# Autonomous Blimp Control using Deep Reinforcement Learning

Yu Tang Liu[1,2], Eric Price[1,2], Pascal Goldschmid[2,1], Michael J. Black[1], Aamir Ahmad[2,1]

*Abstract*—Aerial robot solutions are becoming ubiquitous for an increasing number of tasks. Among the various types of aerial robots, blimps are very well suited to perform long-duration tasks while being energy efficient, relatively silent and safe. To address the blimp navigation and control task, in our recent work [1] we have developed a software-in-the-loop simulation and a PID-based controller for large blimps in the presence of wind disturbance. However, blimps have a deformable structure and their dynamics are inherently non-linear and time-delayed, often resulting in large trajectory tracking errors. Moreover, the buoyancy of a blimp is constantly changing due to changes in the ambient temperature and pressure. In the present paper, we explore a deep reinforcement learning (DRL) approach to address these issues. We train only in simulation, while keeping conditions as close as possible to the real-world scenario. We derive a compact state representation to reduce the training time and a discrete action space to enforce control smoothness. Our initial results in simulation show a significant potential of DRL in solving the blimp control task and robustness against moderate wind and parameter uncertainty. Extensive experiments are presented to study the robustness of our approach. We also openly provide the source code of our approach[1].

*Index Terms*—nonlinear control systems; reinforcement learning; blimp navigation; blimp control; aerial robotics; aerial vehicles

## I. INTRODUCTION

Autonomous unmanned aerial vehicles (UAVs) are becoming increasingly popular for various tasks, such as search and rescue, payload (medicine, food) delivery in difficult-to-reach areas, aerial cinematography and wildlife monitoring [2]–[8]. Current solutions rely on quadcopters and fixed-wings. Although quadcopters can hover in a fixed position, they are not able to accomplish long-term missions due to their short battery life. The situation is opposite for the fixed-wings, which have to move constantly to stay airborne. Therefore, for tasks involving long flight times, more payload and hovering over a small region, blimps provide an attractive solution.

A blimp is an airship without a rigid hull structure. Filled with helium, it becomes lighter than air and can hover for long periods without spending much energy. A blimp's weight is usually concentrated at its gondola, which creates a huge inertia to stabilize itself. From a control perspective, this makes the blimp an inherently stable plant [9] and allows it to recover from undesired states.

For blimp controller design, classic approaches usually rely on PID controllers [10]–[14] and nonlinear control [15]–[22]. PID controller suffers from nonlinearity, and nonlinear control

[1]Max Planck Institute for Intelligent System, Tübingen, Germany. `{firstname.lastname}@tuebingen.mpg.de`
[2]Institute for Flight Mechanics and Controls, The Faculty of Aerospace Engineering and Geodesy, University of Stuttgart, Stuttgart, Germany. `{firstname.lastname}@ifr.uni-stuttgart.de`
[1]https://github.com/robot-perception-group/AutonomousBlimpDRL

Fig. 1: Our autonomous blimp during a flight. Unlike common designs, our blimp has thrust vectoring increasing its agility. Inset: a gazebo model of our blimp.

methods require a dynamic model of the system which is often difficult to acquire. Deep reinforcement learning (DRL), on the other hand, is a new control framework, which has achieved success in a variety of applications that present similar challenges [23]–[28]. The model-free RL approach is particularly useful when it is difficult to estimate physical parameters such as buoyancy and aerodynamic effects. The learning ability allows the controller to potentially adapt to the dynamic change caused by the environment.

In [1] we developed a software-in-the-loop (SITL) simulation and a manually-tuned PID controller for the blimp control task. There, we demonstrated its ability to follow a waypoint sequence in the real world. To address the previously described issues of PID or nonlinear control approaches, in the current work we explore a DRL based method. Here, an RL agent, QRDQN [29], is deployed in the simulation environment to improve exploration and training stability. We design a variety of training task suites following the OpenAI-Gym framework [30], which is compatible with a variety of off-the-shelf RL platforms [31]–[33]. To achieve DRL training within a reasonable amount of time, we reduce problem difficulties by selecting task-specific features. We choose a discrete action space design to enforce continuity in the actuator commands. An action penalty is included in the reward function to regulate motor usage. To robustify the agent, we corrupt data by injecting noise to all observations and actions during training. We show that this integration address some of the common issues that may appear in a real world scenario.

## II. RELATED WORK

Control methods for blimps and airships, which have similar control schemes, have been well studied in the last decade [34]. Classic approaches usually rely on **PID controllers**. Popular applications include visual servo [10]–[12] and indoor miniature blimp [13], [14]. The PID control class, while simple and robust,

often suffers from plant nonlinearity. To overcome this weakness, advanced approaches have been developed using **nonlinear control theory**. Existing methods in this context include inverse optimal tracking control [15], dynamic inversion control [16], and the more investigated, backstepping control [17]–[19], robust control [20], [22], and model predictive control [21], [22]. However, these methods usually require a dynamic model which can be difficult to acquire. The buoyancy of the blimp is heavily dependent on the constantly-changing surrounding environment. When temperature or pressure changes, buoyancy also changes and a controller has to adapt on the fly. Unfortunately, this effect has not been addressed in any of the prior works so far.

