

# Deep Reinforcement Learning based Collision Avoidance in UAV Environment

Siheem Ouahouah<sup>||</sup>, Miloud Bagaa<sup>||</sup>, Jonathan Prados-Garzon<sup>\*\*</sup>, and Tarik Taleb<sup>||§</sup>

<sup>||</sup> Department of Communications and Networking, School of Electrical Engineering, Aalto University, Espoo, Finland. Emails: firstname.lastname@aalto.fi

<sup>\*\*</sup> Department of Signal Theory, Telematics, and Communications, Univ. of Granada, Granada, Spain. Email: jpg@ugr.es

<sup>§</sup> Department of Computer and Information Security, Sejong University, Seoul, Korea.

**Abstract**—Unmanned Aerial Vehicles (UAVs) have recently attracted both academia and industry representatives due to their utilization in tremendous emerging applications. Most UAV applications adopt Visual Line of Sight (VLOS) due to ongoing regulations. There is a consensus between industry for extending UAVs' commercial operations to cover the urban and populated area controlled airspace Beyond VLOS (BVLOS). There is ongoing regulation for enabling BVLOS UAV management. Regrettably, this comes with unavoidable challenges related to UAVs' autonomy for detecting and avoiding static and mobile objects. An intelligent component should either be deployed onboard the UAV or at a Multi-Access Edge Computing (MEC) that can read the gathered data from different UAV's sensors, process them, and then make the right decision to detect and avoid the physical collision. The sensing data should be collected using various sensors but not limited to Lidar, depth camera, video, or ultrasonic. This paper proposes probabilistic and Deep Reinforcement Learning (DRL)-based algorithms for avoiding collisions while saving energy consumption. The proposed algorithms can be either run on top of the UAV or at the MEC according to the UAV capacity and the task overhead. We have designed and developed our algorithms to work for any environment without a need for any prior knowledge. The proposed solutions have been evaluated in a harsh environment that consists of many UAVs moving randomly in a small area without any correlation. The obtained results demonstrated the efficiency of these solutions for avoiding the collision while saving energy consumption in familiar and unfamiliar environments.

**Index Terms**—Unmanned Aerial Vehicles (UAVs), Collision Avoidance, Multi-Access Edge Computing (MEC), Machine Learning, and Deep Reinforcement Learning.

## I. INTRODUCTION

Unmanned aerial vehicles (UAVs), commonly recognized as drones, are small, fast, and mobile cyber-physical entities employed in different industrial verticals, including power supply inspection, parcel and package delivery, disaster management, and traffic monitoring [1]. The utilization of UAVs goes beyond industrial and academic purposes to daily personal use. A UAV operator must always be capable of maintaining the Visual Line of Sight (VLOS) of its UAV that is piloting due to ongoing regulations, unaided by any technology other than prescription glasses or contact lenses.

This work has been partially funded by the Spanish national project TRUE-5G (PID2019-108713RB-C53).

While UAVs are used mostly within VLOS, there is enthusiasm towards their utilization beyond visual line of sight (BVLOS) to enable new emerging applications. Therefore, there is a consensus between industry for attenuation of the regulation by extending UAVs' commercial operations to cover the urban and populated area controlled airspace BVLOS. The latter will be enabled by leveraging a cellular wireless network. 5G system and beyond considers the UAV management BVLOS as one of the essential demonstrators. On the other side, emerging networking paradigms, such as Edge Computing can substitute UAVs to handle high processing flight control applications. Furthermore, GPU vendors allow for realizing different micro-architectures (e.g., Fermi, Maxwell, and Pascal) that might enable real-time and high resourced applications for the UAV's flight control [2].

The UAVs' commercial revenue sees considerable growth by the near future[3].The expected increase in the number of UAVs involves new challenges related to their control and management. Efficient solutions for UAV's collision avoidance is one of the challenges that have been widely tackled in the literature in both ground vehicles [4], [5] context, as well as in the context of UAVs[3], [6]–[24]. Different sensors have been leveraged for scanning and detecting objects surroundings UAVs. Some solutions use cameras for detecting mobile and static obstacles around UAVs[6], [14]. Nevertheless, the information provided by video cameras requires intensive processing to be translated into useful information to control UAVs [25].

Several works [3], [11], [17], [18] have proposed path planning solutions where the UAVs are provided with their whole trajectories before starting their missions to overcome the limitation mentioned above. The path planning fits well in applications with invariable environment scenarios. However, mostly, UAVs fly in unsettled indoor, urban, and confined areas. Indeed, sensing and path planning approaches' success is highly related to the computation capacity of the UAV, the accuracy of the sensor, and the knowledge's degree on the environment. On another side, Reinforcement learning (RL) based approaches got much success in emerging topics, including robotic prediction, Vehicular Ad hoc Networks (VANET), and UAVs. For instance, authors in [26] have provided a heuristic to enhance communication and prevent

jamming attacks in VANET by leveraging UAV. An extended version of this work has been suggested in [27] to prevent the jamming attacks in VANET by leveraging both UAV and RL approaches. This success attracted the researchers to use RL to ensure a self-decision making system for a safe UAV's autonomous flight.

RL consists in providing a kind of knowledge on the environment based on the previous UAV's experiences. Thus, RL builds the knowledge by interacting with the environment based on a Markov Decision Process (MDP) model and following one of the RL methods (e.g., Q\_learning, Deep Q-Networks [DQN], Policy\_gradient, or Actor\_critic). The RL-based UAVs control solutions proposed so far [19]–[24] need large datasets that refer to the abstraction level used to model the RL environment system (e.g., velocity, wind velocity, etc.). The datasets used in RL based solutions might return the same limitations pointed in the classical approaches stated above [19]. To overcome these limitations, this paper suggests two strategies for avoiding the collision in a UAV environment. The first solution, named **Probability distribution based collision avoidance framework (PICA)**, leverages the probabilistic model for avoiding collisions. In contrast, the second solution, dubbed **RL-based collision avoidance framework (RELIANCE)**, leverages Deep Q-Networks (DQNs) for avoiding collisions.

To deal with the disparate UAVs processing capacities and to ease the deployment of the proposed solutions, we suggest two deployment approaches: the UAV's flight controller is deployed onboard or at the Multi-Access Edge Computing (MEC). The first approach convenes UAVs with the new GPU microarchitecture technology is where the agent, either of RELIANCE or PICA, can smoothly make-decision and select the best actions. On the other hand, the second way assembles UAVs with limited computing capacities. In order to ensure close management, services running should be migrated among MECs using Follow Me Edge-Cloud concept. Authors in [28] suggest a MEC architecture that ensures UAVs' resource provisioning. In case the UAV has a limited resource capacity, the same architecture as proposed in [28] can be adopted. In this case, the RELIANCE/PICA agent is responsible for making decisions at the MEC, and then sending the respective actions to the UAV for controlling its motion.

Both Algorithms aim to enable autonomous decision-making for a UAV while a safe and short flight is ensured to save energy consumption. To ensure a fast convergence of RELIANCE and PICA, we have used a detail-less and generalized state by focusing on the closest part of the environment to the agent. In RELIANCE states' design, we have leveraged Partially Observable Markov Decision Process (POMDP) to ensure the generalization and fast convergence. We have used a partially observable state that focuses only on the UAV surrounding to avoid the collisions instead of the whole deployment area. The limitation of the observation at the UAV vicinity helps to reduce the state space by aggregating many observations to a single state. Thanks to this strategy, RELIANCE and PICA avoid over-

loading computation processing by ignoring useless knowledge. Moreover, this strategy helps RELIANCE solution to converge quickly by treating many observations as the same state. Furthermore, to ensure the generalization and that both Algorithms can work in unseen environments, we have used relative target positions. The benefits of this strategy are twofold: it facilitates and speeds up the convergence of the neural networks and, most importantly, improves the generalization of the agent, which is agnostic to the scenario scale. We have evaluated and compared both Algorithms in terms of collision avoidance and energy saving in familiar and unfamiliar environments. The obtained results demonstrate the ability of both Algorithms in the generalization by performing well in unknown environments. Also, the simulation results clearly show the superiority of RELIANCE comparing to PICA.

The rest of the paper is organized as follows. Section II reviews the related works. Section III includes our system model and problem statement. In section IV, PICA solution is described. An overview on DQN and RELIANCE solution are detailed in section V. Section VI presents and discusses the simulation results of PICA and RELIANCE evaluations. Finally, the primary conclusions are drawn in section VII.

## II. RELATED WORK

There is a vast literature to address the collision avoidance problem in the context of both unmanned ground vehicles [4], [5] and UAVs [3], [6]–[24].

Most of the solutions rely on exact methods [3], [6]–[18], i.e., analytical modelling and optimization techniques, to tackle the UAVs collision avoidance problem (UCAP). The existing works, based on exact methods, usually considers part of the UCAP aspects to handle its modelling and computational complexity. However, in order to provide a realistic and practical model of the UCAP, many issues have to be taken into account:

- **Obstacles detection:** In order to detect the static and mobile objects, the UAVs need to be equipped with onboard sensors. The number of these sensors and their precision might be affected and limited due to several external factors, e.g., specific scenario and UAVs' autonomy. For instance, GPS might not work for indoor scenarios like Industry. Other sensors like radars [8] might be too heavy, energy-consuming and expensive. Then, the concrete set of onboard sensors in UAVs depends on the application and scenario.
- **Sensors errors:** Onboard sensors to detect objects are not error-free. All of them have precision errors, which might be affected by external conditions. For instance, GPS error is affected by weather conditions and follows a Gaussian distribution [3], [29].
- **Complex control:** There are several variables to control the UAVs movement, e.g., direction, velocity, and acceleration. Furthermore, these variables strongly depends on external factors like wind velocity.
- **Different approaches:** There are two different approaches to solve UCAP, namely, path planning and

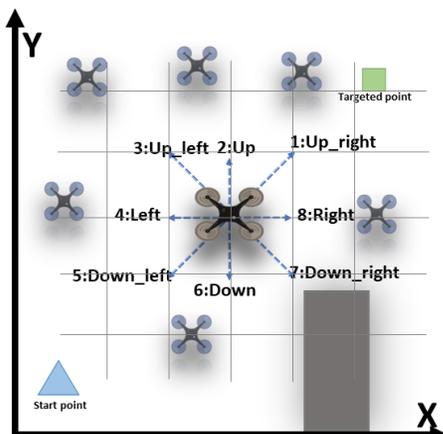


Fig. 1. System model of collision avoidance in UAV environment.

sensing and avoiding [25] methods. Path planning solutions compute the trajectories of the UAVs offline, whereas sensing and avoiding (online) methods determine the movement of the UAVs for small time steps depending on the environment conditions. Sensing and avoiding methods offer higher flexibility and are suitable for a wider range of scenarios. Path planning fits well only for scenarios where the environment remains relatively static during the whole UAVs' mission. The main drawback of the online methods is that typically the UAV has to run the Algorithm (e.g., due to latency constraints), which might exhibit high computational complexity and consume energy.

In the light of the above, the UCAP requires a high domain-knowledge and its modelling leads to complex or even intractable optimization programs. To overcome these problems, Machine Learning techniques are particularly attractive for addressing UCAP, so they have been recently received a lot of attention by the research community [19]–[24].

Choi and Cha in [19] provide a comprehensive survey of ML-assisted solutions for autonomous flight. Specifically, they focus on object recognition and UAV's control strategy. The authors conclude that ML is a promising approach to enable stable flight under uncertain environments, though there are still some open issues that need to be carefully addressed. Among them, existing works do not apply ML in all the UCAP's issues together for autonomous flight. Then, holistic solutions, which cover most of the real world problems and are suitable for a wider spectrum of scenarios, are required. Furthermore, existing solutions need large datasets for training. In this regard, they encourage new less data-hungry proposals with a more lightweight training. Last, they motivate the need for real world tests for an stronger validation of the ML-based UAVs control strategy. Similarly, Fraga-Lamas *et al.* overview the latest advances on IoT UAV systems controlled by deep learning techniques. They analyze the object detection and collision avoidance

TABLE I  
SUMMARY NOTATIONS.

