

# Two-Stage Strategy to Achieve a Reinforcement Learning-Based Upset Recovery Policy for Aircraft

Huanhui Cao, Weifeng Zeng, Hantao Jiang, Hao Hu, Chaoran Li, Wenjie Lu, and Hao Xiong\*

*School of Mechanical Engineering and Automation*

*Harbin Institute of Technology*

Shenzhen, China

xionghao@hit.edu.cn

**Abstract**—Aircraft upset situations are the highest risk to civil aviation. Thus, a reliable upset recovery policy is necessary for aircraft. In this paper, a two-stage strategy to achieve a reinforcement learning (RL)-based upset recovery policy that takes time of recovery and loss of altitude into account is proposed for aircraft to recover from an arbitration upset situation to level flight. Based on the proposed two-stage strategy and Twin Delayed Deep Deterministic Policy Gradient (TD3) algorithm, an algorithm to achieve a TD3-based upset recovery policy for aircraft is developed. Experiments are conducted based on X-Plane 11 to evaluate the effectiveness of the proposed two-stage strategy and the performance of the achieved upset recovery policy in stall recovery and spin recovery.

**Keywords**—aircraft upset recovery, reinforcement learning, twin delayed deep deterministic policy gradient, pre-train, fine-tuning

## I. INTRODUCTION

Aircraft upset situations are the highest risk to civil aviation from decades ago to now [1]. To this end, several researchers have worked on addressing the aircraft upset issue [2]–[4]. Yildiz et al. [2] proposed a finite-state conditional switching structure to achieve recovery from loss-of-control conditions. To deal with stall recovery, Cunis et al. [3] developed a Loss of Altitude (LOA) minimizing approach based on economic model predictive control. By solving a trajectory optimization problem via the direct multiple shooting method, aircraft spin recovery was addressed in [4]. Although the above-mentioned approaches can address nonlinear dynamics of aircraft, these approaches may fail to handle the high complexity of upset situations [5].

Reinforcement Learning (RL) is an approach that can deal with high complexity problems without an explicit model of the problem [6], [7]. RL is appropriate to be applied to develop an upset recovery policy for aircraft suffering complex upset situations [5]. Dutoi et al. combined robust control and RL to address the spin recovery problem [8]. The spin recovery performance of the achieved policy can outperform the policy obtained based on skilled pilots. Nonetheless, this study compressed the action space to ease the computation burden, leading to a limitation on agile spin recovery. Kim et al. [5] developed an RL-based recovery policy including nine actions for unusual attitude recovery and 27 actions for

angular rate arrest for a stable flat spin. By comparing to the optimal solution, the superiority of the performance of the RL-based upset recovery policy was verified. Zhu et al. [9] applied deep deterministic policy gradients (DDPG) and deep Q-network (DQN) to achieve spin recovery policies and proposed an exploring mechanism that dynamically selects between deterministic and stochastic exploration for DDPG. Nonetheless, this study did not consider the uncertainties of the aircraft and the environment.

This paper proposes a two-stage strategy to achieve an RL-based upset recovery policy for aircraft. The major contributions of this paper are as follows.

- A two-stage strategy, which includes a pre-train stage and a fine-tuning stage, to achieve an RL-based upset recovery policy is proposed for aircraft. The achieved upset recovery policy can address different upset situations of an aircraft, rather than stall or spin only.
- Experiments are conducted based on a Cessna172SP aircraft in X-Plane 11 to validate the effectiveness of the proposed two-stage strategy and the performance of the achieved upset recovery policies.

The rest of this paper is organized as follows. Section II demonstrates the preliminaries, including typical upset situations of aircraft, loss of altitude, reinforcement learning, and Twin Delayed Deep Deterministic Policy Gradient (TD3) algorithm. In section III, a two-stage strategy to achieve an RL-based upset recovery policy is presented. Section IV presents experiments to illustrate the proposed two-stage strategy and validate the effectiveness of the proposed two-stage strategy. Finally, section V summarizes this paper.

## II. PRELIMINARIES

### A. Upset Situations of Aircraft

An upset situation of aircraft refers to an abnormal mode of the nonlinear dynamics that shows significantly altered steady-state responses [3], such as stall and spin.

**Stall** [10]. In aerodynamics and aviation, stall is a condition such that lift begins to decrease suddenly if the angle of attack of aircraft increases beyond a certain point. Stall of an aircraft is shown in Fig. 1.

**Spin** [11]. Spin is a nonlinear post-stall phenomenon in which an aircraft develops a high rotational-rate and descends

\*Corresponding author.

