# Smooth Trajectory Collision Avoidance through Deep Reinforcement Learning

Sirui Song, Kirk Saunders, Ye Yue, Jundong Liu*

School of Electrical Engineering and Computer Science,
Ohio University, Athens, OH 45701

*Abstract*—Collision avoidance is a crucial task in vision-guided autonomous navigation. Solutions based on deep reinforcement learning (DRL) has become increasingly popular. In this work, we proposed several novel agent state and reward function designs to tackle two critical issues in DRL-based navigation solutions: 1) smoothness of the trained flight trajectories; and 2) model generalization to handle unseen environments.

Formulated under a DRL framework, our model relies on margin reward and smoothness constraints to ensure UAVs fly smoothly while greatly reducing the chance of collision. The proposed smoothness reward minimizes a combination of first-order and second-order derivatives of flight trajectories, which can also drive the points to be evenly distributed, leading to stable flight speed. To enhance the agent's capability of handling new unseen environments, two practical setups are proposed to improve the invariance of both the state and reward function when deploying in different scenes. Experiments demonstrate the effectiveness of our overall design and individual components.

*Index Terms*—Deep reinforcement learning, collision avoidance, UAV, smoothness, rewards.

## I. INTRODUCTION

Autonomous navigation capability is of great importance for unmanned aerial vehicles (UAVs) to fly in complex environments where communication might be limited. Collision avoidance (CA) is among the most crucial components of high-performance autonomy and thus has been extensively studied. Generally speaking, the existing CA solutions can be grouped into two categories: geometry-based and learning-based solutions. Geometry-based solutions are commonly formulated as a two-step procedure: first to detect obstacles and estimate the geometry surrounding a UAV, followed by a path planning step to identify a traversable route for escape maneuver.

Learning-based CA solutions extract patterns from training data to perceive environments and make maneuver decisions. Such solutions can be broadly divided into two categories: supervised learning-based and reinforcement learning-based. The former performs perception and decision-making simultaneously, predicting control policies directly from raw input images [1]–[5]. Supervised-based methods are straightforward, but they normally require a large amount of labeled training samples, which are often difficult or expensive to obtain. Reinforcement learning [6], on the other hand, relies on a scale

reward function to motivate the learning agent and explores policy through trial and error. Combined with neural networks, deep reinforcement learning (DRL) has been shown to achieve superhuman performance on a number of games by fully exploring raw images [7]–[9]. DRL-based collision avoidance has also been recently proposed [10] [11] [12] [13]. In order to reduce cost and increase effectiveness, such training is often first carried out within a certain simulation environment.

While remarkable progress has been made in DRL-based navigation solutions, insufficient attention has been given to two critical issues: 1) smoothness of the navigation trajectories; and 2) model generalization to handle unseen environments. For the former, Kahn *et al.* [14] proposed a RL-based solution that seeks a tradeoff of collision uncertainty and speed of UAV motion. When collision uncertainty is high, the motion of the robot/UAV is set to slower, and vice versa. The smoothness of the flight trajectories, however, is not directly addressed. Hasanzade *et al.* [15] proposed an RL-based UAV navigation solution based on a trajectory re-planning algorithm, where high order B-splines are used to define and specify flight trajectories. Due to the local support property of B-spline, such trajectories can be updated quickly, allowing the small UAVs to navigate in clutter environments aggressively. However, new knots need to be inserted over the training process for the re-planning procedure to be fully realized, negatively impacting the overall trajectory smoothness.

Model generalization is a critical issue in machine learning, especially for DRL solutions. Many current DRL works, however, were evaluated on the same environments as they were trained on, such as Atari [16], MuJoCo [17] and OpenAI Gym [18]. For UAV training, there is an additional sim-to-real layer, which complicates the problem even more. Kong *et al.* [19] explored the generalization of various DRL algorithm by training them with different (but not unseen) environments. Doukui *et al.* [20] tackle this issue by mapping exteroceptive sensors, robot state, and goal information to continuous velocity control inputs, but their exploration was only tested on unseen targets instead of unseen scenes.