Symbol	Description
UoI	The UAV of Interest.
$X \times Y$	The 2-dimensional geographical area.
$P_S$	The UoI's started position.
$P_T$	The UoI's targeted position.
$\mathcal{A}$	The set of directions controlling UoI motion.
$\mathcal{I}$	The set of mobile and static intruders.
$U(t)$	The position of UoI at time step $(t)$ .
$J(t)$	The position of intruder $j$ at time step $(t)$ .
$\gamma$	The euclidean distance between the positions.
$ A $	The cardinal of the set $A$ .
$\mathcal{Z}$	The set of zones.
$z_i$	The zone surrounding the position $i$ .
$\rho$	The radius of every $z_i$ .
$\eta(i)$	The set of intruders neighboring the position $i$ .
$\mathcal{P}(z_i)$	The probability of collision at the zone $(z_i)$ .
$p_i^j$	The probability of collision between the point $i$ and its neighbor $j$ .
$\theta$	The priority factor rate.
$\Delta$	The squared shaped area surrounds the UoI.
$L_\Delta$	The size of $\Delta$ side.
$\alpha$	The learning rate of DQN.
$\mathcal{L}$	The loss function.
$\epsilon$	A threshold distance before two UAVs collide.
$batch\_size$	The size of each batch.
$Q^\pi$	The policy network.
$Q^T$	The target network.
$\mathcal{M}$	A number of episodes to update $Q^T$ with $Q^\pi$ .
$\omega$	The wight and bias of neural network.
$P_{rel}$	The related address of $P_T$ according to $U(t)$ .
$\xi$	The decay parameter.

problems and present a survey on the state-of-the-art of deep learning techniques to solve them. Also, they detail the most relevant existing datasets and UAVs communication architectures. Finally, they identify the open challenges for UCAP. Interestingly, they extract some similar conclusions to the ones drawn in [19]. For instance, the necessity for large amount of data to generate robust models and the difficulty to produce those data.

In this article, we propose a simple, yet powerful deep reinforcement learning and probability distribution based solutions. Our proposals are suitable for many scenarios, while they need reduced datasets to converge and produce robust models.

### III. SYSTEM MODEL AND PROBLEM STATEMENT

In this paper, we aim to control the movement of a UAV, hereinafter referred to as UoI (UAV of Interest), while avoiding the collisions with static and mobile objects. UoI needs to move within a confined area of dimension  $X \times Y$  while avoiding the collisions. The UoI motion begins from a predefined initial position  $P_S = (x_S, y_S)$  and stops once achieves a predefined targeted destination  $P_T = (x_T, y_T)$ . Let  $\mathcal{A}$  denote the possible action directions of UoI. As mentioned in [30], the possible movement of a UAV is limited and related to the environment that it works in. For the sake of simplicity and without loss of generality, we consider  $\mathcal{A}$  has eight possible directions,  $\mathcal{A} = \{ Up, Down, Right, Left, Up\_right, Up\_left, Down\_right, Down\_left \}$ , as depicted in Fig. 1. For simplicity, we

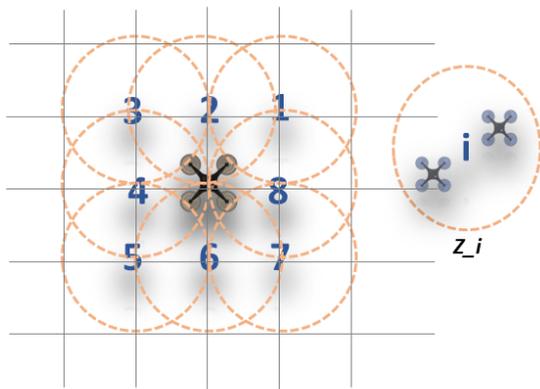


Fig. 2. PICA: zone concept for selecting directions.

assume a constant velocity and altitude for the UoI. Furthermore, we consider that UoI operates autonomously without either any remote ground control or predefined way-points plan. On another side, we consider that the 2-dimensional flying area can include either static (e.g. buildings) and mobile (e.g. birds and other UAVs) obstacles. Therefore, the UoI has to be equipped with an accurate sensor (e.g., Lidar) to precisely detect the surrounding objects' positions. Other ultrasonic-, video-, and radio-based techniques for detecting obstacles and mobile objects have been investigated in the literature [31], [32]. Hereafter we refer to the static and mobile objects as *static\_intruders* and *mobile\_intruders*, respectively. We define  $\mathcal{I}$  as the set of intruders where,  $\mathcal{I} = \{\{static\_intruder\} \cup \{mobile\_intruder\}\}$ .

The UoI might collide with one of the mobile and static obstacles. We assume an arbitrary trajectory for the mobile intruders, which is unknown by the UoI. A collision occurs whenever the euclidean distance between the UoI and any intruders is lower than a predefined threshold distance  $\epsilon$ . The safety distance  $\epsilon$  varies from few centimeters to few meters according to different parameters related to the environment and used sensors. The distance  $\epsilon$  can vary according to the sensor technology used to measure the distances, such as Lidar, depth camera, video, or ultrasonic. Also, the accuracy of the same type of sensors can vary from a manufacturer to another. The UoI can detect this collision by harnessing its onboard sensors at any time  $t$ . Let  $P_u(t) = (x_u(t), y_u(t))$  and  $P_j(t) = (x_j(t), y_j(t))$  denote the position of the UoI and the position of the intruder  $j$  at a given instant  $t$ , respectively. Then, a collision instance is formally defined as follows:

$$Collision = \begin{cases} 1 & \text{if } \exists j \in \mathcal{I} \text{ where } \delta_u^j \leq \epsilon \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

where

$$\delta_u^j = \sqrt{(x_u(t) - x_j(t))^2 + (y_u(t) - y_j(t))^2} \quad (2)$$

In this paper, we also take into consideration the limitation on the battery capacity of the UoI. To that end, the ultimate goal of the algorithm in charge of controlling the UoI movement is to achieve the target  $P_T$  following the shortest path while avoiding collisions with mobile and static objects. The shortest path minimizes the distance traveled by the UoI, thus contributing to energy saving.

#### IV. PROBABILITY DISTRIBUTION BASED COLLISION AVOIDANCE (PICA) FRAMEWORK

In this section, we propose a heuristic, dubbed PICA, to control the movement of the UoI to reach a target position while avoiding collisions. It considers the mission time, i.e., the time from the UoI starts its mission until it reaches its target  $P_T$ , is slotted. At each time step  $t$ , the UoI collects data from different sensors to sense the objects' presence in its vicinity. As depicted in Fig.2, UoI is aware of the surrounding objects in a circular area  $z_i$ , whose extension is limited by the sensors' ranges. Specifically, the circular area  $z_i$  has a radius  $\rho$ . The state of the circular area  $z_i$ , i.e., the spatial distribution of the intruders within it, is updated at every time step after the UoI changes its position according to an action  $a \in \mathcal{A}$  taken by PICA. Let  $\mathcal{Z}$  denotes the set of circular shaped areas  $z_i$  where  $\mathcal{Z} = \{z_i : \forall i \in \mathcal{A}\}$ .

Inside every  $z_i$  it might exist intruders neighboring every  $z_i$ 's center  $i$ . Let  $\eta(i)$  denote the set of *static\_intruders* and *mobile\_intruders* inside the area  $z_i$ . We denote by  $\delta_i^j$  the euclidean distance between  $j^{th}$  intruder belonging  $\eta(i)$  and the position of  $i$ . The distance between each intruder and the center of  $z_i$  can be computed by PICA using the triangulation method. The density distribution of the intruders inside  $z_i$  could refer to the likelihood of a collision if the UoI moves to position  $i$ . In other words, the denser  $z_i$  is, the higher the probability that the UoI experiences a collision. Let us define  $\mathcal{P}(z_i)$  as the probability of collision inside  $z_i$ , formally defined as follows:

$$\mathcal{P}(z_i) = \frac{\sum_{\forall j \in \eta(i)} p_i^j}{\sum_i \sum_{\forall j \in \eta(i)} p_i^j} \quad \forall i \in \{1, 2, \dots, |\mathcal{A}|\} \quad (3)$$

where

$$p_i^j = 1 - \frac{\delta_i^j}{\rho} \quad (4)$$

Where,  $p_i^j$  represents the likelihood that UoI collides with  $j$  if action  $i \in \mathcal{A}$  is chosen. Formally, the closer  $j$  is to UoI, the higher the probability of collision. Note that the value of  $\delta_i^j/\rho$  is within the interval  $[0, 1]$ . In order to avoid the collision, the position with the lower  $\mathcal{P}(z_i)$  should be chosen.

In addition to the safety factor, PICA aims to go through the shortest path by seeking at each time step  $t$  the direction

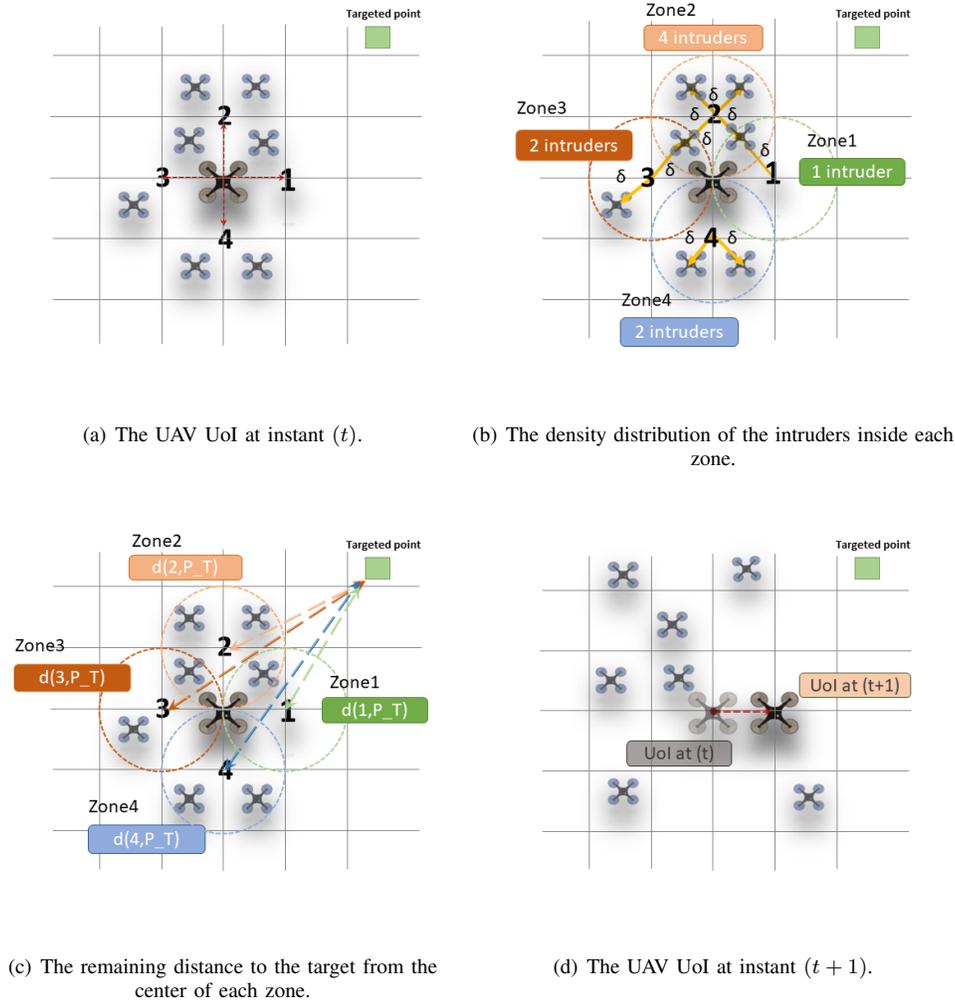


Fig. 3. PICA descriptive example.