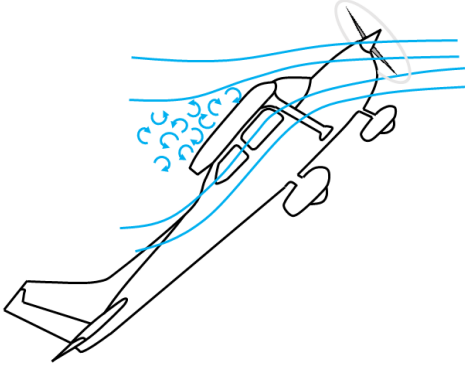


Fig. 1. Stall of an aircraft

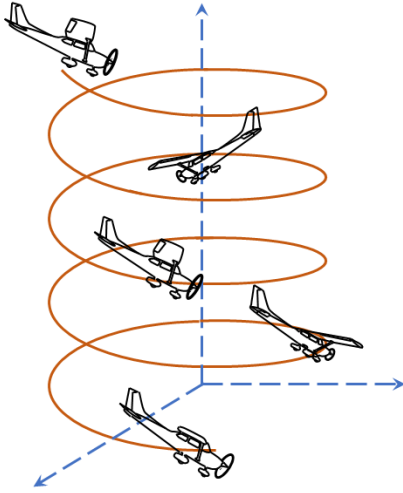


Fig. 2. Spin of an aircraft

almost vertically in a helical trajectory, as shown in Fig. 2. For an aircraft in spin, one wing is stalled more than the other. The more stalled wing is with less lift and more drag than the other, leading to the auto-rotation and subsequent rapid descent of the aircraft.

### B. Loss of Altitude

LOA [3] is a significant performance index for an upset recovery policy. LOA can be exploited to enlarge the operational envelope of aircraft, particularly at low altitudes. The LOA has been used in several previous researchers [12], [13] to evaluate control policies of aircraft. In this study, the LOA is a significant index to measure the performance of the achieved upset recovery policy.

### C. Reinforcement Learning and Twin Delayed Deep Deterministic Policy Gradient Algorithm

Reinforcement learning studies the paradigm of an agent interacting with the environment aiming to learn policies that maximize accumulated rewards. At time step  $t$ , the agent

selects an action  $a \in \mathcal{A}$  based on the current state  $s \in \mathcal{S}$  with respect to its policy  $\pi : \mathcal{S} \mapsto \mathcal{A}$ . The agent receives a reward  $r$  and the state transfers to a new state  $s'$ . The agent aims to maximize the accumulated rewards  $R_t = \sum_{i=t}^T \gamma^{i-t} r(s_i, a_i)$ , where  $\gamma$  is a discount factor.

TD3 algorithm [14] is an RL algorithm proposed for agents with continuous states and actions based on an actor-critic architecture. TD3 applies three approaches to address the overestimation issue of the critic network.

**Smooth target policy.** In order to reduce the variance caused by over-fitting, TD3 adds a clipped random noise to the value of the target critic network. After adding the clipped noise, the value of the target action network is clipped to lie in a feasible action range.

**Apply a pair of critic networks.** TD3 includes two critic networks. In the Bellman error loss functions, we should take the smaller value of these two target critic networks. In this way, the overestimation of the critic network is alleviated.

**Delay actor updates.** TD3 updates the actor network less frequently than the critic network. The actor network is updated after a certain number of time steps, while the critic network is updated at every time step.

## III. REINFORCEMENT LEARNING-BASED UPSET RECOVERY POLICY ACHIEVED

In this section, a two-stage strategy to achieve an RL-based upset recovery policy is proposed for aircraft.

### A. Problem Formulation

In aerodynamics and aviation, upset situations can refer to a variety of abnormal situations. An upset recovery policy is expected to recover aircraft from an upset situation with minimum LOA within the shortest possible time in practice. Assume that an aircraft can adjust its state by aileron, elevator, rudder, and throttle, as shown in Fig. 3. According to studies of multi-objective reinforcement learning [15], The problem addressed by an upset recovery policy can be expressed as

$$\begin{aligned} \min \quad & \{w_{LOA} \|LOA(a_i)\| + \sum_{m=1}^M \omega_m \|s_m(a_i) - s_m^*\|\} \\ \text{s.t.} \quad & a_i^L \leq a_i \leq a_i^H, \quad i = 1, 2, 3, 4 \end{aligned} \quad (1)$$

where  $a_i (i = 1, 2, 3, 4)$  denotes the action of the elevator, aileron, rudder, throttle of an aircraft, respectively.  $a_i^L$  and  $a_i^H$  are the lower bound and the upper bound of  $a_i$ , respectively.  $s_m(a_i) (m = 1, 2, \dots, M)$  denotes the state of the aircraft.  $s_m^*$  denotes a target state.  $LOA(a_i)$  denotes the LOA of the aircraft.  $w_m (m = 0, 1, \dots, M)$  are the weight coefficients of the state  $s_m(a_i)$ .

For an aircraft, an upset recovery policy is not supposed to control the aircraft to certain locations usually, so latitude and longitude can be ignored in the study of upset recovery policies. We can define the state space of an upset recovery policy as  $\mathcal{S} = \{\omega, \kappa, \xi, p, q, r, h, v\}$  and has  $s_m(a_i) \in \mathcal{S}$ .  $\omega, \kappa, \xi, p, q, r, h, v$  denote the pitch, roll, yaw, pitch angular velocity, roll angular velocity, yaw angular velocity, altitude, and airspeed of the aircraft, respectively.