On the other hand, recently there has been a surge of interest in applying RL to robotics [35]. The earliest attempts in the **classic RL** use Gaussian processes (GPs) for system identification [36] and policy learning [37], [38]. Despite sample efficiency, GPs are hard to scale up with problem dimensions and demand higher computational resources. As a result, they are able to achieve success only on low dimensional tasks, such as 1-D altitude control. **DRL**, on the other hand, leverages NNs for policy approximation and has achieved much success. This policy class can interpret rich representations and derive diverse behavior. For example, Nie et al. [39] train two DQN agents for rudder and elevator control of a blimp, respectively, and demonstrate a better performance than a PID controller. In the field of autonomous underwater vehicle (AUVs)[2], Carlucho et al. [40] use a Nessie-VII model with continuous action/observation space using DDPG.

The main challenge with DRL is the lack of sample efficiency. In order to scale up the DRL formulation with the problem dimension, a highly increased amount of environment interactions is needed by the agent. Other challenges include adapting a trained policy to real-world scenario [41], action smoothness [42], etc. Furthermore, issues such as partial observability, disturbances and noise could also lead to unexpected behavior. Issuing stability certificates to RL agents is also an ongoing research topic. In case of AUVs, the robustness issue is addressed in [43], by training a PPO agent in an adversarial fashion at the cost of conservative behavior.

To increase sample efficiency and finish training within a reasonable time budget, in this work we train the policy network with a value-based RL agent, QRDQN, which is a more stable variant in the DQN family. Action and observation space are injected with noise during training to robustify the agent. Lastly, to enforce actuator smoothness, we choose a discrete action space (Sec.III-B3).

## III. METHODOLOGY

### A. Preliminary

Blimp shape and architecture can vary a lot from one another. Therefore, we first describe our blimp, which has 8 actuators (but our approach is agnostic to different configurations). The two main motors (for thrust) are attached to a servo which

---

[2]Due to the lack of existent work, we include AUVs but only focus on those that have a fairly similar shape and task specification
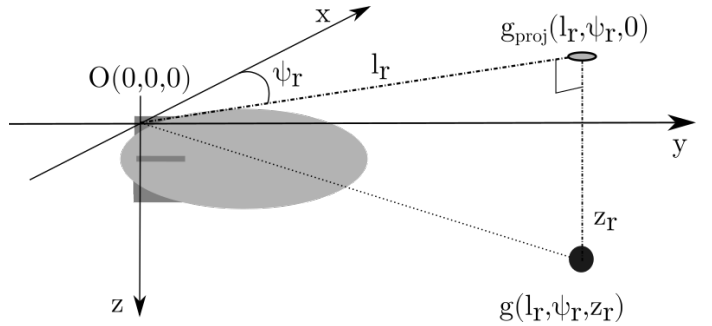


Fig. 2: The bodyframe in NED cylindrical coordinate system. $O$ is at the top fin of the blimp where GPU and IMU sensors are mounted. $g_{proj}$ is the projection on xy-plane of the waypoint $g$.

allows thrust vectoring. At the tail of the blimp, there are four fins controlling yaw and pitch angle and a tail motor, attached to the bottom fin, generating horizontal thrust allowing further yaw controllability. The state of the actuators can be denoted as

$$s_t^{act} = (m_{(0:2)}, s, f_{(0:3)})_t \in \mathbb{R}^8, \qquad (1)$$

where $m_{0:2}, s, f_{0:3}$ stand for motor, servo, and fin states, respectively. Our goal is to navigate this blimp to any given waypoint in the space by controlling these actuators.

### B. Formulation

*1) Control:* We formulate the problem as a path following task as seen in previous works [14], [39], [44]. In this setting, an imaginative path reference is generated based on waypoints for the controller to follow. Casting the path following task as a DRL problem, in this section we show how we reduce the state space size and maintain a reasonable training time. Since the blimp does not have a lateral movement control, we only need to consider longitudinal and altitude control. This allows us to easily decompose the problem into a planar navigation control task and an altitude control task. The objective of the planar navigation control is to control the blimp to arrive at any waypoint in the xy-plane whereas the altitude control is to reach the desired altitude of the waypoint.

Given the blimp position at $O = (0, 0, 0)$ and a target waypoint at $g = (l_r, \psi_r, z_r)$ in body frame cylindrical coordinates (Fig. 2), the control objective of the planar navigation control is the minimization of the relevant distance and yaw angle, or $\min_{a \in A}(|l_r|, |\psi_r|)$. The objective of the altitude control is to minimize the relevant altitude, or $\min_{a \in A} |z_r|$. The spatial information between the target and the blimp can be fully contained in $g = (l_r, \psi_r, z_r)$. Although it is possible to train the DRL method only using $g$, this minimal setting ignores the velocity and pitch state of the blimp, leading to instability in training and an uncontrolled behavior when reaching the waypoint.

