tracking problem with a sophisticated trajectory but also minimize the discount factor-based cost function Eq. (8). Moreover, it should be noted that in the proposed control system (Fig. 2), the optimal control design approach is implemented for both systems in Eqs. (5) and (6) by observing system states with no prior knowledge of the object.

Remark 1. It is different from the classical trajectory tracking control objective to be solved by nonlinear control using Lyapunov stability theory [9], the control objective in this article requires not only the trajectory tracking effectiveness but also the minimization of the given cost function. Moreover, the optimal control problem is necessary to satisfy both systems (Fig. 2), which is unified with trajectory tracking requirements.

III. RL CONTROL DESIGNS

In this section, a novel control structure illustrated in Fig. 2 is proposed with two RL-based control loops to be Position Controller and Attitude Controller after separating the Quadrotor into two systems as mentioned in Eqs. (5) and (6). The Position Controller is considered as an outer loop controller to obtain the desired trajectory and its outputs are the desired Euler angles

$(r_n^{des}, r_e^{des}, r_d^{des})^T$. The Attitude Controller has a duty to track the desired Euler angles. The two RL-based controllers are investigated under the influence of exponential function and data collection to satisfy the complete uncertainty.

A. RL Control Design with Discount Factor

The RL technique is investigated for the following affine system:

$$\frac{d}{dt} J = F J + G(J) u(t) \quad (7)$$

The optimal control signals are designed to guarantee the minimization of the following infinite horizon cost function with a positive discount factor $\alpha > 0$.

$$V^*(t, u(t)) = \int_t^\infty e^{-\alpha(s-t)} U(J(s), u(s)) ds \quad (8)$$

where $U(J, u(s)) \triangleq J^T Q J + u^T R u$, Q and $R \in \mathbb{R}^{n \times n}$ are two positive-definite symmetric constant matrices. It can be seen that the addition of discount factor in cost function Eq. (8) leads to the existence of input signals vector as time converges to infinity.

Hence, the admissible control policy is not necessary to discuss as previous references [11]. The constraint set of control signals is only necessary to satisfy the finite cost function Eq. (8). Thanks to the autonomous property of affine model Eq. (7), the Bellman function obtained from Dynamic programming is the following static function:

$$V^*(J(t)) = \min_{u \in \mathcal{U}} \int_t^\infty U(J(s), u(s)) ds \quad (9)$$

The modified Hamiltonian function associated with discount factor is established by considering two different computation methods of the time derivative of the Bellman function $V^*(J(t))$. First, it is computed directly as:

$$\frac{d}{dt} V^*(J(t)) = \frac{V^*}{dt} \frac{dJ}{dt} = \frac{wV^*}{wJ} F(J) + G(J) u^w \quad (10)$$

where u is the optimal control signal.

The second method to compute the time derivative of the Bellman function $V^*(J(t))$ is developed in the view of the dynamic programming principle with the following static Bellman function Eq. (9):

$$\begin{aligned}
 V^*(t) &= \int_t^\infty e^{-\alpha(s-t)} U(J(s), u(s)) ds \\
 &= \int_t^\infty e^{-\alpha(s-t)} U(J(s), u(s)) ds + \int_t^\infty e^{-\alpha(s-t)} U(J(s), u(s)) ds \\
 &= \int_t^\infty e^{-\alpha(s-t)} U(J(s), u(s)) ds + e^{-\alpha(t-t)} V^*(t) \\
 &= \int_t^\infty e^{-\alpha(s-t)} U(J(s), u(s)) ds + V^*(t)
 \end{aligned} \quad (11)$$

Therefore, Eq. (11) implies that:

$$\frac{d}{dt} V^*(t) = -\gamma [U^*(t) + \int_t^{\infty} e^{-\gamma(s-t)} \dot{V}^*(s) ds] \quad (12)$$

Observing Eqs. (12) and (10) it implies that the static Bellman function $V^*(J(t))$ can be computed from the optimal control signal $u^*(t)$ after solving the following partial derivative equation:

$$U^*(J, t, u^*(t)) - \gamma \left(\frac{\partial}{\partial u} V^*(J, t) \right)^T F(J, G(J)u^*(t)) = 0 \quad (13)$$