that brings the UoI closest to the target. To achieve this goal, PICA measures the remaining distance to the target from every  $i$ . Indeed, to decide the UoI's next direction, PICA ranks every  $Z_i$ 's center,  $i$ , using the following equation:

$$R_i = \theta P(z_i) + (1 - \theta) \frac{\delta_i^{P_T}}{\sqrt{X^2 + Y^2}} \quad (5)$$

, where  $\theta \in [0, 1]$  is a parameter used to favor either safety or energy.  $\delta_i^{P_T}$  denotes the distance between the current position and the target. The concept of application integrating the UoI has an immediate impact on selecting either  $\theta$  or its complement  $(1 - \theta)$ . In our case, we give more priority to the safety of the UoI agent. For this reason, we have selected higher values of  $\theta$ . Furthermore,  $\delta_i^{P_T}$  refers to the euclidean distance between the  $Z_i$ 's center and the targeted point  $P_T$ . Since the unity of both the probabilities and the distances values are in different scales. To prevent distance domination, we have normalized the value of  $\delta_i^{P_T}$  to be between 0 and 1 by dividing it by the maximum possible distance  $\sqrt{X^2 + Y^2}$ . Indeed, the UoI will choose to move in the direction  $a$  that has the lowest rank using the following formula:

$$a = \arg \min_{i \in \mathcal{A}} \{R_i\} \quad (6)$$

Figure. 3 illustrates the PICA functionality, which is detailed in Algorithm 1. For simplicity, in this example, we consider that UoI can move only on four directions  $\{Up, Down, Left \text{ and } Right\}$  corresponding to the positions  $\{1, 2, 3, 4\}$ , respectively. We also consider that all the distances  $\delta_i^j$  are the same and equal to  $\delta$ . Hereafter, based on the density of intruders in each zone, PICA computes the probability of collisions  $\mathcal{P}(z_i)$  for every  $z_i$  (Algorithm 1:line 10) using (3) and (4). For example, the probability of collision  $\mathcal{P}(z_2)$  at the zone 2 is the summation of the probabilities that UoI collides at the  $Z_i$ 's center 2 with every  $j \in \eta(2)$ . Indeed, as shown in Fig.3(b),  $\mathcal{P}(z_2) = \frac{\sum_{\forall j \in \eta(2)} p_2^j}{\sum_{\forall i \in \{1, \dots, 4\}} \sum_{\forall j \in \eta(i)} p_i^j}$  where,  $p_2^1 = p_2^2 = p_2^3 = p_2^4 = (1 - \frac{\delta}{\rho})$ . In this case, the probability collision distribution of the positions 1, 2, 3 and 4 is  $\frac{1}{9}$ ,  $\frac{4}{9}$ ,  $\frac{2}{9}$  and  $\frac{2}{9}$ , respectively. As depicted in Fig.3(b), the zone 2 is the most dense in term of

**Algorithm 1: Probability distribution based collision avoidance (PICA)**

```

Input :
    X: The x_axis limit of the geographical area X.
    Y: The y_axis limit of the geographical area X.
    ρ: The radius.
    θ: The priority rate.
    A: The set of actions.
    PT: The started point of UoI in the environment.
    PS: The targeted point of UoI in the environment.

Output:
    done: The UoI reaches PT or collides with one of the intruders.
1  done ← False;
2  while (done = False) do
3      Z ← ∅;
4      R ← ∅;
5      foreach (a ∈ A) do
6          z ← Circle(a, ρ);
7          Z ← Z ∪ z;
8      end
9      foreach (z ∈ Z) do
10         Compute P(z);
11         Compute δiPT;
12         r ← θP(z) + (1 - θ)  $\frac{\delta_i^{P_T}}{\sqrt{X^2 + Y^2}}$ ;
13         R ← R ∪ r;
14     end
15     a ← arg mini ∈ A {Ri};
16     UoI applies action a;
17     if (U(t) = PT or Collision) then
18         done ← True;
19     end
20 end
    
```

intruders compared to the other ones. By contrast, the zone 1 is the less dense since it contains only one intruder with probability collision  $\frac{1}{9}$ . To rank the candidates' directions  $i \in \{1, \dots, 4\}$  of PICA, based on (5) and (6), chooses the best action that has the lowest probability collision and the lowest remaining distance to the target (Algorithm 1: line 12 and 15). Following the same example and as shown in Fig. 3(b) and Fig. 3(c), the candidate direction (1) has the smallest rank since it has the lowest probability and the lowest missing distance to the target (Algorithm 1: line 15). Finally, at the time step  $(t + 1)$ , the UoI moves into the position (1) as depicted in Fig. 3(d).

**V. RELIANCE: REINFORCEMENT LEARNING BASED COLLISION AVOIDANCE SOLUTION**

In this section, we provide an RL-based solution for avoiding the collision. In contrast to the model-based approach, Markov Decision Process (MDP), which requires full knowledge about the environment (i.e., transition probabilities), RL does not require any prior knowledge. This makes RL a more suitable framework for dealing with unsuspected and uncorrelated mobility of objects around UoI. Thanks to sampling and bootstrapping in RL, RELIANCE can forecast the next movement of each mobile object and then avoid the collision. Fig. 4 depicts the main overview of the RELIANCE solution. In this section's balance, we will give first some background on RL, and, more precisely, DQN employed in this paper. Then, we will give a detailed description of RELIANCE.

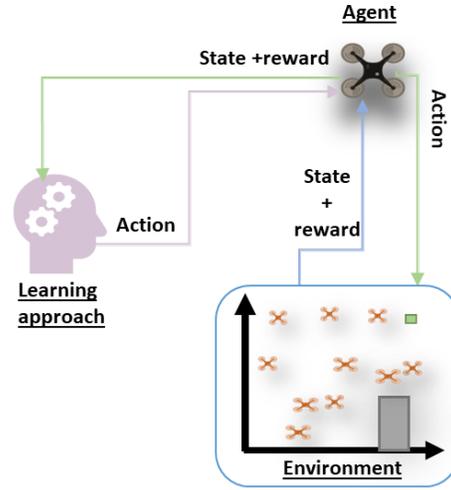


Fig. 4. Reinforcement learning based collision avoidance system overview

**A. Background on RL and DQN**

The reinforcement learning (RL) technique has been widely used in the literature in various applications and services, such as robotics and industry 5.0. RL's ultimate goal is to endow vertical industry with the ability to learn, improve, and adapt according to the environment's changes. With the new trend towards the self-optimized and the cognitive network, industry and academia shifted their attention to employing RL. An RL system mainly consists of four elements, as depicted in Fig. 4, which are: *i*) Environment  $\mathcal{E}$ ; *ii*) States  $\mathcal{S}$ ; *iii*) Agent, in our case, is the motion controller of the UoI; *iv*) Actions  $\mathcal{A}$ ; and *v*) rewards  $r$  received after the execution of each action  $a \in \mathcal{A}$ .

While RL works either for episodic or continuous tasks, in this work, our environment is episodic. Each episode presents UoI mission that starts at  $P_S = (x_S, y_S)$  and ends either when UoI attends its destination  $P_T = (x_T, y_T)$  or collides with a mobile or static obstacle. As depicted in Fig. 4, UoI discretely interacts with the environment by taking different actions, and then accordingly receiving an observation and rewards that reflect the action taken. The agent UoI keeps interacting with the environment  $\mathcal{E}$  and receiving reward  $r_t$  on steps  $t \in \{1, 2 \dots T\}$ . While  $T = \infty$  for continuous tasks, it is finite in the case of an episodic task. The objective of UoI agent is to increase cumulative reward  $G_t$  received after time step  $t$  until the end of the episode.

$$G_t \doteq \sum_{k=0}^T \gamma^k r_{t+k+1} = r_{t+1} + \gamma G_{t+1} \quad (7)$$

, such that  $\gamma \in [0, 1]$  is the discount factor and  $r_t$  denotes the immediate reward received at the instant  $t$ .

Many RL techniques have been proposed in the literature, including policy-based (e.g., REINFORCE), actor-critic (e.g., A3C and DDPG), and value-based approaches (e.g., QN, DQN, and DDQN). While the two formal methods

aim to provide the policy that estimates the state's action probabilities, the latter approach estimates the state-action value. Then, this value is used to deliver the optimal policy. Considering that the space of actions is discrete and limited, in this work, we opt for a value-based approach and, more precisely, DQN. Particularly, we have chosen DQN due to the size of the action-state space, as explained later.

The state-action value of a state  $s \in \mathcal{S}$  using the action  $a \in \mathcal{A}$  under the policy  $\pi$  is defined with the equation 8 [33]:

$$Q_{\pi}(s, a) \doteq E_{\pi}[G_t | S_t = s, A_t = a] \\ \leftarrow E_{\pi}\left[\sum_{k=0}^T \gamma^k r_{t+k+1} | S_t = s, A_t = a\right] \quad (8)$$

The optimal action-state value  $Q^*(s, a)$  can be delivered from  $Q_{\pi}(s, a)$  by choosing the optimal policy. Formally,  $Q^*(s, a) = \max_{\pi} Q_{\pi}(s, a)$  for  $a \in \mathcal{A}$  and  $s \in \mathcal{S}$ .  $Q^*(s, a)$  can be also delivered using Bellman optimality equation using the following formula (9) [33]:

$$Q_*(s, a) \leftarrow E[r_{t+1} + \gamma \max_{a' \in \mathcal{A}} Q_*(S_{t+1}, a') | S_t = s, A_t = a] \quad (9)$$

The basic idea behind many value-based algorithms is sampling and bootstrapping. The sampling is leveraged for enabling the Algorithm to learn by exploring the environment thanks to the trial and error approach. Meanwhile, bootstrapping is a technique used to estimate the state-action value in order to speed up the Algorithm convergence [33]. Q-Learning Algorithm is one of the widely used Algorithm in the literature. Q-Learning Algorithm leverages sampling and bootstrapping methods to converge to the optimal policy. During the learning step, Q-Learning Algorithm updates the state value action using the following formula:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \times [r_{t+1} + \gamma \times \max_{a \in \mathcal{A}} Q(s_{t+1}, a) - Q(s_t, a_t)] \quad (10)$$

, such that  $\alpha$  is the learning rate.

Thanks to bootstrapping, Q-learning repeatedly updates  $Q(s_t, a_t)$  by shifting it towards the optimal value using TD error ( $r_{t+1} + \gamma \times \max_{a \in \mathcal{A}} Q(s_{t+1}, a) - Q(s_t, a_t)$ ) and learning rate  $\alpha$ . This approach enables to gradually increasing  $Q(s_t, a_t)$  towards the optimal value. The optimal policy can be delivered from the optimal state action using the following formula:

$$\forall s \in \mathcal{S} : \pi^*(s) \leftarrow \arg \max_a Q(s, a) \quad (11)$$

In Q-learning Algorithm, the state action value  $Q$  is presented as a table, where the states are the lines and actions are the columns. Unfortunately, Q learning is unfeasible for large action-state spaces as ours. Fortunately, DQN has been

suggested to overcome that limitation [34] by creating an estimator of Q table by leveraging the neural network. In fact, the Q table is approximated with a neural network with parameter  $\omega$ .