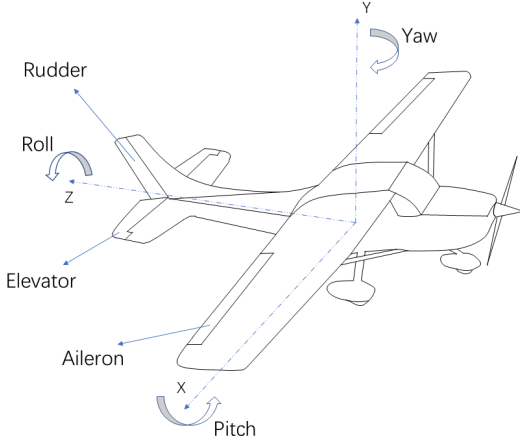


Fig. 3. The control surfaces of an aircraft

### B. Pre-train Stage and Fine-tuning Stage to Achieve an RL-Based Upset Recovery Policy

The dynamics of an aircraft is with complexity and uncertainty caused by aerodynamics, weather, fuel, and payload. Upset recovery policies based on robust control and the knowledge of dynamics of the aircraft cannot achieve optimal performance [5]. RL is a promising approach to achieve upset recovery policies in the sense of complexity and uncertainty addressing. In view of the uncertainties of aircraft such as aerodynamics, weather, fuel, and payload, a two-stage strategy to achieve an RL-based upset recovery policy is proposed, as shown in Fig. 4. The proposed two-stage strategy includes a pre-train stage and a fine-tuning stage.

**Pre-train stage.** In the pre-train stage, a general upset recovery policy is learned based on RL for an aircraft with parameters (e.g., aerodynamics, weather, fuel, and payload) disturbed in certain ranges randomly. The general upset recovery policy is sub-optimal but is with the adaptability to aircraft with parameter uncertainties.

**Fine-tuning stage.** In the fine-tuning stage, a specific upset recovery policy is learned via RL by the current aircraft whose aerodynamics, fuel, and payload are determined, based on the general upset recovery policy. Initializing the general upset recovery policy achieved in the pre-train stage as the initial specific upset recovery policy for the current aircraft can speed up the training of the specific upset recovery policy.

### C. Reward Function

The design of the reward function is crucial for RL since it affects not only the convergence but also the quality of the convergent point of the RL. Inspired by the idea of reward shaping [16], we define the reward function of an RL-based upset recovery policy as

$$R(s, a, s') = T(s') + P(s) + C(s, a) \quad (2)$$

where  $T(s')$  is a termination reward.  $P(s)$  is a state-related reward.  $C(s, a)$  is an attitude-related reward.

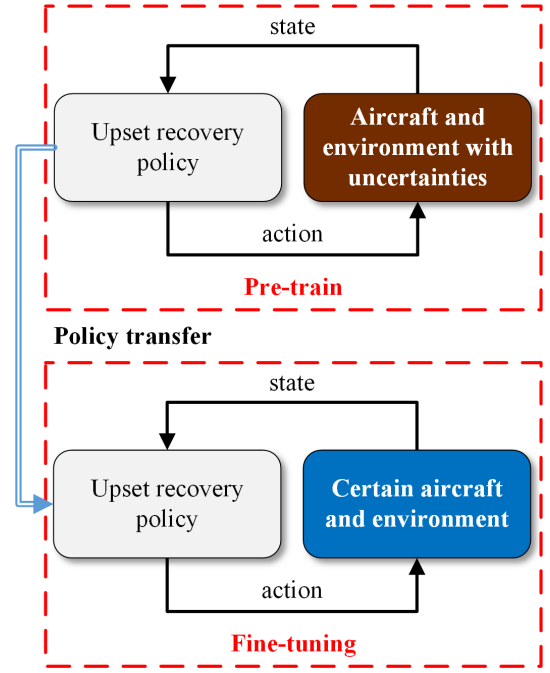


Fig. 4. Two-stage to achieve an RL-based upset recovery policy.

The termination reward is the reward obtained when an upset recovery policy recovers an aircraft successfully or makes the aircraft crashed eventually. The termination reward is defined as

$$T(s') = \begin{cases} T_1, & \text{if the aircraft is crashed} \\ T_2, & \text{if success to recover the aircraft} \\ 0, & \text{else} \end{cases} \quad (3)$$

where  $T_1$  is a negative reward while  $T_2$  is a positive reward.