In this work, we address the afore-mentioned issues with novel designs for agent state and reward functions. To ensure the smoothness of the learned flight trajectories, we integrate two curve smoothness terms, based on first-order and second-order derivatives respectively, into the agent reward

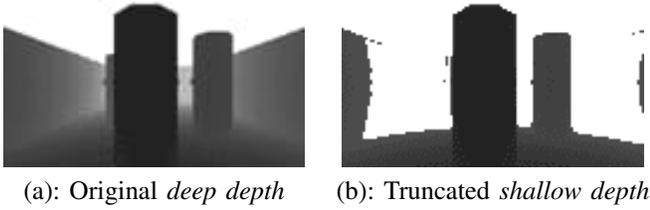(a): Original *deep depth*    (b): Truncated *shallow depth*

Fig. 1: An example pair of deep and shallow depth maps.

functions. To improve the agent's generalization capability, two practical setups, *shallow depth* and *unit vector towards the target*, are adopted to boost the robustness of the state and reward function in dealing with new environments. The proposed designs are trained and tested in simulation scenes with large geometric obstacles. Experiments demonstrate the effectiveness of our overall design and individual components.

## II. METHOD

In this work, a multirotor UAV takes off at a designated starting point and navigates autonomously towards a destination. The line segment connecting the start and end points is regarded as a predefined path, along which certain objects have been put as obstacles.

### A. Design and environment setup

Our overall design goal is to fly the UAV mostly along the predefined route while being able to avoid the obstacles. Such capability is trained through DRL with the following considerations. Firstly, to ensure the UAV to fly along the predetermined path, we minimize the distance of the drone's trajectory away from such a path. Secondly, to ensure the drone avoids collisions while flying smoothly, we set up a variety of rewards, including those for *margin*, *arrival*, and *penalty for collisions* and *smoothness*. In addition, we aim to design DRL agent with a good generalization capability in handling unseen environments.

*States*, *actions*, and *rewards* are three basic components of most DRL algorithms. In this work, the state $s_t$ at time $t$ is defined to include three components: 1) the depth map of the current view facing the camera; 2) the current velocity of the UAV, and 3) a unit vector pointing from the UAV's current position to the target.

Our choices of 1) and 3) are both with model generalization in mind. At each time point, a depth image is obtained from the onboard monocular camera of the UAV. In order to limit the impact of environment changes and thus improve the generalization capability of our model, we focus on nearby objects and ignore those beyond a certain distance. We call this truncated depth image as *shallow depth*, in contrast with the original *deep depth*. An example pair are shown in Fig 1. We include *unit vector to the target* as part of the agent's state, which later will also be used in our proposed reward function. This is in contrast to the Euclidean distance of the UAV away from the destination. Compared with the distances, our unit vectors are scale invariant, and therefore have better

generalization potentials to deal with new environments of different sizes. Each action $a_t$ at time $t$ is defined as $(v_{x_t}, v_{y_t})$, a velocity vector with x-axis and y-axis components. The proposed reward functions will be explained in the next subsection.

We choose *Deep Deterministic Policy Gradient* (DDPG) [21] as the DRL algorithm to train the flight policy of the UAV. DDPG uses an actor-critic method in which the critic network learns the value function (Q value), and the actor network decides how the policy model should be updated. The output of the Actor network can be real-valued vectors, which enables the DDPG model to directly learn actions in continuous space. The detailed network structure and state composition of our DDPG model can be found in Fig. 2.

Furthermore, for each state, we keep historical information and stack together depth images from several consecutive time points. This is designed to alleviate blind spot issues in flight and allow us to monitor the flight trajectory for smoothness control. To reduce the dimensionality of the depth map, we use an LSTM network on the depth stack, as shown in Fig. 2, to capture basic information before feeding it into the actor and critic networks.

### B. Reward functions

In this work, the overall reward $r_t$ at time $t$ is designed to include multiple components, each of which corresponds to a desired system condition. Note that our design goals include: 1) avoiding collisions, 2) enhancing model generalization, and 3) encouraging smooth flights. The overall $r_t$ is given as follows,

$$r_t = R_{\text{margin}} + R_{\text{towards}} + R_{\text{smooth}} + \begin{cases} R_{\text{g}} & \text{if at destination} \\ R_{\text{c}} & \text{if collision} \\ R_{\text{f}} & \text{if normal flight} \end{cases}$$

where $R_{\text{margin}}$ and $R_{\text{smooth}}$ denote the rewards to ensure margin and smoothness respectively; $R_{\text{towards}}$ is aimed to attract the UAV to fly towards the target; $R_g$ is rewarded if the UAV reaches the end point; $R_c$ is a a penalty (negative reward) if collisions happen; and $R_f$ contains a reward for flying forward and a penalty for any deviation from the predefined route.