To obtain the optimal control signal $u^*(t)$ from the Bellman function $V^*(J(t))$, based on the Dynamic Programming principle, it yields the following optimization problem:

$$V^*(J, t) = \min_{u \in \mathcal{U}} \int_t^{\infty} U^*(J, u(s)) ds + e^{-\gamma(t-\infty)} V^*(J, \infty) \quad (14)$$

As (14) implies the modified optimization problem is

$$\min_{u \in \mathcal{U}} [U^*(J, u(t)) - \gamma \left(\frac{\partial}{\partial u} V^*(J, t) \right)^T F(J, G(J)u(t))] = 0 \quad (15)$$

Denoting the modified Hamiltonian function to be associated with a discount factor γ as

$$H(J, u(t), V, \dot{V}) = \left(\frac{\partial}{\partial u} V^*(J, t) \right)^T F(J, G(J)u(t)) + U^*(J, u(t)) - \gamma \left(\frac{\partial}{\partial u} V^*(J, t) \right)^T F(J, G(J)u(t)) \quad (16)$$

where it implies that the optimal control is then achieved from Eq. (15) as

$$u^*(J) = \underset{u \in \mathcal{U}}{\operatorname{argmin}} H(J, u(t), V^*(J), \dot{V}^*(J)) = \frac{1}{2} R^{-1} G^T(\cdot) \left(\frac{\partial}{\partial u} V^*(J) \right) \quad (17)$$

Moreover, substituting the optimal control $u^*(J)$ Eq.(17) into Eq. (15) obtains the Partial Derivative Equation(PDE) as

$$\begin{aligned} H^*(J, u^*, V, \dot{V}) &= \left(\frac{\partial}{\partial u} V^*(J, t) \right)^T Q \left(\frac{\partial}{\partial u} V^*(J, t) \right) \\ &+ \frac{1}{4} \left(\frac{\partial}{\partial u} V^*(J, t) \right)^T G(J) R^{-1} G^T(J) \left(\frac{\partial}{\partial u} V^*(J, t) \right) \\ &- \gamma \left(\frac{\partial}{\partial u} V^*(J, t) \right)^T F(J, G(J)u^*(t)) = 0 \end{aligned} \quad (18)$$

However, it is impossible to analytically solve the PDE Eq. (18) to find the Bellman function from the optimal

control signal $u^*(J)$. Hence, the data driven algorithm is mobilized to seek the optimal control signal $u^*(J)$ in Sections B and C.

Remark 2. The discount factor is added to the cost function Eq. (8) for keeping the control signals as the tracking problem is satisfied. However, it is necessary to address the impact of the terms in Eq. (18), which establishes some modifications in actor-critic-based RL control design Section III. B and C.

B. Data-Driven PI Position Controller

In this section, after achieving the separation of the quadrotor model, the control design of each system is developed by the data driven RL technique as follows. The first step is to rewrite the position system described in Eq. (5):

$$\ddot{r} = m^{-1} k_f u_f R \hat{e}_3 + g \hat{e}_3 \quad (19)$$

$$m^{-1} k_f u_f \quad (20)$$

where the term can be known as an offset element and eliminated by determining a bias $m k_f^{-1} g \hat{e}_3$ to lift the quadrotor off the ground. To develop the RL technique as described in Section A, we employ the states vector $x_r = (r_n, \dot{r}_n, r_e, \dot{r}_e, r_d, \dot{r}_d)^T \in \mathbb{R}^6$. Thus, the position system Eq. (19) can be transformed as:

$$\dot{x}_r = A_r x_r + B_r u_f \quad (21)$$

where $A_r = \operatorname{diag}(a, a, a) \in \mathbb{R}^{6 \times 6}$, $a = (0, 1, 0)$ and $B_r = m^{-1} k_f \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}$. Assume that, the desired trajectory is bounded and is the Lipschitz function. On the other hand, defining the reference $x_r^{des} = (r_n^{des}, \dot{r}_n^{des}, r_e^{des}, \dot{r}_e^{des}, r_d^{des}, \dot{r}_d^{des})^T \in \mathbb{R}^6$ and $x_r^{des}(t)$ can be completely represented as $x_r^{des}(t) = A_{rd} x_r^{des}(t)$. Define $e_r = x_r - x_r^{des}$, Eq. (21) can be written as:

$$\dot{x}_r = \begin{pmatrix} a & 0 & 0 & 0 & 0 & 0 \\ 0 & a & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} x_r + \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} u_f \quad (22)$$