The state-related reward  $P(s)$  is expressed as

$$P(s) = w_{LOA} ||LOA^n|| + \left( \sum_{m=1}^8 w_m ||s_m^n - s_m^{*,n}|| \right) \quad (4)$$

where  $s_m^n$ ,  $s_m^{*,n}$ ,  $LOA^n$  are the normalization of state  $s_m$ , target state  $s_m^*$ , LOA of an aircraft, respectively.  $w_m$  ( $m = 1, \dots, 8$ ) is a non-positive weight coefficient of the difference between  $s_m^n$  and  $s_m^{*,n}$ .  $w_{LOA}$  is a non-positive weight coefficient of LOA.

$C(s, a)$  depends on the attitude and angular velocity of an aircraft. If the angular velocity tends to decrease the difference between the attitude of the aircraft and the target attitude, a positive reward  $c_1$  will be received. Otherwise, a negative reward  $c_2$  will be received.

### D. Algorithm to Achieve a TD3-Based Upset Recovery Policy

Based on the two-stage strategy to achieve an RL-based upset recovery policy presented in section III-B and the reward function presented in section III-C, we can develop an algorithm to achieve TD3-based upset recovery policy by applying

the TD3 algorithm in both the pre-train stage and in the fine-tuning stage. The pseudo-code of the TD3 algorithm applied to achieve upset recovery policy is shown in Algorithm 1.

Although the TD3 algorithm is used in both the pre-train stage and in the fine-tuning stage, applications of the TD3 algorithm are different from two aspects - initialization and parameter uncertainty. For the initialization aspect, the upset recovery policy is initialized randomly in the pre-train stage while, in the fine-tuning stage, the upset recovery policy is initialized to the general upset recovery policy achieved in the pre-train stage. For the parameter uncertainty aspect, the environment is with parameters distributed randomly in certain ranges in the pre-train stage, but the environment is with parameters determined by the current aircraft.

#### IV. EXPERIMENTS

##### A. Environment Setups

To evaluate the proposed two-stage strategy to achieve upset recovery policy, we conduct experiments based on X-Plane 11, a flight simulation environment. X-Plane 11 is a professional flight simulation environment, involving various types of aircraft ranging from general aircraft to commercial aircraft. X-Plane has been used by several researchers [17]–[19] to study the design, control, and guidance of aircraft. A Cessna172SP is included in experiments to study the performance of the achieved upset recovery policies, as shown in Fig. 5. Moreover, the wind condition in the experiments is adjusted to achieve a general upset recovery policy in the pre-train stage.

X-Plane 11 is integrated with TensorFlow-gpu 2.4 in this study. A personal computer with 64-bit Windows 10 operation system, 32 gigabytes memory, a 3.2 GHz Core i7 CPU, and an RTX3070 GPU is used to validate the proposed two-stage strategy and the performance of achieved upset recovery policies.

##### B. Training Setups

With the above-mentioned environment set up, we set the TD3 algorithm included in the TD3-based upset recovery policy as follows. The learning rate is  $1.0 \times 10^{-4}$ . The discount factor  $\gamma$  is set to 0.99. The actor network has three hidden



Fig. 5. Aircraft included in experiments - Cessna172SP

---

#### Algorithm 1: Algorithm to achieve a TD3-based upset recovery policy for aircraft

---

**Input:** Empty replay buffer  $\mathcal{D}$

---

```

1 if pre-train stage then
2   Initialize actor network  $\theta$ , critic networks  $\phi_1$  and  $\phi_2$ ;
3   Initialize target networks:
4      $\theta_{targ} \leftarrow \theta$ ,  $\phi_{targ,1} \leftarrow \phi_1$ ,  $\phi_{targ,2} \leftarrow \phi_2$ ;
5   Set parameters of aircraft and environment
6     randomly in certain ranges;
7
8 if fine-tuning stage then
9   Load actor network  $\theta$ , critic networks  $\phi_1$  and  $\phi_2$ ,
10    and target networks  $\theta_{targ}$ ,  $\phi_{targ,1}$ ,  $\phi_{targ,2}$ 
11    achieved in pre-train stage;
12   Set parameters according to current aircraft and
13    environment;
14
15 for episode=1 to  $M$  do
16   Receive initial observation state  $s_{ini}$ ;
17   Normalize the initial state  $s_{ini}$ ;
18   for  $t=1$  to  $T$  do
19     Select action with exploration noise
20      $a \sim \pi_\theta(s) + \epsilon$ , where  $\epsilon \sim \mathcal{N}(0, \delta)$ ;
21     Convert action  $a$  into the aircraft maneuver  $a_m$ ;
22     Execute  $a_m$ , and obtain reward  $r$  and new state  $s'$ ;
23     Normalize the state  $s'$ ;
24     Store transition  $(s, a, r, s', d)$  in  $\mathcal{D}$ ;
25     Sample a mini-batch of transitions
26      $\mathcal{B} = \{(s, a, r, s', d)\}$  from  $\mathcal{D}$ ;
27     Compute action  $a'(s') =$ 
28        $clip(\pi_{\theta_{targ}}(s') + clip(\epsilon, -c, c), a_{Low}, a_{High})$ ;
29     Compute the value of target critic network
30      $y(r, s', d) =$ 
31        $r + \gamma(1 - d) \min_{i=1,2} Q_{\phi_{i,targ}}(s', a')$ ;
32     Update critic network by one step of gradient
33     ascent using:
34      $\nabla_{\phi_i} \frac{1}{|\mathcal{B}|} \sum_{(s,a,s',r,d) \in \mathcal{B}} (Q_{\phi_i}(s, a) - y(r, s', d))^2$ ;
35
36   if  $t \bmod p_{delay} = 0$  then
37     Update actor network by one step of
38     gradient ascent using:
39      $\nabla_{\theta} \frac{1}{|\mathcal{B}|} \sum_{s \in \mathcal{B}} Q_{\phi_1}(s, \pi_\theta(s))$ ;
40     Update target networks:
41      $\phi_{targ,i} \leftarrow \tau \phi_{targ,i} + (1 - \tau) \phi_i \quad i = 1, 2$ 
42      $\theta_{targ} \leftarrow \tau \theta_{targ} + (1 - \tau) \theta$ 