$R_{\text{margin}}$ is design to penalize the UAV for getting too close to the obstacles. Two margin zones, *soft margin* and *hard margin* are set up, as shown in Fig. 3. When the drone flies into the soft margin zone, it will be pushed back with a moderate force. If it enters the hard margin zone, the system should provide a rapidly increasing repulsive force to prevent the drone from getting closer to the obstacle. Computationally, this two-margin design is implemented as:

$$R_{\text{margin}} = \begin{cases} -C_1(d_{\text{soft}} - d_{\text{obs}})/(d_{\text{soft}} - d_{\text{hard}}) & \text{in soft-margin} \\ -C_2/d_{\text{obs}} & \text{in hard-margin} \\ 0 & \text{otherwise} \end{cases}$$

where $C_1$ and $C_2$ are positive constants; $d_{obs}$ represents the minimum distance from the drone to the nearest obstacle; $d_{\text{soft}}$,
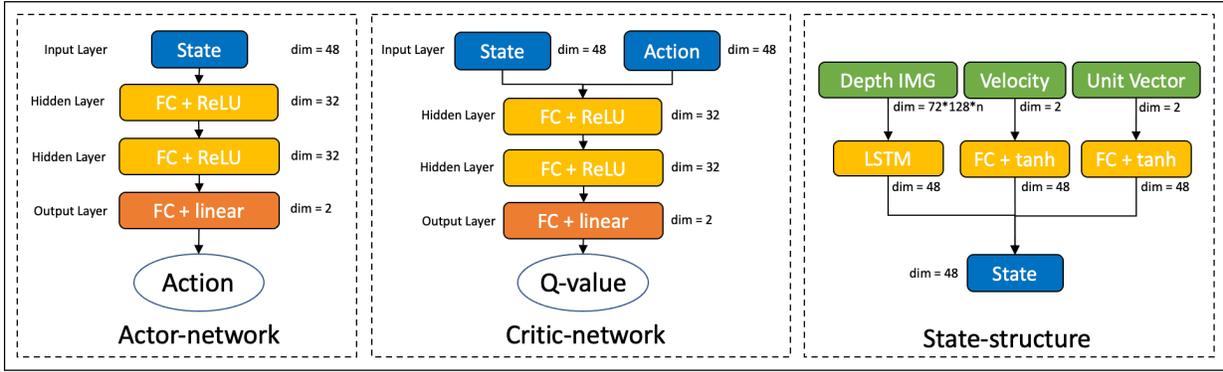
Fig. 2: Overall network architecture of our DDPG-based method. Left: actor-network; Middle: critic-network; Right: information fusion and dimensional reduction through LSTM.

$d_{\text{hard}}$ are constant parameters denoting the range of soft zone and hard zone respectively.
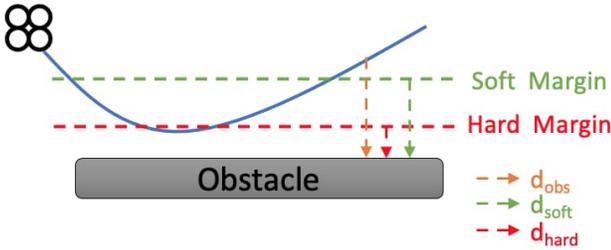


Fig. 3: Illustration of the soft margin and hard margin zones.

$R_{\text{towards}}$ can be designed in many different ways. For instance, the Euclidean distance away from the destination can be used as a negative reward to attract the UAV. However, absolute distances would have poor generalization as they are not scale invariance. In this work, we use *unit vector to the target* as a better alternative, which is also used as part of the agent state. Let $v_d$ be the direction (a unit vector) from current position to the destination and $v_{vel}$ be the direction of the current speed. We design $R_{\text{towards}} = cos(v_d, v_{vel})$. Comparing with absolute distances, our design is scale invariant, which leads to two major benefits: 1) potentially better generalization for unseen environments; and 2) enhanced training and convergence performance. The latter comes from the fact that angle to the destination always provides strong guidance, even when the UAV is far away from the destination.