The tracking performance function is chosen as:

$$V_r(x_r(t)) = \int_t^{\infty} e^{-\gamma(t-s)} [x_r^T(t) Q x_r(t) + u_f^T(t) R u_f(t)] ds \quad (23)$$

where $Q = \begin{pmatrix} Q_{er} & 0_{6,6} \\ 0_{6,6} & 0_{6,6} \end{pmatrix}$ with and to be both positive

symmetric definite matrices. Note that, the term is added to penalize both tracking error cost and input control energy cost. It is easy to determine that although the desired trajectory is bounded, the control input does not converge to zero as it comes to infinity. According to the work in

Section III. A, the Position Controller can be deployed as described in the following Algorithm 1:

Algorithm 1 (Fig. 3): Data-driven PI Position Control

Step 1 (Initialization): Starting the stabilizing control policy and the disturbance term to satisfy the PE condition. Implementing the data collection and establishing the threshold ϵ_p

Step 2 (Policy Evaluation): For each control signal $u_p^i(X_p)$, solve simultaneously the $V_p^{i+1}(X_p)$ and $u_p^{i+1}(X_p)$ by the following equation:

$$\begin{aligned} & V_p^{i+1}(X_p(t+\theta)) - V_p^i(X_p(t)) \\ & \int_t^{t+\theta} X_p^T(Q_p X_p(\tau) + W_p u_p^i(X_p(\tau)))^T R_p^{-1} \dot{u}_p^i(X_p(\tau)) d\tau \\ & - \int_t^{t+\theta} \dot{V}_p^i(X_p(\tau)) d\tau - \int_t^{t+\theta} u_p^{i+1}(X_p(\tau))^T R_p^{-1} \dot{u}_p^i(X_p(\tau)) d\tau \\ & - 2 \int_t^{t+\theta} u_p^{i+1}(X_p(\tau))^T R_p^{-1} u_p^i(X_p(\tau)) d\tau \end{aligned} \quad (24)$$

Step 3: (Policy Improvement): Update the control policy $u_p^i(X_p)$ to $u_p^{i+1}(X_p)$, $i = i + 1$ and come back to Step 2 until $\|u_p^{i+1} - u_p^i\| < \epsilon_p$

After obtaining the control signal in the control structure (Fig.1), the reference of the attitude control scheme can be obtained (u_m, u_e, u_d) . According to Ref. [1], the desired yaw angle ψ^{des} can be chosen as zero and ϕ^{des} can be easily solved as follows:

$$\begin{aligned} u_{rp} &= \sqrt{u_m^2 + u_e^2 + (u_d - u_0)^2}, \psi^{des} = 0, \\ \phi^{des} &= \arcsin\left(\frac{u_m \sin(\psi^{des}) + u_e \cos(\psi^{des})}{u_{rp}}\right), \\ \theta^{des} &= \arctan\left(\frac{u_m \cos(\psi^{des}) - u_e \sin(\psi^{des})}{u_d - u_0}\right) \end{aligned} \quad (25)$$

C. Data-driven PI Attitude Controller

After the desired attitudes are obtained from Eq (25), the objective of the attitude controller in Fig.1 is to design the input signal for satisfying the optimal control problem. The model in Eq. (6) can be written as:

$$\ddot{X}_d = J^{-1} T J^{-1} C \ddot{X}_d \quad (26)$$

Based on the states vector $[\phi, \dot{\phi}, \theta, \dot{\theta}, \psi, \dot{\psi}]^T$, the attitude control scheme (Fig.1) is similar to the Position Controller in subsection A. According to the model in Eq. (6), the attitude model can be rewritten as:

$$\dot{X}_d = \begin{bmatrix} \dot{\phi} \\ \ddot{\phi} \\ \dot{\theta} \\ \ddot{\theta} \\ \dot{\psi} \\ \ddot{\psi} \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} X_d + \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad (27)$$

It implies that the attitude control design can be developed in the following algorithm:

Algorithm 2 (Fig. 3): Data-driven PI Position Control

Step 1(Initialization): Choosing the stabilizing control policy and the disturbance term $w_e(t)$, the threshold to satisfy the PE condition, and collecting the data.