```

---

layers that have 64, 64, and 32 units, respectively. The critic network has two hidden layers with 64, 64 units, respectively. The time step  $T$  is defined as 500 and the minibatch size  $B$  is 64. The number of episodes  $M$  is 800 in the pre-train stage and is 50 in the fine-tuning stage. The feasible ranges of the states of the aircraft included in experiments are listed in Table I. The parameters of the reward function used in experiments are presented in Table II.

Table I. The feasible ranges of the states of aircraft

variable	value	variable	value
$\omega, \kappa$	$[-180^\circ, 180^\circ]$	$\xi$	$[0^\circ, 360^\circ]$
$p, q, r$	$[-90^\circ/s, 90^\circ/s]$	$h$	$[0\text{ m}, 5000\text{ m}]$
$v$	$[0\text{ m/s}, 120\text{ m/s}]$		

Table II. The parameters of reward function

variable	value	variable	value
$T_1$	-1000	$T_2$	500
$p_0$	-1.0	$p_1, p_2, p_3, p_7$	-0.5
$p_4, p_5, p_6, p_8$	-0.1	$c_1$	1.0
$c_2$	-1.0		

### C. Training of Upset Recovery Policies

Based on the environment set up and training set up presented above, we conduct the training of upset recovery policies according to the proposed two-stage strategy. In the pre-train stage, it is assumed that the fuel and payload of aircraft and the weather are not determined. Thus, the weight of the aircraft distributes randomly from 781 kilograms to 1160 kilograms and wind speed distributes randomly from 0 knots to 60 knots. In the fine-tuning stage, it is assumed that the weight of aircraft and the weather are determined. The weight of the aircraft is randomly set to 917 kilograms. The initial upset recovery policy is the upset recovery policy achieved in the pre-train stage. The accumulated rewards achieved in the pre-train stage and the fine-tuning stage are shown in Fig. 6 and Fig. 7.

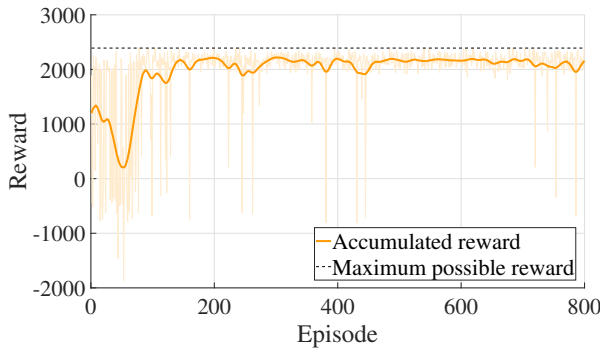


Fig. 6. Reward achieved in the pre-train stage

As shown in Fig. 6, in the pre-train stage, upset recovery policy converges in 100 episodes. The accumulated reward

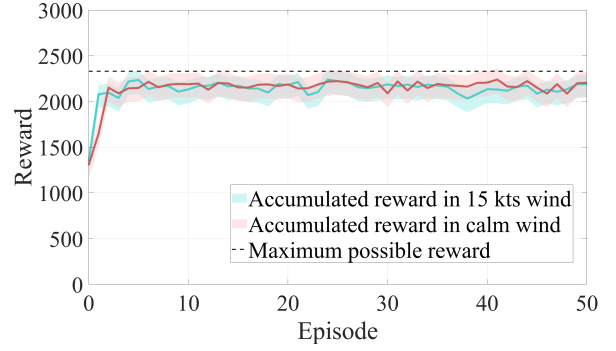


Fig. 7. Rewards achieved in the fine-tuning stage

reaches approximately 2300 on average after convergence. It is shown that the upset recovery policy achieved in the pre-train stage is sub-optimal with respect to the maximum possible reward. Besides, we can see that the achieved reward is anomalously low in some episodes after the convergence, suggesting that the upset recovery policy has unstable performances.