*1) Smoothness rewards:* The design of our smoothness rewards is inspired by the classic active contour model (also known as snakes) proposed by Kass *et al.* for image segmentation [22]. Let $v(\mathbf{s})$ be an evolving contour. Snake model drives the contour to capture a target object, with an internal energy term to ensure the smoothness of the evolution curve:

$$J(v) = J_1(v) + J_2(v) = \int \alpha |v(\mathbf{s})'|^2 + \beta |v(\mathbf{s})''|^2 d\mathbf{s} \quad (1)$$

where $\alpha$ and $\beta$ are contributing weights. The Euler-Lagrange (E-L) equations to minimize $J_1(v)$ and $J_2(v)$ are linear

functions and cubic functions, respectively. The former (linear functions) would force the curve to stretch up into straight line segments, while the latter (cubic functions) ensures the curve to keep smooth at all points.

Inspired by the snake model, we take flight trajectories as evolving curves and design a smoothness reward $R_{\text{smooth}}$ based on the combination of first-order and second-order derivatives of the trajectories. More specifically, we intend to minimize the square of the first-order derivative to stretch up the trajectory, so the UAV would be forced to fly directly towards the destination. Minimizing the squared second-order derivative would make the flight path generally smooth.

Let $\mathbf{p}_t$ denote the location of the UAV at time $t$. The first-order term will be minimized when $\mathbf{p}_{t-1}$, $\mathbf{p}_t$, $\mathbf{p}_{t+1}$ are on the same line. The difference between $|\mathbf{p}_t - \mathbf{p}_{t-1}| + |\mathbf{p}_{t+1} - \mathbf{p}_t|$ and $|\mathbf{p}_{t+1} - \mathbf{p}_{t-1}|$ can be used as an indictor to measure the deviation of the points from collinearity. The norm of the second-order derivative can be numerically approximated with $|\mathbf{p}_{t+1} + \mathbf{p}_{t+1} - 2\mathbf{p}_t|$, which would be minimized by uniform and gradual changes in direction. With these two observations, we come up with our smoothness reward at time $t$:

$$\begin{aligned} R_{\text{smooth}} = &- C_3(|\mathbf{p}_t - \mathbf{p}_{t-1}| + |\mathbf{p}_{t+1} - \mathbf{p}_t| - |\mathbf{p}_{t+1} - \mathbf{p}_{t-1}|) \\ &- C_4|\mathbf{p}_{t-1} - 2\mathbf{p}_t + \mathbf{p}_{t+1}| \end{aligned}$$
(2)

where $C_3$ and $C_4$ are positive constants, which can be set manually or empirically in experiments. It should be noted that, for computational convenience and numerical stability, we use the norms of the first-order and second-order derivatives instead of their squared values as in the original snake model. In addition, our setup of the first-order derivative term, inspired by the Greedy snake model [23], would encourage the points to be evenly distributed along the trajectory, leading to stable flight speed.

## III. SIMULATION ENVIRONMENTS AND EXPERIMENT DESIGN

In this work, we train our DRL solutions in one training scene and test them under three different test scenes, which

(a): Training scene       (b): Test scene 1
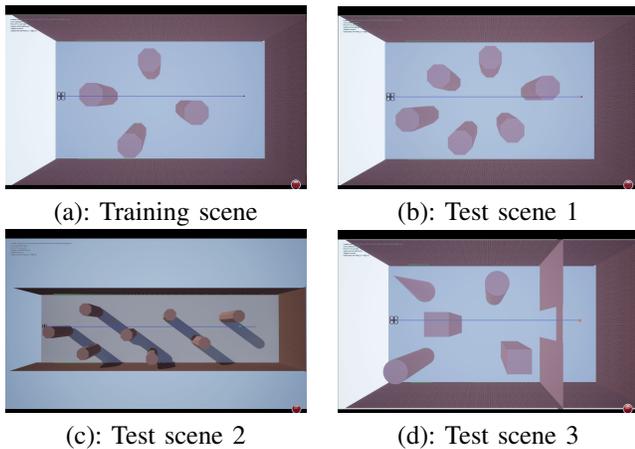
(c): Test scene 2       (d): Test scene 3