Step 2(Policy Evaluation): For each control signal $u^i(X)$, solve simultaneously the $V^{i+1}(X)$ and $u^{i+1}(X)$ by the following equation:

$$\begin{aligned} & V^{i+1}(X(t+\theta)) - V^i(X(t)) \\ & \int_t^{t+\theta} X^T(Q(X(\tau)) + W u^i(X(\tau)))^T R^{-1} \dot{u}^i(X(\tau)) d\tau \\ & - \int_t^{t+\theta} \dot{V}^i(X(\tau)) d\tau - \int_t^{t+\theta} u^{i+1}(X(\tau))^T R^{-1} \dot{u}^i(X(\tau)) d\tau \\ & - 2 \int_t^{t+\theta} u^{i+1}(X(\tau))^T R^{-1} u^i(X(\tau)) d\tau \end{aligned} \quad (28)$$

Step 3(Policy Improvement): Update the control policy $u^i(X)$ to $u^{i+1}(X)$, $i = i + 1$ and come back to Step 2 until $\|u_p^{i+1} - u_p^i\| < \epsilon_p$.

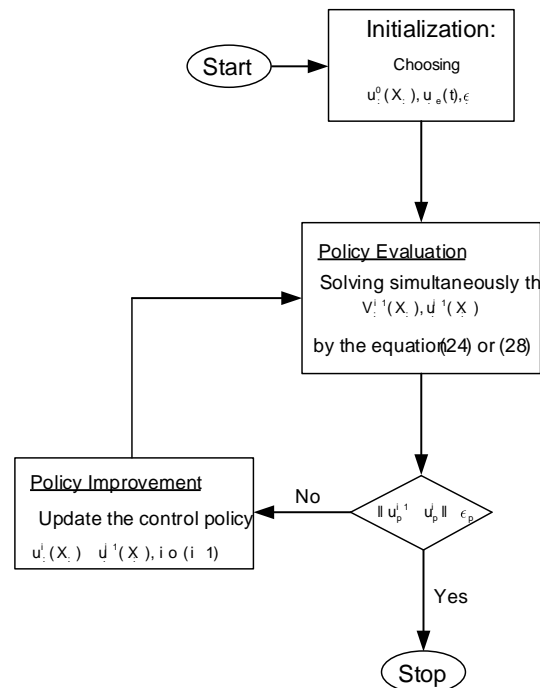


Fig. 3. The flowchart of Algorithm 1 or 2

On the other hand, it is worth emphasizing that the training time is considered from the initial time to the time of obtaining the optimal controllers by Algorithms 1 and 2. Furthermore, due to the dynamic uncertainties in models in Eq. (5) and (6), the proposed Algorithms 1 and 2 are developed by data collection of input and state signals in the practical system, which are employed to solve the Eqs (24) (28). However, it is worth emphasizing that the existence of root Eqs (24) (28), (Algorithm 1), (Algorithm 2) requires the Persistence of Excitation (PE) condition as shown in [10, 11].

IV. SIMULATION RESULTS

Consider a perturbed quadrotor with the parameters as follows: $m = 2.0(\text{kg})$, $k_w = 1(\text{N}\cdot\text{s}^2)$, and input disturbances in force are $0.1 \sin(\omega t)$ (Nm).

The desired trajectory is a spiral trajectory: $r^{des}(t) = [2\sin(\omega t), 2\cos(\omega t), 0.8t]^T$ where $\omega = 0.5$.

Initially, at the stage of collecting data, two PD (Proportion Derivative) controllers are implemented for both position and attitude. These non-optimal controllers are tuned manually to keep the quadrotors stable with the position and attitude errors illustrated in Figs. 4 and 5.

Additionally, to guarantee the PE conditions, noises with $u_{de} = \frac{500}{m} \cdot 0.002 \sin(\omega_m t)$ (ω_m is a frequency chosen in range randomly), are added to the position and attitude controllers respectively. The next stage is to apply the two algorithms in the previous sections after obtaining data. The algorithms' parameters are chosen as:

$$Q_e = 100I_6, R_p = I_3, Q_{de} = 100I_6, R_4 = I_3, \sigma = 0.01, T_{step} = 0.01$$

The activation functions of the critic and actor neural networks are second-order polynomials and first-order polynomials respectively. It can be seen that the convergence of the weights Algorithms 1 and 2 is shown in Fig. 4. Moreover, it leads to the tracking problem satisfied as shown in Fig. 5.