In the fine-tuning stage, the weight of the aircraft and wind condition are determined. We train an upset recovery policy based on the upset recovery policy achieved in the pre-train stage, both in 15 knots wind and in calm wind. As shown in Fig. 7, upset recovery policy converges in three episodes both in wind and in calm wind. Also, we can find that upset recovery policies perform stably and can approach the maximum possible reward both in wind and in calm wind with fine-tuning.

### D. Evaluation of Upset Recovery Policies

The performance of upset recovery policies achieved according to the proposed two-stage strategy is evaluated based on stall recovery and spin recovery. Proportional-Integral-Differential (PID) controllers that have been applied on the upset recovery of aircraft [20] are well-tuned and included in the evaluation of RL-based upset recovery policies. PID-based upset recovery policies are regarded as references.

**Stall recovery.** TD3-based upset recovery policy and PID-based upset recovery policy are used to address the stall recovery of an aircraft in wind and in calm wind. The roll, pitch, and airspeed of the aircraft are shown in Figs. 8, 9, and 10. Both the TD3-based upset recovery policy and the PID-based upset recovery policy can recover the aircraft from stall to level flight both in wind and in calm wind. However, according to Fig. 11, the TD3-based upset recovery policy can recover the aircraft with a smaller LOA than the PID-based upset recovery policy can both in wind and in calm wind.

**Spin recovery.** In this study, TD3-based upset recovery policy and PID-based upset recovery policy are used to deal with the spin recovery of an aircraft in wind and in calm wind. As shown in Figs. 12, 13, and 14, both the TD3-based upset recovery policy and the PID-based upset recovery policy can recover the aircraft from spin to level flight both in wind and in calm wind. The roll, pitch, and airspeed of the aircraft can



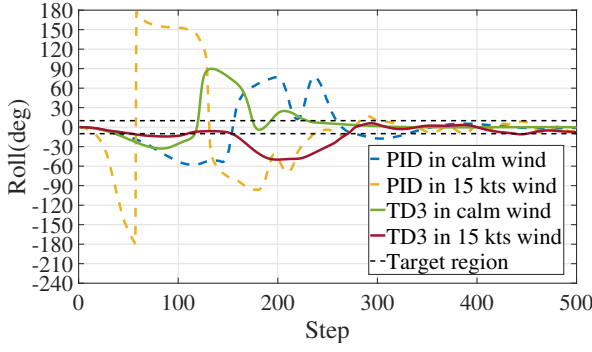


Fig. 8. The roll of the aircraft in stall recovery

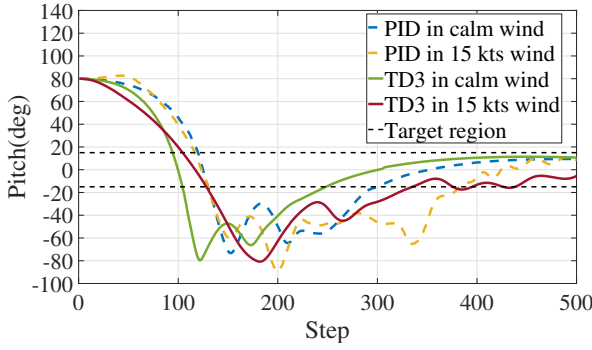


Fig. 9. The pitch of the aircraft in stall recovery

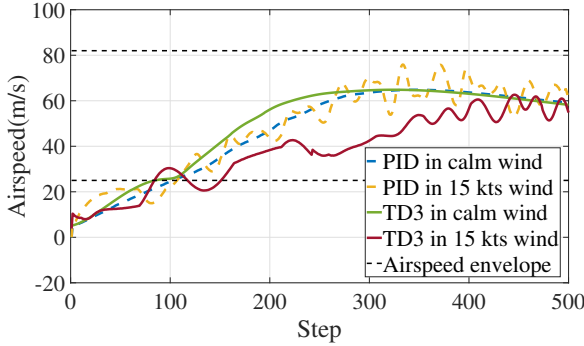


Fig. 10. The airspeed of the aircraft in stall recovery

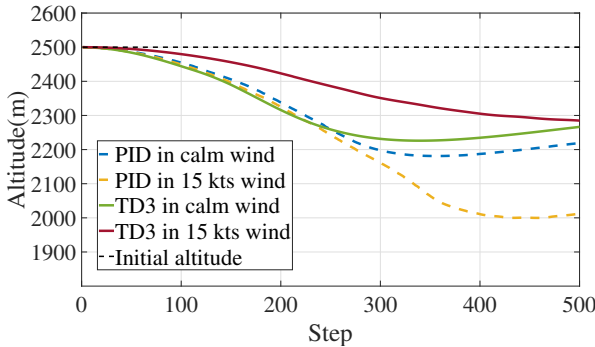


Fig. 11. The altitude of the aircraft in stall recovery

be confined into certain target regions. According to Fig. 15, the TD3-based upset recovery policy can recover the aircraft with a smaller LOA than the PID-based upset recovery policy can in 15 kts wind. the TD3-based upset recovery policy can recover the aircraft with almost the same LOA as the PID-based upset recovery policy can in calm wind.