Fig. 4: Simulation scenes of training and testing environments under Airsim + UE4. (a) is the training scene; (b) has more obstacles; (c) is a test scene of larger size and (d) has unseen shapes. In each scene, the solid line is the predefined route and the orange point at the right side is the destination. Refer to text for details.

are previously unseen. The scenes are designed to evaluate our proposed components in boosting the generalization capability and trajectory smoothness, respectively.

The training scene is designed as a 3D container, 20 meters long, 15 meters high and 15 meters wide. Four obstacles are put inside, as shown in Fig. 4.(a). The drone takes off at the start point, and then flies towards the destination, which is marked with an orange dot. The straight line connecting the two points is the predefined route. Four thick hexagonal prism (diameter = 2.55m) obstacles are on the predefined route, blocking the drone flying to the destination.

Fig. 4(b - d) show the three test scenes. The first test scene (TS1) has the same dimension as the training scene, but has two additional (totally 6) hexagonal prism obstacles of the same shape and size. This setup is aimed to test whether the UAV can successfully fly to the end point in an environment with more crowded obstacles. The second test scene (TS2) has eight hexagonal prism obstacles of the same size, and the total scene length is set to 45 meters. Compared with the training scene, this scene is much longer, aiming to evaluate model generalization under different scene sizes. The third test scene (TS3) has the same dimension as the training scene, however, is filled with obstacles of different shapes. This is to test the learned policy in handling obstacle shapes that have not been trained with.

## IV. EXPERIMENTS AND RESULTS

We build the training and testing simulation scenes on Airsim and UE4, simulating the UAV collision avoidance in a real environment. The simulator runs on a Nvidia GeForce RTX 2080 graphics card and the UAV in the simulator is equipped with a LiDAR, a GPS and an onboard monocular camera. Experiments are conducted to evaluate the effectiveness of

our designs for 1) boosting the model generalization and 2) ensuring the smoothness of flight trajectories.

### A. Results for smoothness setups

Our designs for improving the trajectory smoothness comes from the reward component $R_{smooth}$ in Eqn. (2). For the convenience of discussion, we use $W_{smooth}$ to denote the full DRL model with $R_{smooth}$ and $WO_{smooth}$ to denote the model without this term. Evaluations are conducted based on visual inspections and statistical analysis, under both the training and testing environments.

Fig. 5 shows an example result of the models under the training environment. Fig. 5(a) and Fig. 5(b) record the trajectories of the UAV, where the red dashed line is the predefined route. The trajectory produced by $W_{smooth}$ is smoother and closer to the predefined route than that of $WO_{smooth}$.



(a) With $R_{smooth}$       (b): Without $R_{smooth}$
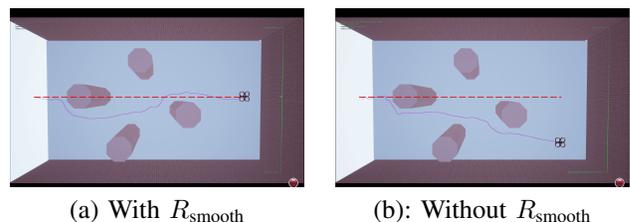
Fig. 5: An example pair of training results from the models of (a) $W_{smooth}$ and (b) $WO_{smooth}$.

For the inference performance in the unseen test environments, Fig. 6 shows an example of the trajectories generated by the two models. In each figure, the dashed line is the predefined route and the purple line shows the trajectory of the UAV. One can see $W_{smooth}$ clearly outperforms $WO_{smooth}$ in all three test scenes, producing smoother and more direct tracks towards the destination. In the third test scene, $WO_{smooth}$ can not find the path to the destination, as the UAV collides with the cube and the door, indicating that it fails to bypass obstacles of unseen shapes. In contrast, $W_{smooth}$, with the proposed smoothness rewards, can successfully find a rather smooth path to reach the destination.