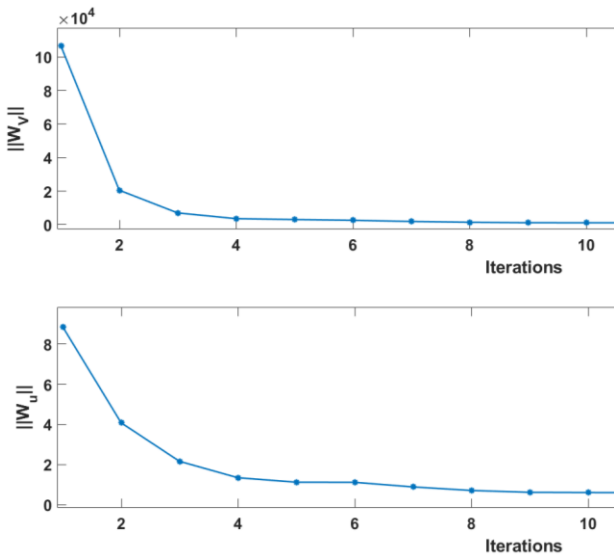


Fig. 4. The weight convergence of the learning stage

On the other hand, Fig. 6 shows the comparisons between the evolutions of tracking error depending on the discount factor. It can be seen that if the discount factor increases then the convergence speed of tracking error increases. However, it follows that the fluctuation will increase as the discount factor increases (Fig. 6). Additionally, unlike the classical controllers like Fuzzy

and PID only considering the position control problem, the proposed Algorithms 1 and 2 develop not only the trajectory tracking control performance but also model free RL strategy in the presence of dynamic uncertainties. Moreover, it can be seen that the advantage of these proposed methods is the extension of adding the discount factor and two RL algorithms, which have not been considered in the recent references [10, 11].

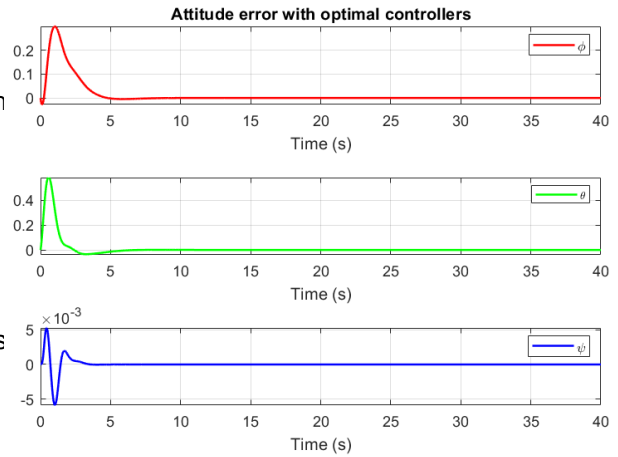


Fig. 5. The tracking of orientation angles

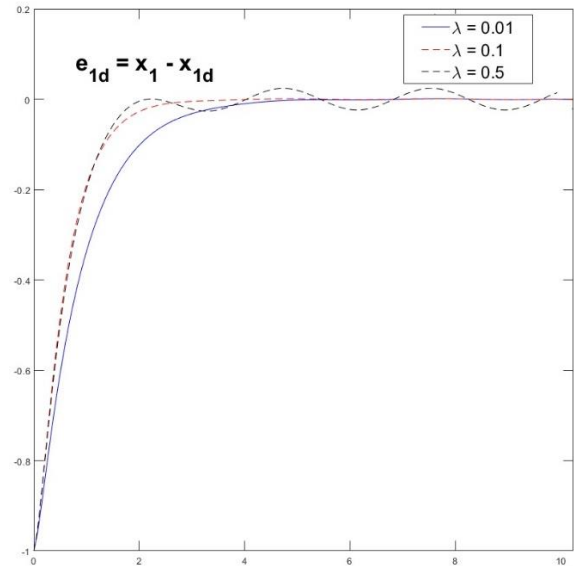


Fig. 6. The tracking error depends on the discount factor

V. CONCLUSION

This paper has designed the model-free data RL strategies for both attitude and position systems in cascade UAV control structure to achieve the unification of trajectory tracking problem and optimality purpose. The main idea is to establish the Off-Policy RL algorithm with a discount factor to satisfy the existence of control signal as time comes to infinity and obtain the model-free consideration without UAV model knowledge. Moreover, data collection and computation techniques are considered to achieve simultaneously the optimal value function and optimal control policy. Finally, an illustrative system is employed to validate the effectiveness of the proposed