Unfortunately, the basic DQN Algorithm suffers from overestimations of action value due to using the same neural network in the update and estimation of the next Q value used to compute the TD error. This approach creates lots of noise and makes it hard to find the action with maximum expected/estimated Q-value. To prevent this issue, authors in [35] have suggested DQN (DQN) that uses two different Q neural networks. The first one, called policy Q neural network  $Q^{\pi}$ , is used for estimating the action. Meanwhile, the second one, called target Q network  $Q^T$ , is used to generate the target action values. To mitigate the noises in the update, while  $Q^{\pi}$  is updated at each iteration,  $Q^T$  is updated from  $Q^{\pi}$  only after a specific number of episodes.

In this case, the policy  $\pi$  during the exploitation, either during the training or inference modes, is generated from the approximation  $Q^{\pi}(s_t, a_t, \omega)$  neural network using the following equation (12):

$$\pi_{s_t} \leftarrow \arg \max_{a \in \mathcal{A}} Q^{\pi}(s_t, a, \omega) \quad (12)$$

Meanwhile, the parameter  $\omega$  (bias and weights) of the estimator  $Q^{\pi}$  is updated periodically during the training step using the following formula:

$$\omega_{t+1} = \omega_t + \alpha \times [r_{t+1} + \gamma \times \max_a Q^T(s_{t+1}, a; \omega') - Q^{\pi}(s_t, a_t; \omega)] \times \nabla_{\omega} Q^{\pi}(s_t, a_t; \omega) \quad (13)$$

By leveraging different gradient descent methods (e.g., stochastic gradient descent, RMSprop, or ADAM), the DQN Algorithm keeps updating  $\omega$  during the training step.  $\omega$  is updated from replay memory ( $\mathcal{B}$ ) that consists of transitions observed during the exploration or exploitation. Each transition  $\langle s_t, a_t, r, s_{t+1} \rangle$  consists of the current state  $s_t$ , the taken action  $a_t$ , the immediate received reward  $r$  and the next state  $s_{t+1}$ . To break the correlation between transitions and to allow a stable learning curve, a batch of transitions (i.e., batch\_size) are randomly selected from the replay memory  $\mathcal{B}$  [33].

To ensure a balance between the exploration and exploitation during the training to update  $\omega$ , an epsilon greedy method is used. The Algorithm keeps randomly switching between the exploration and exploitation modes. At the end of the training, the DQN Algorithm should favor exploitation than exploration to assist its convergence. For this purpose, an epsilon decay strategy has been adopted by decreasing the epsilon decay  $\xi$  parameter during the training.  $\xi$  initially starts by 1, and it should converge to zero at the end of the training. To switch between the exploration and exploitation, a random number (i.e.,  $[0, 1]$ ) is generated and compared to  $\xi$ . If the generated number is lower than  $\xi$ , the exploration procedure is executed. Otherwise, the exploitation procedure is considered.

## B. RELIANCE Model Overview

The autonomous flight control system of the UoI is realized as a Reinforcement Learning (RL) agent. To model the energy-aware collision avoidance problem using the RL framework, as mentioned in the previous subsection, the following elements need to be formally defined: i) environment, state, agent, actions, and reward.

- *Agent*: The RL agent is instantiated and run within the UoI to control its trajectory in order to avoid collisions while minimizing the energy consumption by taking the shortest path until its destination.
- *Environment*: Geographical are of dimensions  $X \times Y$  that include a set of mobile (e.g., other UAVs) and static objects (e.g., walls). The agent moves within this 2D confined area.
- *Actions*: The action space comprises a set of eight directions, i.e.,  $\mathcal{A} = \{Up, Down, Right, Left, Up\_right, Up\_left, Down\_right, Down\_left\}$ , as previously mentioned.
- *Reward*: If the agent succeeds and reaches its targeted destination, it is positively rewarded with 100. If the UoI experiences a collision during its trajectory to the destination, the agent is penalized with a negative reward of  $-100$ . Finally, in order to encourage the agent to take the shortest path, there is a penalty of  $-0.1$  for each step taken by the agent until reaching its destination.
- *State*: The state (agent's observations) consists of two parts:
  - i) the distance vector  $P_{rel} = (x_{rel}, y_{rel})$  defined as the vector from the current UoI's position  $(x_C, y_C)$  to the UoI's destination  $(x_T, y_T)$ . That is:

$$x_{rel} = \frac{x_T - x_C}{X}$$

$$y_{rel} = \frac{y_T - y_C}{Y}$$

Please observe that  $(x_{rel}, y_{rel}) = (0, 0)$  means the UoI is at the destination. Also, note that  $x_{rel}$  and  $y_{rel}$  have been normalized by  $X$  and  $Y$  (flight area dimensions), respectively. In this way, the agent is agnostic to the scenario scale, which makes the solution more general, i.e., the same trained model can be used in many environments.

- ii) The number of mobile and static objects distribution across a grid square centered around the UoI. Specifically, a  $\Delta = \|\|L_\Delta \times L_\Delta\|\|$  grid square is considered. The UoI dimensions give the size of each tile of the grid. The grid is formally described as a binary matrix. Each element of this matrix indicates whether there is any static or mobile object (intruder) within the respective cell (tile) ( $=1$ ) or not ( $=0$ ) (see Fig. 5). This approach enables us to consider the UoI vicinity, which is the most relevant to avoid collisions and reduce the state space by aggregating many observations to a single

state. Thus, faster learning and convergence will be perceived. As depicted in Fig. 5, thanks to the aggregation method adopted by RELIANCE, different observations depicted in Figures 5(a), 5(b) and 5(c) can be presented by the same state shown in Fig. 5(d).

It is remarkable that both components of the state considered are agnostic of the scenario, which makes the solution more general.

## C. RELIANCE example description

As stated previously, we provide the agent with an RL-Algorithm that adopts one of the existing RL approaches. The principal role of this Algorithm consists of giving the agent the ability to self-decide in which direction has to move following the defined goals. Before starting the flight, we provide the agent with the coordinates of its started and targeted points,  $P_S, P_T$ , respectively. Thus, at each step, the agent needs to do the following:

- Receives the status of the area from the sensing equipment.
- Traces the square-shaped area surrounded the agent where the current agent position centers the square.
- Ignores the area beyond the square and uses the received sensing status to update every cell inside the square by (1) if it contains an intruder and (0) otherwise.
- Computes the relative address  $P_{rel}$  of the targeted point  $P_T$  in proportion to the current position of the agent.
- Normalize the value of the targeted point relative position.
- Generates the current state  $S_t$  where  $S_t = \{P_{rel} = (x_{rel}, y_{rel}), \Delta\}$ .
- Uses the prior learned knowledge to choose the best action (direction) that might allow the agent to not collide with any nearby intruder and get closer to its target.
- Apply the selected action and observe the impact of the agent's dynamic on the environment by recording the new knowledge in terms of reward earned and new agent position.
- The agent repeats the previous steps until reaching the targeted point  $P_T$  or an instance of collision occurs.

Indeed, as shown in Fig.5 the state  $S_t$  introduced at each step highly impacts the learning process and the action selection. Thus, the state's definition needs to be done based on clear and logical arguments. In what follows, we detail our logic behind the definition of the state  $S_t$ .

First, instead of using the 2-d coordinate of the target, the use of the relative address allows the agent to move in the direction of the targeted position wherever its position is in the environment. Furthermore, the use of the relative address allows the agent to involve the remaining distance to the target in the learning process of the agent. Thus, we keep the agent seeking the shortest way to move on.

On the other side, the focus on the square surrounding the agent position instead of considering the whole environment

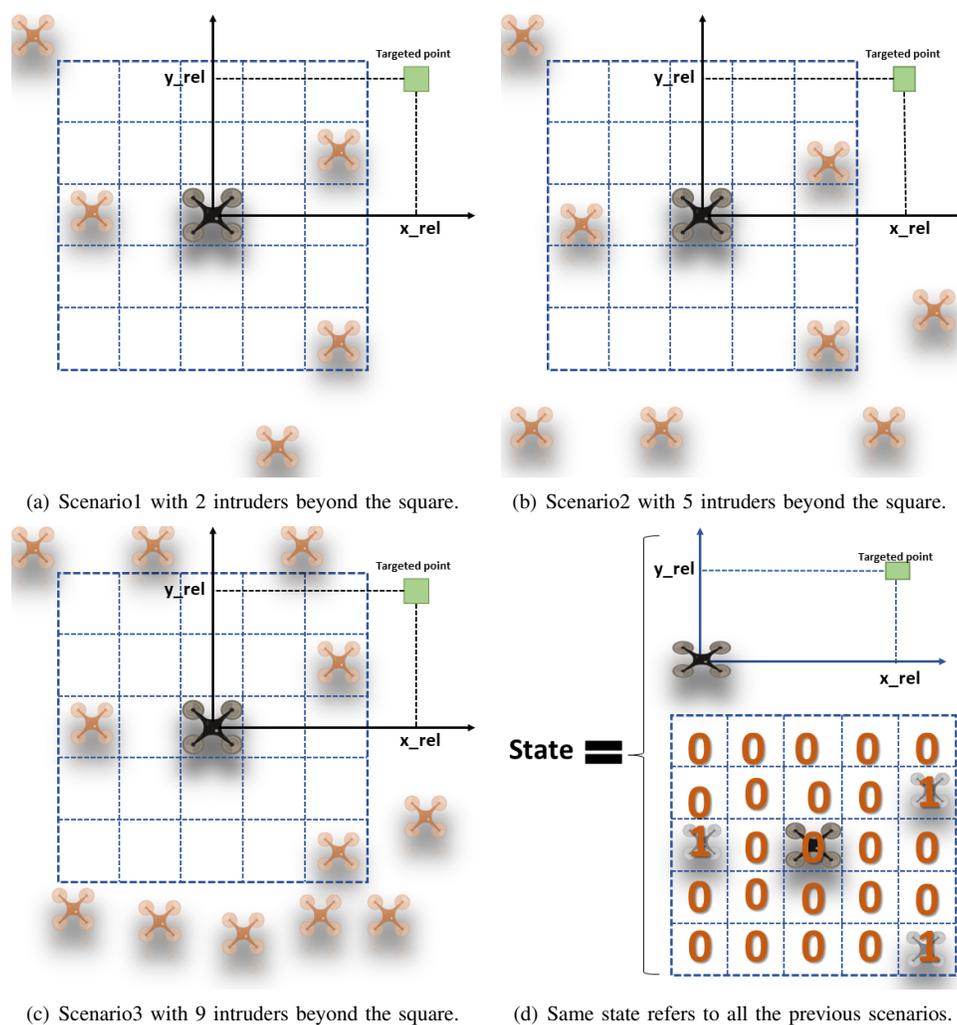


Fig. 5. State aggregation process adopted by RELIANCE.

area aims to aggregate many environment states to one state at the agent. Let  $\Delta$  denote the surrounding area of the agent. The size of the surrounding area is a hyperparameter that can be tuned during the training step. As shown in Fig.5, three different environment states can be presented by the same state. However, they are the intruders' position beyond the square, and the state is the same since it considers only the intruders positioned inside the square and the relative targeted position. Indeed, the agent will know a limited number of states. Furthermore, the non-use of related intruders features (e.g., intruders coordinates) aims to have a generalized view that might be used in similar status even with different intruders' positions. Finally, the use of such state might help improve the learning convergence time of the agent. Moreover, this approach helps to train the agent on a limited number of intruders and can also be functional in an environment with a high number of intruders.

#### D. RELIANCE DQN Algorithm

Throughout this section, we detail the RL approach used by our based RL agent. Several RL approaches exist in

the literature, such as the Q-learning, Deep Neuron network (DQN), Policy gradient, and Actor-critic. However, the elements state-space and the action-space from the RL modeled system are either deterministic or continuous, highly impacting the selection of the RL approach. In our case, eight deterministic actions compose the action-space. On the other side, the state-space contains two deterministic elements: the relative address and the square shaped area. The Q-learning approach requires that the agent is within a deterministic limited space. Indeed, the Q-learning approach seems to be the most suitable for our problem. However, the size of the square-shaped area might be considerable. Then the agent could take a long time to converge, and consequently, the computation process will consume more resources. Furthermore, the learning process can be less efficient by getting a useless action.