We conduct multiple inference tests of the two models under the three test scenes. The numbers of the tests are 100, 50 and 100, respectively, for test scenes 1, 2 and 3. Three evaluation metrics are employed to evaluate the policy performance: *average linear acceleration* over the flights, *average path curvature* of the flight trajectories, and *success rate* to reach the destinations. The comparisons are summarized in Table I. Our $W_{smooth}$ model has lower average curvatures and accelerations than $WO_{smooth}$, demonstrating that the produced trajectories are generally smoother and more stable. Our model also has higher success rates to reach destinations, indicating the smooth terms also make positive contributions to training convergence.

It should be noted that our setup of the first-order derivative in the smoothness term in Eqn. (2) not only encourages the flight trajectory to be tighten up, but also drives the points to be distributed evenly [23]. This would lead to more stable

flight speed, which is partly reflected in the smaller values of the average acceleration. Both first-order and second-order derivatives contribute to the overall smoothness. The smaller average curvatures produced by $W_{\text{smooth}}$ indicate we have achieved our design goals.



(a): Test 1: with $R_{\text{smooth}}$     (b): Test 1: without $R_{\text{smooth}}$

(c): Test 2: with $R_{\text{smooth}}$

(d): Test 2: without $R_{\text{smooth}}$

(e): Test 3: with $R_{\text{smooth}}$     (f): Test 3: without $R_{\text{smooth}}$
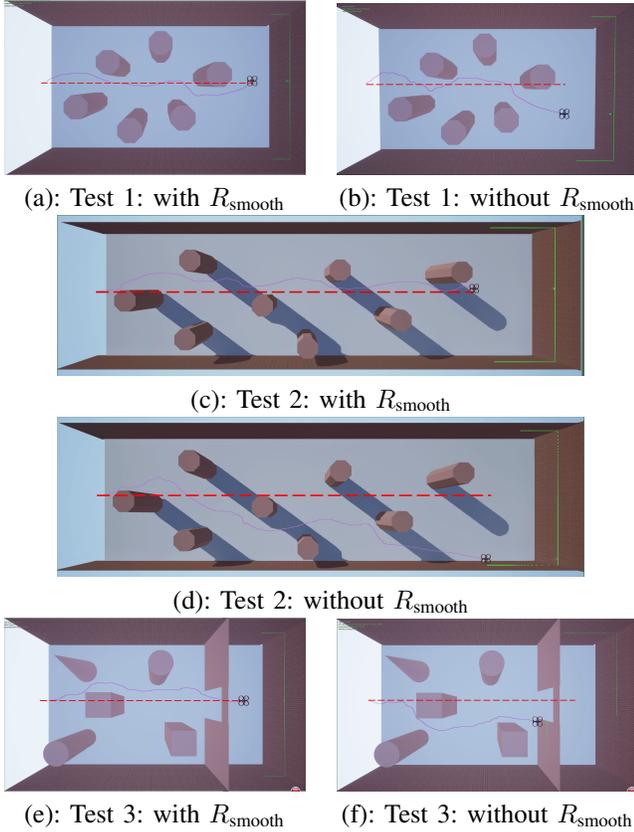
Fig. 6: (a) and (b): final trajectories of $W_{\text{smooth}}$ and $WO_{\text{smooth}}$ in Test scene 1. (c) and (d): final trajectories of the two models in Test scene 2. (e) and (f): final trajectories of the two models in Test scene 3.

### B. Results for model generalization

In this work, we propose two practical designs for improving the DRL agent's handling of environment changes: 1) using *shallow depth* instead of *deep depth* as a part of agent's state; and 2) using *unit vector towards destination* (denoted as $\overrightarrow{vec}$) instead of the distance (denoted as $dist$) in agent's state and reward functions.

In order to test the effectiveness of our designs, we design four DRL models based on: 1) *deep depth* + $\overrightarrow{vec}$; 2) *shallow depth* + $\overrightarrow{vec}$; 3) *deep depth* + $dist$ and 4) *shallow depth* + $dist$. It should be noted that $\overrightarrow{vec}$ is used in both the state and reward of our models. In models 3) and 4), both $\overrightarrow{vec}$ terms are replaced with $dist$. To simplify the analysis, smoothness reward is not integrated into the models. For each of the combination, we train the policy under the training environment to convergence first, and then directly perform inferences in the three unseen test environments described in the previous section. Note that

test scene TS1 has more crowded obstacles than the training scene; test scene TS2 has a larger size; and test scene TS3 is filled with objects that have not been seen in the training.