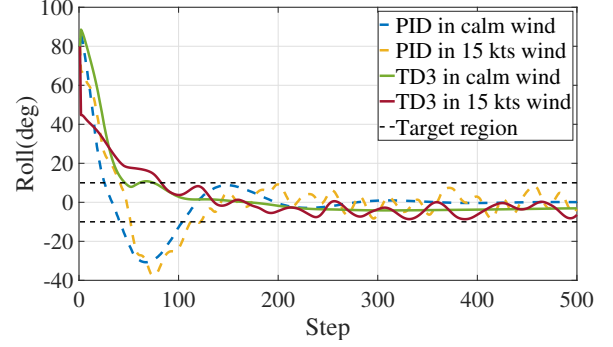


Fig. 12. The roll of the aircraft in spin recovery

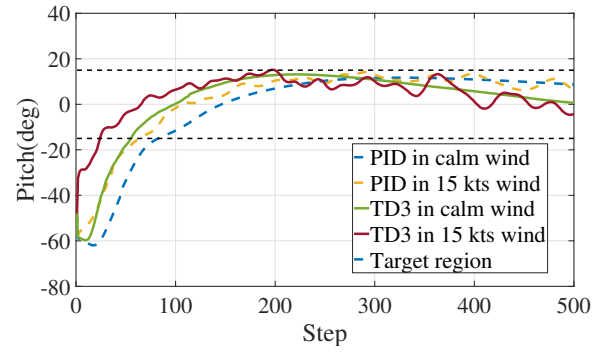


Fig. 13. The pitch of the aircraft in spin recovery

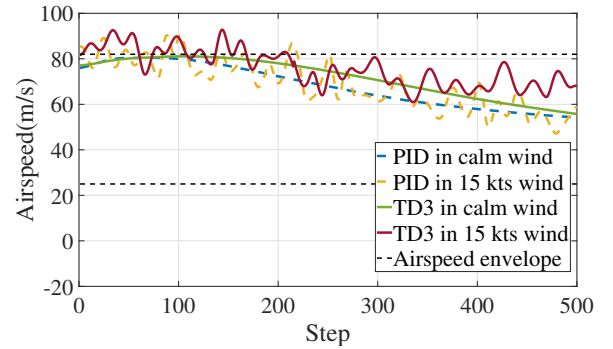


Fig. 14. The airspeed of the aircraft in spin recovery

According to the results of the experiments of stall recovery and spin recovery, we can see that TD3-based upset recovery policy can recover the states of an aircraft into target regions. The effectiveness of the RL-based upset recovery policy is validated. Moreover, TD3-based upset recovery policy achieves

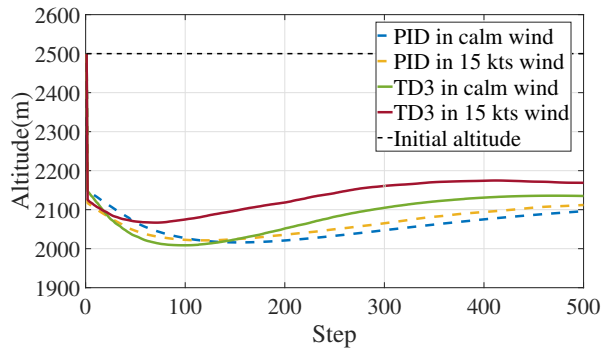


Fig. 15. The altitude of the aircraft in spin recovery

better performances in LOA and time of recovery than well-tuned PID-based upset recovery policy in general.

## V. CONCLUSION

In this paper, a two-stage strategy to achieve an RL-based upset recovery policy was proposed to recover an aircraft from an arbitrary upset situation to level flight taking the time of recovery and LOA into account. According to the proposed two-stage strategy, an algorithm to achieve a TD3-based upset recovery policy was developed. To demonstrate the implementation and verify the effectiveness of the proposed two-stage strategy, training of upset recovery policies and experiments of stall recovery and spin recovery were conducted based on a Cessna172SP in X-Plane 11. The results of training show that an RL-based upset recovery policy for a certain aircraft can be obtained in a few episodes based on a general RL-based upset recovery policy. Experiments validate that the achieved RL-based upset recovery policy can recover an aircraft from either stall or spin to level flight and has a better or equal performance in LOA than a well-tuned PID-based upset recovery policy.