To mitigate this problem, we opted to use the Deep Reinforcement Learning (DQN) approach. Thus, more details about our based DQN autonomous UAV collision avoidance and energy-aware agent are summarized in Algorithm 2. First, the agent starts by instantiating two different neural

---

**Algorithm 2:** RELIANCE: reinforcement learning based collision avoidance solution

---

**Input :**  
 $Q^\pi$  and  $Q^T$ : The initialized policy and target networks using Xavier-uniform.  
 $\mathcal{B}$ : The batch replay memory size to size  $N$ .  
 $batch\_size$ : The size of each batch.  
 $\mathcal{M}$ : A number of episodes to update  $Q^T$  with  $Q^\pi$ .  
 $\xi_0$ : The initial value of epsilon greedy.  
 $\mathcal{N}$ : Number of episodes.

**Output:**  
 $done$ : The UoI reaches  $P_T$  or collides with one of the intruders.

```

1 episode = 1;
2 while episode ≤ N do
3     done ← False;
4      $\xi \leftarrow \frac{\xi_0}{\xi_0 + episode}$ ;
5      $S_0 = \mathcal{E}.init()$ ;
6     while done = False do
7          $S_t = \{P_{rel} = (x_{rel}, y_{rel}), \Delta\}$ ;
8         if random() ≤  $\xi$  then
9              $a \leftarrow randint(\mathcal{A})$ ;
10        else
11             $a \leftarrow \arg \max_{a \in \mathcal{A}} Q^\pi(S_t, a)$ ;
12        end
13         $S_{t+1}, reward, done \leftarrow \mathcal{E}(S_t, a)$ ;
14         $\mathcal{B} \leftarrow (S_t, a, reward, S_{t+1}, done)$ ;
15         $t \leftarrow t + 1$ ;
16        if size( $\mathcal{B}$ ) ≥ batch_size then
17            mini_batch ← random( $\mathcal{B}$ , batch_size);
18            foreach
19                 $(S_i, a_i, reward_i, done_i, S'_i) \in mini\_batch$  do
20                if ( $done_i = True$ ) then
21                     $y_i \leftarrow reward_i$ ;
22                else
23                     $y_i \leftarrow reward_i + \gamma \max_{a' \in \mathcal{A}} Q^T(S'_i, a'_i)$ ;
24                end
25            end
26             $\mathcal{L} = \frac{1}{N} \sum_{i=0}^{N-1} (Q^\pi(S_i, a_i) - y_i)^2$ ;
27            Update  $\omega$  of  $Q^\pi$  using  $\mathcal{L}$ ;
28        end
29        if episode %  $\mathcal{M} = 0$  then
30             $Q^T \leftarrow Q^\pi$ ;
31        end
32        episode = episode + 1;
33    end

```

---

networks (NNs), named policy-network,  $Q^\pi$ , and target-network,  $Q^T$ . To ensure the fast convergence of the Algorithm, we have used Xavier initialization to initialize the weights of both neural networks  $Q^\pi$  and  $Q^T$ . The Xavier initialization helps to converge fast and prevent the exploding and vanishing gradients during the training process. To give the agent more time to explore the behavior of the actions set, we opted to use the decayed  $\epsilon$ -greedy strategy. The Algorithm starts from the first episode and ends at the last episode  $N$  (Algorithm 2:lines 1 – 2). For each episode (Algorithm 2:lines 2 – 33), RELIANCE does the next steps. Initially, the episode sets to an undone state (Algorithm 2:line 3). Then,  $\xi$  is initialized to enable either the exploration or exploitation (Algorithm 2:line 4). Later, the environment is initialized by creating a new mission to train the agent (Algorithm 2:line 5).

While the episode is not completed (UoI achieves the tar-

get or collides), we do the following steps (Algorithm 2:lines 6 – 28): First, the agent generates and normalize the current state (Algorithm 2:line 7). Then, according to the decayed value of  $\xi$  and a randomly generated number, we select either exploration or exploitation (Algorithm 2:lines 8 – 12). If the exploration is selected, a random action is issued from  $\mathcal{A}$  (Algorithm 2:lines 8 – 10). Otherwise, the agent of the UoI chooses the action with the maximum reward previously earned using the policy network (Algorithm 2: lines 10–12). After the agent applies the selected action and saves the transition to  $\mathcal{B}$  (Algorithm 2:lines 13 – 14), the agent moves to the new observed state (Algorithm 2:line 15). However, when the number of experiences exceeds the  $batch\_size$ , the agent selects a random batch of transitions from  $\mathcal{B}$  to update the  $Q^\pi$  following  $TD(0)$  (Algorithm 2:lines 16–27). The agent keeps updating the  $Q^T$  using  $Q^\pi$  every  $\mathcal{M}$  steps. The agent repeats the previous steps until the end of all the episodes in the training.

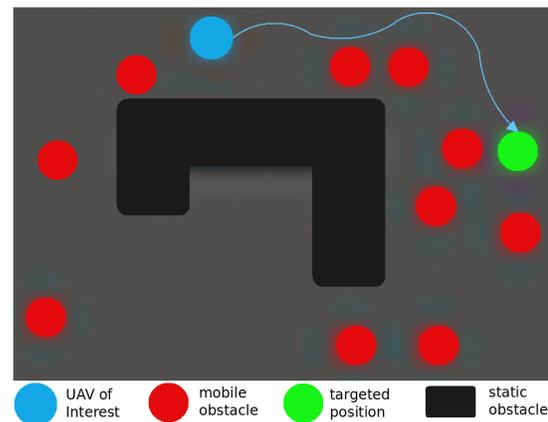


Fig. 6. OpenAI Gym compliant simulator

## VI. EXPERIMENTATION AND RESULTS

In this section, we evaluate the performances of our two solutions PICA and RELIANCE. In the balance of this section, we first present the simulation setup; then, we present the convergence of RELIANCE during the training mode. Last but not least, we conclude this section by evaluating the performances of RELIANCE in inference mode to PICA.

### A. Simulation Setup

Existing UAV simulators, such as Air Sim and Software in the Loop (SITL), use telemetry data to control the motion of a single UAV in a closed and well-controlled environment. These simulators mainly focus on telemetry data to maintain a single UAV for landing and flying. Still, they did not consider mobile objects, which is a handicap facing their utilization to evaluate PICA and RELIANCE solutions. In order to overcome these limitations, we have developed an OpenAI Gym [36] compliant simulator with graphical rendering capability using Python language and OpenCV

<https://youtu.be/7UcRxfAREw>

library. This simulator provides a customizable environment that considers both static (e.g., building) and dynamic (e.g., UAVs and birds) obstacles. The static obstacles can be included in a JSON format to the simulator. In the simulator, we have adopted a discrete-time implementation of the events (e.g., UAVs mobility). This strategy helps to reduce the simulation time significantly by considering only the counted events rather than using real execution time. To make the proposed framework orthogonal on agents (PICA, RELIANCE...) implementation, we have developed a complete framework that consists of an abstraction layer of the agent and environment. We have also designed RELIANCE and PICA to be transparent in the environment and easily adapted to other simulators or real experiments later. Thus, we believe the suggestion of this simulator will have an added value to the scientific community. While PICA has been implemented using Python and Numpy, the neural network model of RELIANCE is implemented with Python and Pytorch library.

Besides the UoI, the environment also consists of a set of customizable number of *static\_intruders* and *mobile\_intruders*. As aforementioned, both UoI and *mobile\_intruders* move in a 2D plan using eight possible actions. The *mobile\_intruders* move in the simulator using a random walk technique. In contrast to *mobile\_intruders*, the UoI moves under the control of either PICA or RELIANCE agents. The rendering environment consists of a gray screen with black rectangles and blue, green, and red circles. As depicted in Fig. 6, the black rectangles and red circles refer to the static obstacles and the intruder(s), respectively. Meanwhile, the blue and green circles refer to the UoI and its targeted position, respectively. The simulation runs in episodes, such that each of which ends when UoI collides or reaches the target.

### B. RELIANCE Training mode

The training of the RELIANCE model happens using 14 Dual Intel Xeon E5-2680 v3 @ 2.5 GHz, with 117 GB of RAM, one Nvidia P100 GPU, and running CentOS 7. During the training process, we have fixed the size of the simulation area by  $20 \times 20$  and considered 10 *mobile\_intruders* besides three static obstacles with different shapes and sizes. We have fixed the hyper-parameter surrounding area  $\Delta$  of the agent by  $5 \times 5$  after performing a set of different tests. To ensure fast convergence without underfitting or overfitting, we have tuned the neural network hyper-parameters used by RELIANCE. We have performed many experimental tests before fixing the hyper-parameters. We have fixed the discount factor  $\gamma$  by 0.95 and the learning rate  $\alpha$  by  $10^{-4}$ . We have also used two fully-connected hidden layers in which the number of units (i.e., activation functions) is 40. We have also tested with 400 units in each layer. However, similar convergence rate is perceived. We have also observed similar performances during the inference mode. We adopted the Rectified Linear Unit (ReLU) activation function for both hidden and output layers. A Xavier initialization has been adopted to initialize the neural network units in the model.

This initialization helps to converge fast and prevent the exploding and vanishing gradients during the training process. During the training, we have used *batch\_size* = 1024, replay buffer size = 500000 and target update = 8 to update the weight of target network  $Q^T$  from the policy network  $Q^\pi$ .

As depicted in Fig. 7, we have conducted two sets of experiments. Initially, we have trained one RELIANCE agent as depicted in Fig. 7(a) for a period of 2000 episodes. In this figure, while the blue curve shows the cumulative reward gained at the end of each episode, the red one shows the average of the last 50 cumulative rewards. From this figure, we observe that the RELIANCE agent converges at 600 episodes. Starting from that point, the RELIANCE agent succeeds in most of cases to achieve the target without any collision. A live video has been recorded that shows the convergence of RELIANCE.

Meanwhile, in Fig. 7(b), we have evaluated RELIANCE agent's stability. The neural network's bias and weights are randomly initialized in the RELIANCE agent, affecting the training convergence. Moreover, at each episode, the starting and target point of UoI, and the mobility of *mobile\_intruders* are randomly generated. In Fig. 7(b), we have trained 40 RELIANCE agents, simultaneously. In this figure, we have evaluated both the average and the cumulative variance reward achieved. We observe that all the agents succeeded in converging by getting almost the total possible reward after only 400 episodes. Also, we observe that the variance between the trained agents is close to zero, which confirms the algorithm's convergence.