The four models are all trained successfully in the training scene. When moving into the three testing environments, we conducted 100 trials for each of the TS1, TS2 and TS3 environments.

To measure the performance of the trained policies in the unseen testing environments, we design two evaluation metrics: 1) *success rate* (SR), which is the percentage of the trials where the UAV can successfully arrive at the destination; and 2) *capability of collision avoidance* (CAC), which is defined as the average length that the UAV can fly prior to its first collision. The results are shown in Table II.

Comparing with the other three models, *shallow* + $\overrightarrow{vec}$ demonstrates a much stronger capability of finding the targets in unseen environments. It also shows the best capability in avoiding collisions. *deep depth* + $dist$, on the other hand, show the worst performance in CAC.

From the above experiments, we can observe that both *shallow depth* and *unit vector* are helpful for the learning policy to be directly deployed in changed environments. Their combination demonstrates the best generalization capability in the experiments.

## V. CONCLUSIONS

In this work, we propose a DRL-based collision avoidance solution for UAVs to fly smoothly and stably. We also seek to develop our system to have good generalization performance to handle unseen environments. A number of novel designs are made regarding the agent's state, as well as its reward function. The major innovation is a smoothness reward term, based on the minimization of a combined first-order and second-order derivatives. Experiments demonstrate this term can indeed lead to smoother trajectories, as well as flights with stable speeds.

Our design for boosting model generalization is based on shallow depth and unit vector towards target. While simple, this combined setup reduces the input space and increases system invariance, leading to greatly enhanced robustness. Experimental results demonstrate the effectiveness of our overall design and individual components. To transplant proposed solutions and deploy them into real UAVs is our ongoing effort.

### REFERENCES

[1] D. K. Kim and T. Chen, "Deep neural network for real-time autonomous indoor navigation," *arXiv preprint arXiv:1511.04668*, 2015.

[2] A. Giusti, J. Guzzi, D. C. Cireşan, F.-L. He, J. P. Rodríguez, F. Fontana, M. Faessler, C. Forster, J. Schmidhuber, G. Di Caro *et al.*, "A machine learning approach to visual perception of forest trails for mobile robots," *IEEE Robotics and Automation Letters*, vol. 1, no. 2, pp. 661–667, 2015.

[3] L. Tai, S. Li, and M. Liu, "A deep-network solution towards model-less obstacle avoidance," in *2016 IEEE/RSJ international conference on intelligent robots and systems (IROS)*. IEEE, 2016, pp. 2759–2764.

[4] X. Dai, Y. Mao, T. Huang, N. Qin, D. Huang, and Y. Li, "Automatic obstacle avoidance of quadrotor uav via cnn-based learning," *Neurocomputing*, vol. 402, pp. 346–358, 2020.

[5] S. Back, G. Cho, J. Oh, X.-T. Tran, and H. Oh, "Autonomous uav trail navigation with obstacle avoidance using deep neural networks," *Journal of Intelligent & Robotic Systems*, vol. 100, no. 3, pp. 1195–1211, 2020.

TABLE I: Quantitative results for $WO_{\text{smooth}}$ and $W_{\text{smooth}}$ under the unseen test scenes. Acc. is short for *Acceleration*; Cur. is short for *Curvature* and Succ. is for *Success Rate*.

| Scene | $WO_{\text{smooth}}$ | | | $W_{\text{smooth}}$ | | |
|---|---|---|---|---|---|---|
| | Acc. $(m/s^2)$ | Cur. $(m^{-1})$ | Succ. (%) | Acc. $(m/s^2)$ | Cur. $(m^{-1})$ | Succ. (%) |
| Test scene 1 | 0.07 | 0.08 | 99 | 0.03 | 0.06 | 90 |
| Test scene 2 | 0.04 | 0.08 | 45 | 0.03 | 0.05 | 57 |
| Test scene 3 | NA | NA | 0 | 0.04 | 0.06 | 100 |

TABLE II: Comparisons of the 4 policies in the unseen testing environments, measured with *success rate* (SR) and *capability of collision avoidance* (CAC).

| | deep + dist | shallow + dist | deep + $\overline{vec}$ | shallow + $\overline{vec}$ |
|---|---|---|---|---|
| SR | 3 % | 18.33% | 3 % | 63.67% |
| CAC | 50.18% | 55.10% | 57.25% | 87.33% |

[6] J. Michels, A. Saxena, and A. Y. Ng, "High speed obstacle avoidance using monocular vision and reinforcement learning," in *Proceedings of the 22nd international conference on Machine learning*, 2005, pp. 593–600.