## REFERENCES

- [1] Christine M Belcastro, John V Foster, Gautam H Shah, Irene M Gregory, David E Cox, Dennis A Crider, Loren Groff, Richard L Newman, and David H Klyde. Aircraft Loss of Control Problem Analysis and Research Toward a Holistic Solution. *Journal of Guidance, Control, and Dynamics*, 40(4):733–775, 4 2017.
- [2] Anil Yildiz, M Ugur Akcal, Batuhan Hostas, and N Kemal Ure. Switching Control Architecture with Parametric Optimization for Aircraft Upset Recovery. *Journal of Guidance, Control, and Dynamics*, 42(9):2055–2068, 4 2019.
- [3] T Cunis, D Liao-McPherson, J Condomines, L Burlion, and I Kolmanovsky. Economic Model-Predictive Control Strategies for Aircraft Deep-stall Recovery with Stability Guarantees. In *2019 IEEE 58th Conference on Decision and Control (CDC)*, pages 157–162, 2019.
- [4] D.M.K.K. Venkateswara Rao and Tiauw H Go. Optimization of aircraft spin recovery maneuvers. *Aerospace Science and Technology*, 90:222–232, 2019.
- [5] Donghae Kim, Gyeongtaek Oh, Yongjun Seo, and Youdan Kim. Reinforcement Learning-Based Optimal Flat Spin Recovery for Unmanned Aerial Vehicle. *Journal of Guidance, Control, and Dynamics*, 40(4):1076–1084, 6 2016.
- [6] Zhuangdi Zhu, Kaixiang Lin, and Jiayu Zhou. Transfer learning in Deep Reinforcement Learning: A survey. *arXiv*, pages 1–22, 2020.
- [7] Hao Xiong and Xiumin Diao. Safety Robustness of Reinforcement Learning Policies: A View from Robust Control. *Neurocomputing*, 422:12–21, 2021.
- [8] Brian Dutoi, Nathan Richards, Neha Gandhi, David Ward, and John Leonard. Hybrid Robust Control and Reinforcement Learning for Optimal Upset Recovery. In *AIAA Guidance, Navigation and Control Conference and Exhibit*, Guidance, Navigation, and Control and Co-located Conferences. American Institute of Aeronautics and Astronautics, 8 2008.
- [9] Y Zhu, H Liu, B Ren, H Duan, X She, and Z Wu. A Model-free Flat Spin Recovery Scheme for Miniature Fixed-wing Unmanned Aerial Vehicle. In *2019 IEEE International Conference on Unmanned Systems (ICUS)*, pages 623–630, 2019.
- [10] Ashraf Omran and Ayman Kassem. Optimal task space control design of a Stewart manipulator for aircraft stall recovery. *Aerospace Science and Technology*, 15(5):353–365, 2011.
- [11] Bilal Malik, Jehanzeb Masud, and Suhail Akhtar. A review and historical development of analytical techniques to predict aircraft spin and recovery characteristics. *Aircraft Engineering and Aerospace Technology*, 92(8):1195–1206, 1 2020.
- [12] D W Sparks and D D Moerder. Optimal aircraft control upset recovery with and without component failures. In *Proceedings of the American Control Conference*, volume 5, pages 3644–3649, 2002.
- [13] Roberto Bunge and Ilan Kroo. Automatic Spin Recovery with Minimal Altitude Loss. In *2018 AIAA Guidance, Navigation, and Control Conference*, AIAA SciTech Forum. American Institute of Aeronautics and Astronautics, 1 2018.
- [14] Scott Fujimoto, Herke van Hoof, and David Meger. Addressing Function Approximation Error in Actor-Critic Methods. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 1587–1596. PMLR, 2018.
- [15] C Liu, X Xu, and D Hu. Multiobjective Reinforcement Learning: A Comprehensive Overview. *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, 45(3):385–398, 2015.
- [16] Andrew Y Ng, Daishi Harada, and Stuart J Russell. Policy Invariance Under Reward Transformations: Theory and Application to Reward Shaping. In *Proceedings of the Sixteenth International Conference on Machine Learning, ICML '99*, page 278–287, San Francisco, CA, USA, 1999. Morgan Kaufmann Publishers Inc.
- [17] Kyunghwan Cho, Jinok Shin, and Taeyong Kuc. Design of quadrotor controller for emergency situation using Xplane. In *2015 12th International Conference on Ubiquitous Robots and Ambient Intelligence (URAI)*, pages 311–314, 2015.
- [18] G He, L Yu, S Jia, and X Wang. Simulation verification of Flight Control of a tilt tri-rotor UAV Using X-plane. In *2020 39th Chinese Control Conference (CCC)*, pages 7008–7013, 2020.
- [19] A Bittar, H V Figuereido, P A Guimaraes, and A C Mendes. Guidance Software-In-the-Loop simulation using X-Plane and Simulink for UAVs. In *2014 International Conference on Unmanned Aircraft Systems (ICUAS)*, pages 993–1002, 2014.
- [20] A Khalid, K Zeb, and A Haider. Conventional PID, Adaptive PID, and Sliding Mode Controllers Design for Aircraft Pitch Control. In *2019 International Conference on Engineering and Emerging Technologies (ICEET)*, pages 1–6, 2019.