### C. PICA and RELIANCE performance evaluation during the inference mode

In this subsection, we evaluate the performances of RELIANCE in the inference mode against the PICA solution. We simulate  $10^3$  episodes and compare the two solutions in terms of the following metrics:

- **Percentage of collision:** is defined as the percentage of times that the UAV agent collides with *static\_intruders* or *mobile\_intruders*. This metric shows the percentage of time that the UAV agent fails to achieve its final destination;
- **PDF of extra traveled distance:** shows the extra distance needed by a UAV to prevent collisions. This metric shows the probability of a distribution function (PDF) of the extra distance traveled to avoid collisions. In fact, the energy consumption in the UAVs is proportional to the traveled distance before attending the target location. Overall, the more the traveled distance is, the higher energy consumption becomes;
- **PDF of the number of success before a failure:** It shows the PDF of the number of successes arriving at the target before the failure, i.e., the UoI collides with any object. In other words, this metric shows the

<https://youtu.be/5ULpSuMdrSE>

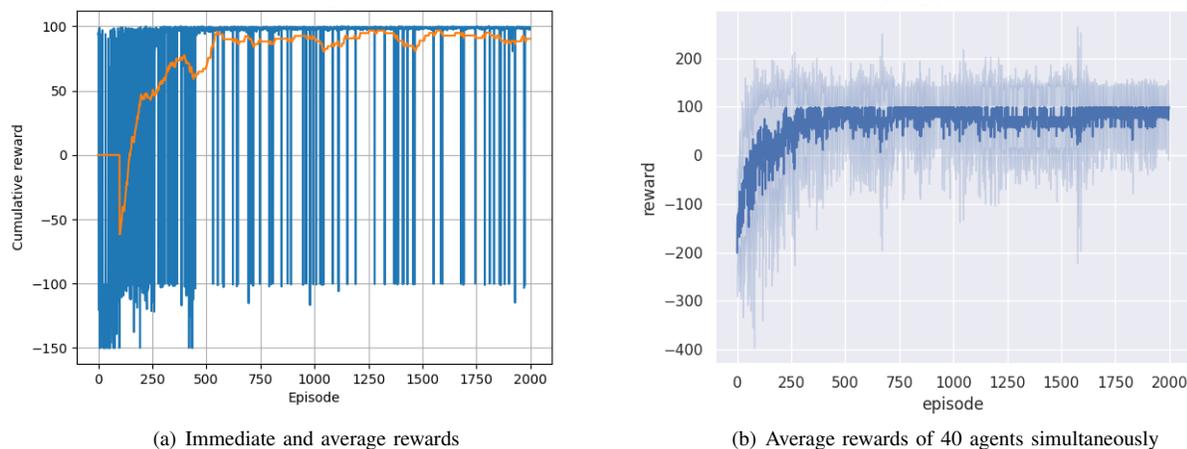


Fig. 7. Convergence evaluation of RELIANCE during the training mode

capability of each solution for traveling consecutive missions without any collision.

To assess the generalization capability of RELIANCE, we have considered three different scenarios during the inference mode. As aforementioned, we have trained RELIANCE agent against static obstacles and 10 *mobile\_intruders*. In contrast, during the inference mode, besides the three static obstacles, we have evaluated the performance of PICA and RELIANCE agents against 5, 10, and 20 *mobile\_intruders*, respectively. The idea behinds these three scenarios is to show the capability of RELIANCE to outlive in unfamiliar environments by leveraging the effectiveness of surrounding area  $\Delta$  and aggregated state.

1) *Percentage of collision*: Figure 8 shows the percentage of collisions as a function of the number of episodes. Both solutions have been evaluated in harsh conditions by including static obstacles and many *mobile\_intruders* in a small area with dimensions  $20 \times 20$ . Moreover, the *mobile\_intruders* move randomly without any correlation, which makes hard to predict their next movement<sup>1</sup>. The first observation that we can draw from this figure is that RELIANCE ensures the generalization by behaving well in unseen environments (e.g., 5 and 20 *mobile\_intruders*). We also observe that whatever the scenarios (5, 10, or 20 *mobile\_intruders*), the RELIANCE offers better performances than PICA. As expected, we also observe that the number of *mobile\_intruders* has a negative impact on the number of collisions as shown in figures 8(a), 8(b) and 8(c), respectively.

For 5 *mobile\_intruders* as depicted in Fig. 8(a), regardless of the number of episodes, the PICA agent has arrived at the target without collision with 70% of success. Whereas, RELIANCE agent has succeeded with 95% to reach the target while avoiding the collisions. Increasing the number of *mobile\_intruders* to 10 hurts the collision percentage in the network as depicted in Fig. 8(b). We observe that the percentage of cases that the UAV agent arrives at the target without collisions drooped out from 70% and 95% to 65%

and 90% for PICA and RELIANCE, respectively. Finally, as depicted in Fig. 8(c), we observe that the increase of the number of *mobile\_intruders* to 20 leads to reduce the percentage of success to arrive at the destination without collisions to 60% and 80% for PICA and RELIANCE, respectively.

The better performances achieved by RELIANCE compared to PICA can be explained as follow. In both solutions, the algorithm controlling the UoI makes the decisions relying on the snapshot from the environment to avoid collisions. The environment snapshot refers to the surrounding area of UoI that is defined by  $\Delta$  and  $\mathcal{Z}$  in RELIANCE and PICA, respectively. On the one hand, based on this snapshot, PICA takes the action that minimizes the likelihood of collisions in  $\mathcal{Z}$ . Nonetheless, PICA does not consider the dynamics of the mobile intruders within  $\mathcal{Z}$ . In contrast, RELIANCE using the DRL approach can learn the temporal correlation between different snapshots  $\Delta$  and therefore make more effective decisions to avoid mobile intruders. This fact explains why PICA exhibits a higher number of collisions compared to RELIANCE.

2) *PDF of extra traveled distance*: Figure 9 shows the performances of PICA and RELIANCE related to energy saving. Unfortunately, avoiding the collision comes with an unavoidable overhead in terms of the extra distance traveled by the UoI. This figure shows this extra distance compared to the traveled distance in the straight travel, i.e., the euclidean distance between the UoI starting and target points. We have estimated the probability distribution function (PDF) of the PICA and RELIANCE extra distances from the results from  $10^3$  episodes. To that end, we employed the KernelDensity function from `sklearn.neighbors`. We observe that regardless of the number of mobile intruders (5, 10, or 20), the percentage of extra distance does not exceed 60%, i.e., the UoI travels 1.6 times the distance of the optimal path.

Figures 9(a) and 9(b) show the PDF of the extra traveled distance for 5 mobile intruders. From Fig. 9(a), RELIANCE succeeded in almost 90% of cases to add only 35% of

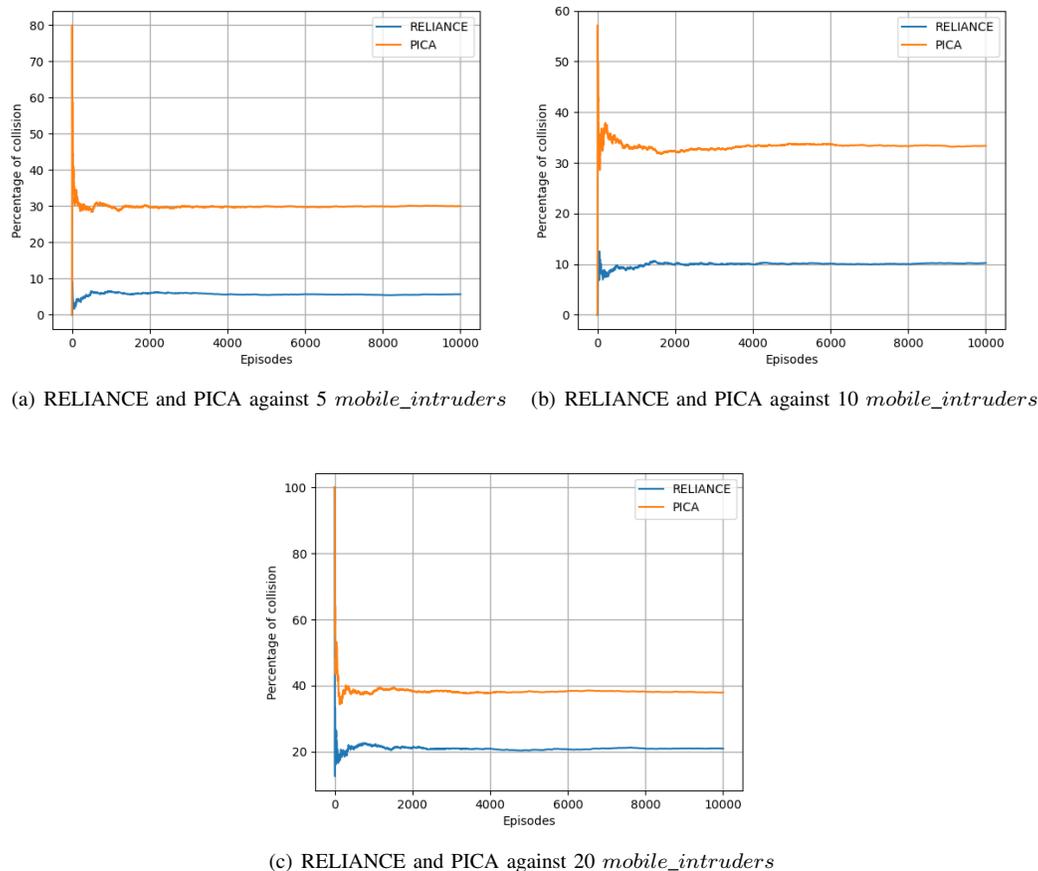


Fig. 8. Percentage of collision in PICA and RELIANCE solutions.

extra distance compared to the optimal one (i.e., 1.35 times). With more than 0.08 probability, RELIANCE succeeded in traveling the distance with less than 10% extra distance. We also observe from 9(b) that PICA succeeded in reaching the target in 70% of cases without adding any extra distance. Also, most of the extra distance of PICA does not exceed 40%. Observe that PICA offers shorter extra traveled distances than RELIANCE, which, overall, translates into energy saving. However, this reduction in the traveled extra distance offered by PICA is at the cost of a higher probability of collision, as discussed previously.

Figures 9(c) and 9(d), and 9(e) and 9(f) show the PDF of extra traveled distance for 10 and 20 mobile intruders, respectively. Similar to the case of 5 mobile intruders, we observe that the PICA algorithm succeeded in most of the cases without adding any extra distance. Interestingly, we observe that increasing the number of mobile intruders reduces the extra traveled distance offered by the PICA solution. This can be explained as follow, in the simulation, the extra distance of incomplete mission are filtered (not considered). At each episode, the starting and target point (i.e, mission) of UoI are randomly generated. In fact, increasing the number of intruders will create more collisions on the long distance missions comparing to the short ones. Hence, more short distance mission will participate for generating the PDF

of extra distance. Usually, the probability of adding extra distance in shorter mission is lower than the longer ones, which positively affects the PDF of extra distance metric. Meanwhile, from Figures 9(a), 9(c) and 9(e), we observe that similar behavior in terms of extra traveled distance. The RELIANCE solution succeeded to save long distance mission, however with unavoidable extra distance.

The extra traveled distance and percentage of collision are two contradictory objectives. The lower percentage of collision is, the higher likelihood of extra traveled distance becomes. While the PICA solution leverages a probabilistic approach by considering only one snapshot of the environment, RELIANCE employs DRL to make the correlation between snapshots and then takes the decisions that consider the mobility of intruders. The safety level, i.e., low probability of collision with surrounding intruders, offered by RELIANCE is at the expense of traveling longer extra distances.