[7] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski *et al.*, "Human-level control through deep reinforcement learning," *nature*, vol. 518, no. 7540, pp. 529–533, 2015.

[8] L. Xie, S. Wang, A. Markham, and N. Trigoni, "Towards monocular vision based obstacle avoidance through deep reinforcement learning," *arXiv preprint arXiv:1706.09829*, 2017.

[9] L. He, N. Aouf, J. F. Whidborne, and B. Song, "Integrated moment-based lgmd and deep reinforcement learning for uav obstacle avoidance," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2020, pp. 7491–7497.

[10] A. Singla, S. Padakandla, and S. Bhatnagar, "Memory-based deep reinforcement learning for obstacle avoidance in uav with limited environment knowledge," *IEEE transactions on intelligent transportation systems*, vol. 22, no. 1, pp. 107–118, 2019.

[11] Z. Xue and T. Gonsalves, "Vision based drone obstacle avoidance by deep reinforcement learning," *AI*, vol. 2, no. 3, pp. 366–380, 2021.

[12] S. Song, Y. Zhang, X. Qin, K. Saunders, and J. Liu, "Vision-guided collision avoidance through deep reinforcement learning," in *NAECON 2021 - IEEE National Aerospace and Electronics Conference*, 2021, pp. 191–194.

[13] S. Ouahouah, M. Bagaa, J. Prados-Garzon, and T. Taleb, "Deep reinforcement learning based collision avoidance in uav environment," *IEEE Internet of Things Journal*, 2021.

[14] G. Kahn, A. Villaflor, V. Pong, P. Abbeel, and S. Levine, "Uncertainty-aware reinforcement learning for collision avoidance," *arXiv preprint arXiv:1702.01182*, 2017.

[15] M. Hasanzade, O. Shadeed, and E. Koyuncu, "Deep reinforcement learning based aggressive collision avoidance with limited fov for unmanned aerial vehicles," in *AIAA SCITECH 2022 Forum*, 2022, p. 2043.

[16] M. G. Bellemare, Y. Naddaf, J. Veness, and M. Bowling, "The arcade learning environment: An evaluation platform for general agents," *Journal of Artificial Intelligence Research*, vol. 47, pp. 253–279, 2013.

[17] E. Todorov, T. Erez, and Y. Tassa, "Mujoco: A physics engine for model-based control," in *2012 IEEE/RSJ international conference on intelligent robots and systems*. IEEE, 2012, pp. 5026–5033.

[18] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba, "Openai gym," *arXiv preprint arXiv:1606.01540*, 2016.

[19] J. Yu, Y. Su, and Y. Liao, "The path planning of mobile robot by neural networks and hierarchical reinforcement learning," *Frontiers in Neurorobotics*, p. 63, 2020.

[20] O. Doukhi and D.-J. Lee, "Deep reinforcement learning for end-to-end local motion planning of autonomous aerial robots in unknown outdoor environments: Real-time flight experiments," *Sensors*, vol. 21, no. 7, p. 2534, 2021.

[21] D. Silver, G. Lever, N. Heess, T. Degris, D. Wierstra, and M. Riedmiller, "Deterministic policy gradient algorithms," in *International conference on machine learning*. PMLR, 2014, pp. 387–395.

[22] M. Kass, A. Witkin, and D. Terzopoulos, "Snakes: Active contour models," *International journal of computer vision*, vol. 1, no. 4, pp. 321–331, 1988.

[23] D. J. Williams and M. Shah, "A fast algorithm for active contours and curvature estimation," *CVGIP: Image understanding*, vol. 55, no. 1, pp. 14–26, 1992.