3) *PDF of the number of success before a failure*: Figure 10 depicts the PDF of the number of success before a failure happens. It shows the PDF of the number of hits arriving at the target before the collapse. This metric shows the capability of each solution for traveling consecutive missions without any collision. We have conducted three sets of experiments by varying the number of *mobile\_intruders*

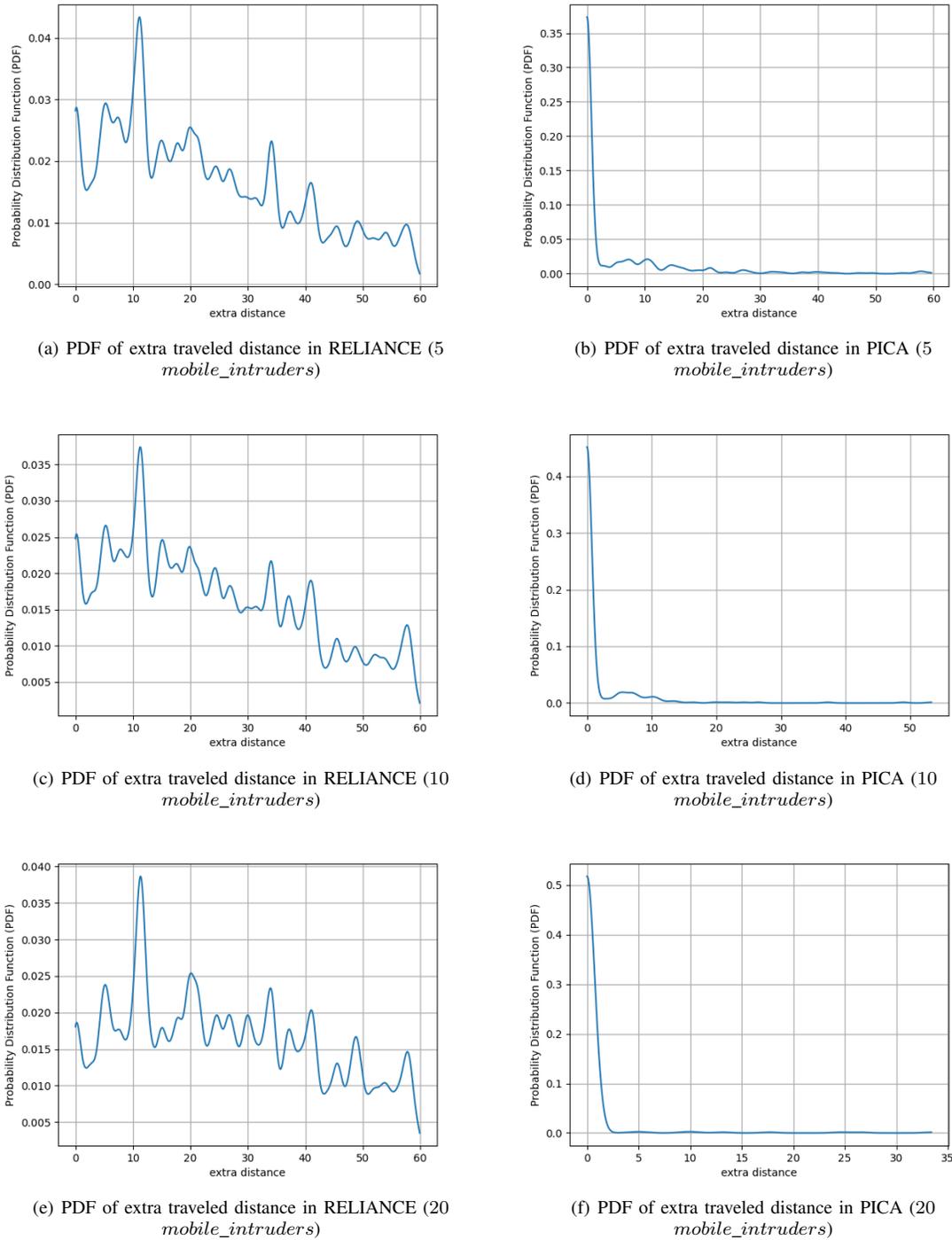
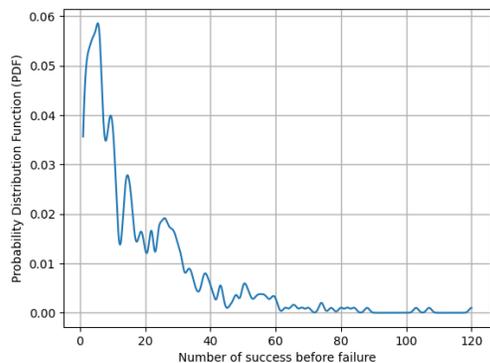


Fig. 9. Probability distribution function of extra traveled distance

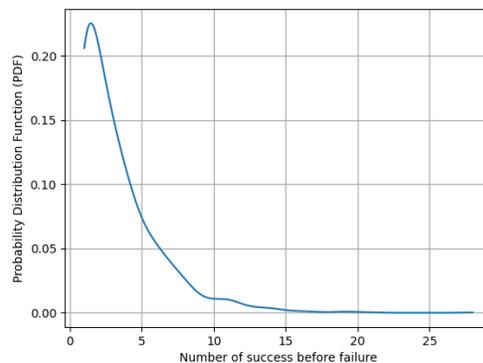
from 5, 10 and 20, respectively. The first observation that we can draw from this figure is that the RELIANCE solution performs better than the PICA solution. Also, we observe that the number of *mobile\_intruders* harms the number of successes before failure.

Figures 10(a) and 10(b) show the performances of PICA and RELIANCE when 5 *mobile\_intruders* is considered. As depicted in these figures, while RELIANCE succeeded in

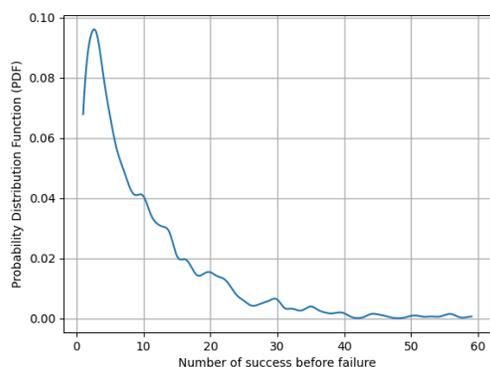
getting 120 successful episodes achieving the target safely, PICA succeeded in achieving the target in 30 episodes without any single failure. We also observe that RELIANCE's probability of fewer than five times consecutively arriving at the target without interruption does not exceed 12%. Meanwhile, in PICA, UoI with a probability of almost 1 does not exceed the 15 episodes consecutively. Figures 10(c) and 10(d) show the impact of 10 *mobile\_intruders* on PICA



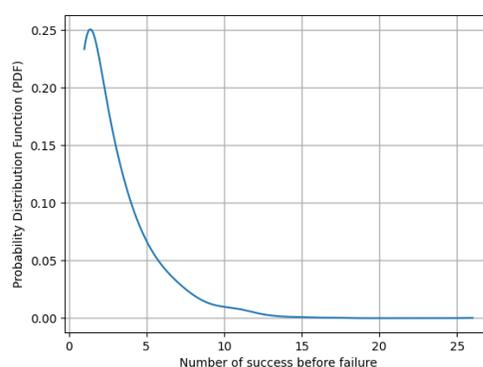
(a) PDF of number of success before a failure in RELIANCE (5 intruders)



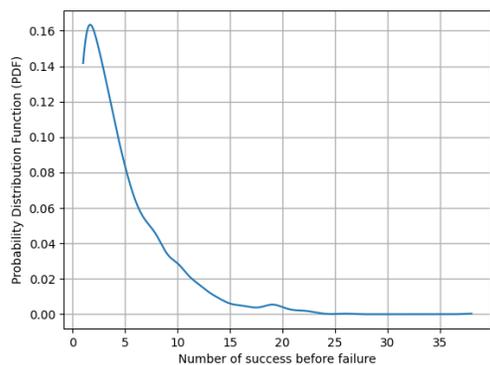
(b) PDF of number of success before a failure in PICA (5 intruders)



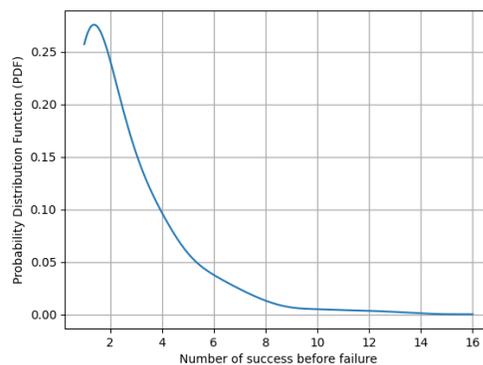
(c) PDF of number of success before a failure in RELIANCE (10 intruders)



(d) PDF of number of success before a failure in PICA (10 intruders)



(e) PDF of number of success before a failure in RELIANCE (20 intruders)



(f) PDF of number of success before a failure in PICA (20 intruders)

Fig. 10. Probability distribution function of extra traveled distance

and RELIANCE. We can observe that the increase in the number of *mobile\_intruders* hurts the number of successes before failure. In RELIANCE, the number of successful episodes before a failure is dropped from 120 to 60. Also, the probability of five consecutive times arrive at the target without interruption does not exceed 20%. Finally, Figures 10(e) and 10(f) show the impact of 10 *mobile\_intruders* on the two solutions. We observe that in  $\frac{1}{3}$  of cases RELIANCE,

the number of success before a failure: does not exceed the threshold 5. Meanwhile, for PICA, the agent with almost probability 1 does not succeed to exceed 15 episodes.

## VII. CONCLUSION

The new enthusiasm for extending UAV commercial operations to cover the urban and populated area controlled airspace beyond visual line of sight (BVLOS) comes with

unavoidable challenges related to object detection and collision avoidance. In this paper, we suggested two solutions named: *i*) **Probability distribution based collision avoidance framework (PICA)**; *ii*) **RL-based collision avoidance framework (RELIANCE)**. While the PICA solution leverages the probability density for avoiding collisions, RELIANCE uses the DQN technique to prevent collisions while saving energy consumption. We have also developed an OpenAI Gym [36] compliant environment<sup>1</sup> with graphical rendering capability using Python language and OpenCV library to evaluate these two solutions. We have developed a complete framework that includes an abstraction of the environment and agent. Our plan to make the platform's code source, including PICA and RELIANCE agents, public for the research community.

We have simulated the agent in the context of both PICA and RELIANCE under similar circumstances. The agent behaves successively following PICA or RELIANCE to prevent the collision and save energy consumption. We have evaluated both protocols in known and unknown environments to assist their generalization capability. The obtained results demonstrate their capacity for generalization. Also, they show the superiority of RELIANCE over PICA in terms of collision avoidance. Also, the simulation results demonstrate the convergence of RELIANCE during the training process<sup>2</sup>.

As a future research direction, we plan to consider other RL Algorithms, including *i*) Policy gradient method, such as RELIANCE; Actor-Critic approach including but not limited to A3C, Deep Deterministic Policy Gradient (DDPG), Trust Region Policy Optimization (TRPO) and Proximal Policy Optimization (PPO). Also, we plan to consider more complex scenarios by considering the velocity and the acceleration of UAVs. A real deployment implementation is envisaged of RELIANCE by leveraging the UAVs available in our lab.

## REFERENCES

- [1] N. Hossein Motlagh, T. Taleb, and O. Arouk, "Low-altitude unmanned aerial vehicles-based internet of things services: Comprehensive survey and future perspectives," *IEEE Internet of Things Journal*, vol. 3, no. 6, pp. 899–922, 2016.
- [2] D. Jaiswal and P. Kumar, "Real-time implementation of moving object detection in UAV videos using gpus," *J. Real Time Image Process.*, vol. 17, no. 5, pp. 1301–1317, 2020.
- [3] S. Ouahouah, J. Prados-Garzon, T. Taleb, and C. Benzaid, "Energy-aware collision avoidance stochastic optimizer for a uavs set," in *2020 International Wireless Communications and Mobile Computing (IWCMC)*, 2020, pp. 1636–1641.
- [4] T. Taleb, K. Ooi, and K. Hashimoto, "An efficient collision avoidance strategy in its systems," in *Proc. of IEEE Wireless Communications and Networking Conference, WCNC, Las Vegas, USA*, March 2008.
- [5] T. Taleb, S. Dutta, A. Ksentini, M. Iqbal, and H. Flinck, "Toward an effective risk-conscious and collaborative vehicular collision avoidance system," *IEEE Communications Magazine*, vol. 55, no. 3, pp. 38–43, March 17.
- [6] F. Giancarmine, A. Domenico, M. Antonio, C. Ciro, C. Umberto, C. Federico, and L. Salvatore, "Multi-sensor-based fully autonomous non-cooperative collision avoidance system for unmanned air vehicles," *Journal of Aerospace Computing Information and Communication*, vol. 5, no. 10, pp. 338–360, 2008.
- [7] R. Sharma and D. Ghose, "Collision avoidance between uav clusters using swarm intelligence techniques," *International Journal of Systems Science*, vol. 40, no. 5, pp. 521–538, 2009.
- [8] Y. K. Kwag and C. H. Chung, "Uav based collision avoidance radar sensor," in *Proc. of 2007 IEEE International Geo-science and Remote Sensing Symposium (IGARSS 2007)*, pp. 639–642, July 2007.
- [9] J. W. Park, H. D. Oh, and M. J. Tahk, "Uav collision avoidance based on geometric approach," in *Proc. of 2008 SICE Annual Conference*, pp. 2122–2126, August 2008.
- [10] J. P. K. Kim and M. Tahk, "Uav collision avoidance using probabilistic method in 3-d," in *Proc. of 2007 International Conference on Control, Automation and Systems*, pp. 826–829, August 2007.
- [11] M. Shanmugavel, A. Tsourdos, and B. A. White, "Collision avoidance and path planning of multiple uavs using flyable paths in 3d," in *2010 15th International Conference on Methods and Models in Automation and Robotics*, Aug 2010, pp. 218–222.
- [12] Z. Chao, L. Ming, Z. Shaolei, and Z. Wenguang, "Collision-free uav formation flight control based on nonlinear mpc," in *2011 International Conference on Electronics, Communications and Control (ICECC)*, Sept 2011, pp. 1951–1956.
- [13] J. G. Manathara and D. Ghose, "Reactive collision avoidance of multiple realistic uavs," *Aircraft Engineering and Aerospace Technology*, vol. 83, no. 6, pp. 388–396, 2011. [Online]. Available: <https://doi.org/10.1108/00022661111173261>
- [14] M. C. P. Santos, C. D. Rosales, M. Sarcinelli-Filho, and R. Carelli, "A novel null-space-based uav trajectory tracking controller with collision avoidance," *IEEE/ASME Transactions on Mechatronics*, vol. 22, no. 6, pp. 2543–2553, Dec 2017.
- [15] L. A. Tony, D. Ghose, and A. Chakravarthy, "Avoidance maps: A new concept in uav collision avoidance," in *2017 International Conference on Unmanned Aircraft Systems (ICUAS)*, June 2017, pp. 1483–1492.
- [16] R. He, R. Wei, and Q. Zhang, "Uav autonomous collision avoidance approach," *Automatika*, vol. 58, no. 2, pp. 195–204, 2017. [Online]. Available: <https://www.tandfonline.com/doi/abs/10.1080/00051144.2017.1388646>
- [17] Y. Lin and S. Saripalli, "Sampling-based path planning for uav collision avoidance," *IEEE Transactions on Intelligent Transportation Systems*, vol. 18, no. 11, pp. 3179–3192, Nov 2017.
- [18] S. Ouahouah, J. Prados-Garzon, T. Taleb, and C. Benzaid, "Energy and delay aware physical collision avoidance in unmanned aerial vehicles," in *2018 IEEE Global Communications Conference (GLOBECOM)*, 2018, pp. 1–7.
- [19] S. Y. Choi and D. Cha, "Unmanned aerial vehicles using machine learning for autonomous flight; state-of-the-art," *Advanced Robotics*, vol. 33, no. 6, pp. 265–277, 2019.
- [20] 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*, pp. 1–12, 2019.
- [21] S.-Y. Shin, Y.-W. Kang, and Y.-G. Kim, "Obstacle avoidance drone by deep reinforcement learning and its racing with human pilot," *Applied Sciences*, vol. 9, no. 24, p. 5571, Dec 2019. [Online]. Available: <http://dx.doi.org/10.3390/app9245571>
- [22] P. Fraga-Lamas, L. Ramos, V. Mondéjar-Guerra, and T. M. Fernández-Caramés, "A review on iot deep learning uav systems for autonomous obstacle detection and collision avoidance," *Remote Sensing*, vol. 11, no. 18, p. 2144, Sep 2019. [Online]. Available: <http://dx.doi.org/10.3390/rs11182144>
- [23] D. Wang, T. Fan, T. Han, and J. Pan, "A two-stage reinforcement learning approach for multi-uav collision avoidance under imperfect sensing," *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 3098–3105, 2020.
- [24] K. Wan, X. Gao, Z. Hu, and G. Wu, "Robust motion control for uav in dynamic uncertain environments using deep reinforcement learning," *Remote Sensing*, vol. 12, no. 4, p. 640, Feb 2020. [Online]. Available: <http://dx.doi.org/10.3390/rs12040640>
- [25] B. M. Albaker and N. A. Rahim, "A survey of collision avoidance approaches for unmanned aerial vehicles," in *2009 International Conference for Technical Postgraduates (TECHPOS)*, Dec 2009, pp. 1–7.
- [26] L. Xiao, W. Zhuang, S. Zhou, and C. Chen, *UAV Relay in VANETs Against Smart Jamming with Reinforcement Learning*. Cham: Springer International Publishing, 2019, pp. 105–129.
- [27] L. Xiao, X. Lu, D. Xu, Y. Tang, L. Wang, and W. Zhuang, "UAV relay in VANETs against smart jamming with reinforcement learning," *IEEE Transactions on Vehicular Technology*, vol. 67, no. 5, pp. 4087–4097, 2018.

- [28] O. Bekkouche, T. Taleb, and M. Bagaa, "Uavs traffic control based on multi-access edge computing," in *IEEE Global Communications Conference, GLOBECOM 2018, Abu Dhabi, United Arab Emirates, December 9-13, 2018*. IEEE, 2018, pp. 1–6.
- [29] M. G. Ting-Hua Yi, Hong-Nan Li, "Experimental assessment of high-rate GPS receivers for deformation monitoring of bridge," *Elsivier*, vol. 46, no. 1, pp. 420 – 432, Aug. 2012.
- [30] O. Bekkouche, T. Taleb, and M. Bagaa, "Uavs traffic control based on multi-access edge computing," in *2018 IEEE Global Communications Conference (GLOBECOM)*, 2018, pp. 1–6.
- [31] S. V. Amarasinghe, H. S. Hewawasam, W. B. D. K. Fernando, J. V. Wijayakulasooriya, G. M. R. I. Godaliyadda, and M. P. B. Ekanayake, "Vision based obstacle detection and map generation for reconnaissance," in *2014 9th International Conference on Industrial and Information Systems (ICIIS)*, 2014, pp. 1–6.
- [32] D. Kim and H. Ryu, "Obstacle recognition system using ultrasonic sensor and duplex radio-frequency camera for the visually impaired person," in *13th International Conference on Advanced Communication Technology (ICACT2011)*, 2011, pp. 326–329.
- [33] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*, 2nd ed. The MIT Press, 2018. [Online]. Available: <http://incompleteideas.net/book/the-book-2nd.html>
- [34] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, "Playing atari with deep reinforcement learning," 2013, cite arxiv:1312.5602Comment: NIPS Deep Learning Workshop 2013. [Online]. Available: <http://arxiv.org/abs/1312.5602>
- [35] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, Feb. 2015. [Online]. Available: <http://dx.doi.org/10.1038/nature14236>
- [36] Gym: [https://gym.openai.com/envs/classic\\_control](https://gym.openai.com/envs/classic_control).



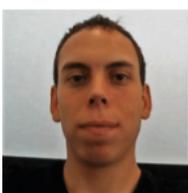
**Sihem Ouahouah** is currently pursuing her doctoral studies at Aalto University, Espoo, Finland. She received her engineer's from the University of Science and Technology Houari Boumediene (USTHB), Algeria and master's degrees in computer science from Ecole Nationale Supérieure d'Informatique (ESI), Algeria. Her research interests include unmanned aerial vehicles, Machine Learning, and the Internet of Things.



**Prof. Tarik Taleb** is currently Professor at the School of Electrical Engineering, Aalto University, Finland. He is the founder and director of the MOSA!C Lab ([www.mosaic-lab.org](http://www.mosaic-lab.org)). He is also working as part time professor at the Center of Wireless Communications, University of Oulu. Prior to his current academic position, he was working as Senior Researcher and 3GPP Standards Expert at NEC Europe Ltd, Heidelberg, Germany. He was then leading the NEC Europe Labs Team working on RD projects on carrier cloud platforms, an important vision of 5G systems. Before joining NEC and till Mar. 2009, he worked as assistant professor at the Graduate School of Information Sciences, Tohoku University, Japan, in a lab fully funded by KDDI. From Oct. 2005 till Mar. 2006, he worked as research fellow at the Intelligent Cosmos Research Institute, Sendai, Japan. He received his B. E degree in Information Engineering with distinction, M.Sc. and Ph.D. degrees in Information Sciences from Tohoku Univ., in 2001, 2003, and 2005, respectively. Prof. Taleb's research interests lie in the field of architectural enhancements to mobile core networks (particularly 3GPP's), network softwarization slicing, mobile cloud networking, network function virtualization, software defined networking, mobile multimedia streaming, inter-vehicular communications, and social media networking. Prof. Taleb has been also directly engaged in the development and standardization of the Evolved Packet System as a member of 3GPP's System Architecture working group. Prof. Taleb is a member of the IEEE Communications Society Standardization Program Development Board. Prof. Taleb is/was on the editorial board of the IEEE Transactions on Wireless Communications, IEEE Wireless Communications Magazine, IEEE Journal on Internet of Things, IEEE Transactions on Vehicular Technology, IEEE Communications Surveys Tutorials, and a number of Wiley journals. Prof. Taleb is the recipient of the 2017 IEEE ComSoc Communications Software Technical Achievement Award (Dec. 2017) for his outstanding contributions to network softwarization. He is also the (co-) recipient of the 2017 IEEE Communications Society Fred W. Ellersick Prize (May 2017), the 2009 IEEE ComSoc Asia-Pacific Best Young Researcher award (Jun. 2009), the 2008 TELECOM System Technology Award from the Telecommunications Advancement Foundation (Mar. 2008), the 2007 Funai Foundation Science Promotion Award (Apr. 2007), the 2006 IEEE Computer Society Japan Chapter Young Author Award (Dec. 2006), the Niwa Yasujiro Memorial Award (Feb. 2005), and the Young Researcher's Encouragement Award from the Japan chapter of the IEEE Vehicular Technology Society (VTS) (Oct. 2003). Some of Prof. Taleb's research work have been also awarded best paper awards at prestigious IEEE-flagged conferences.



**Dr. Miloud Bagaa** is a IEEE senior member and he received his Engineer's, Master's, and Ph.D. degrees from the University of Science and Technology Houari Boumediene (USTHB), Algiers, Algeria, in 2005, 2008, and 2014, respectively. From 2009 to 2015, he was a researcher with the Research Center on Scientific and Technical Information (CERIST), Algiers. From 2015 to 2016, he was granted a postdoctoral fellowship from the European Research Consortium for Informatics and Mathematics, and worked at the Norwegian University of Science and Technology, Trondheim, Norway. He was a postdoc researcher at Aalto University from 2016 to 2019, then a senior researcher from 2019 to October 2020. Currently, he is a senior cloud specialist at IT Center For Science LTD. His research interests include machine learning, optimization, networking modeling and network slicing.



**Dr. Jonathan Prados-Garzon** received his B.Sc., M.Sc., and Ph.D. degrees from the University of Granada (UGR), Granada, Spain, in 2011, 2012, and 2018, respectively. From 2018 to 2020, he worked as a postdoc researcher at MOSA!C Lab, led by Prof. Tarik Taleb, and the Department of Communications and Networking of Aalto University (Finland). Currently, he is a postdoc researcher at WiMuNet Lab, headed by Prof. Juan Manuel Lopez Soler, and the Department of Signal Theory, Telematics and Communications of the University

of Granada (Spain). His research interests include Mobile Broadband Networks, Network Softwarization, Deterministic Networking, and Network Performance Modeling and Optimization.