We denote the velocity of the blimp as $V = (u, v, w)$, and attitude (roll, pitch, yaw) as $\Phi = (\phi, \theta, \psi)$. Assuming near zero lateral movement in the blimp (i.e. $v, \phi \simeq 0$), the velocity and pitch angle information can be encoded by velocity magnitude ($|V| = ||(u, v, w)||_2$) and the altitude velocity ($w = |V| sin\theta$),

alone. We augment our state with this velocity information and derive a compact representation for the overall state as

$$s_t^{blimp} = (l_r, \psi_r, z_r, |V|, w)_t, \quad (2)$$

which encodes all the spatial, velocity, and attitude information. As we do not directly use the pitch angle information, this representation is agnostic to sensor calibration error in pitch angle, which can be easily misaligned from the simulation.

*2) Markov Decision Process:* We consider the RL problem as an infinite horizon discrete time Markov Decision Process, $M$, defined by a tuple $(S, A, P, R, \gamma)$ [45]. At any time step $t \in \mathbb{R}^+$ and state $s_t \in \mathbb{R}^S$, an agent draws an action $a_t$ from a discrete action space $A = \{0, 1, ..., K\}$ given the policy distribution $a_t \sim \pi_\theta(\cdot|s_t)$ parameterized by $\theta$. The environment then samples the next state from an unknown transition distribution, i.e. $s_{t+1} \sim P(\cdot|s_t, a_t)$. A reward is received based on some reward function $r_t = R(s_t, a_t)$. Given the discount factor $\gamma \in [0, 1)$, the goal of the agent is to find the optimal policy parameter $\theta$ that comes with the highest cumulative discounted reward (3),

$$\pi^* = \arg\max_{\pi_\theta} \mathbb{E}_\pi [\sum_t^\infty \gamma^t r_t | a_t \sim \pi(\cdot|s_t), s_{t+1} \sim P(\cdot|s_t, a_t)] \quad (3)$$

*3) Observation and Action Space:* The full actuator state, $s_t^{act}$, is described in (1). Since we do not allow differential thrust, symmetric actuators are always in the same state. Thus, we only need to feedback one of them. The reduced state of actuators is therefore defined as $s_t^{act'} = (m_{(0,2)}, s, f_{(0,2)})_t \in \mathbb{R}^5$. The full state $s_t$ for the DRL formulation, as used in (3), is now obtained below as the concatenation of $s_t^{blimp}$ and $s_t^{act'}$.

$$s_t = (s_t^{act'}, s_t^{blimp}). \quad (4)$$

Note that all states are scaled to the range $[-1, 1]$ and zero-initialized. To prevent significant and sudden changes in the actuator command, we use discrete action space, denoted as $a_t \in A = \{0, 1, ..., K\}$. The action command is then mapped to the actuator command $\delta s_{a_t} \in \mathbb{R}^8$ following Table 1 and then summed up with the actuator state. This process is described below in (5).

$$a_t \sim \pi(\cdot|s_t)$$
$$\delta s_{a_t} = Table\ 1(a_t) \quad (5)$$
$$s_{a_t} \leftarrow s_{a_t} + \delta s_{a_t}$$

*4) Reward Function:* The control tasks of the UAVs usually involve navigation and hover. Navigation requires moving the robot in space by specifying a target position or following a sequence of targets, whereas hovering requires staying near the target position. These two tasks can be combined and trained with the same setup by using appropriate reward functions. When the goal is far away, we use a reward function for navigation only, otherwise a hover reward function. The reward function is defined by (6)

$$r_t = w_0 r_t^{success} + w_1 r_t^{track} + w_2 r_t^{act}, \quad (6)$$

where $w_{0:2} = (1, 0.95, 0.05)$ in this paper. The agent receives a success reward, $r_t^{success}$, if the task is completed. Tracking reward, $r_t^{track}$, indicates the tracking performance. Action reward, $r_t^{act}$, is defined to regularize actuator commands.

$$r_t^{success} = \begin{cases} 1 & \text{if } d(s_{blimp}, s_{target}) \leq \epsilon \\ 0 & \text{otherwise} \end{cases}, \quad (7)$$

$$r_t^{track} = \begin{cases} r_t^{hover} & \text{if } d(s_{blimp}, s_{target}) \leq \epsilon \\ r_t^{navigate} & \text{otherwise} \end{cases}, \quad (8)$$

where $d(s_{blimp}, s_{target})$ measures euclidean distance between the blimp and the target position. (7-8) indicate if this distance is short enough, the reward will be switched from navigation reward, $r_t^{navigate}$, to hover reward, $r_t^{hover}$, which does not take yaw component into account (9-10). Note that we could also use $|V|$ or $w$ in the reward function to address other tasks.

$$r_t^{navigate} = -i_0|z_r| - i_1|l_r| - i_2|\psi_r|, \quad (9)$$
$$r_t^{hover} = -j_0|z_r| - j_1|l_r|, \quad (10)$$
$$r_t^{act} = -k_0||m_0, m_1, m_2||_2, \quad (11)$$

where $i_{0:2} = (0.1, 0.7, 0.2), j_{0:1} = (0.3, 0.7), k_0 = 1$ in this paper.

| $A$ | Name | $\delta s_{a_t} = \delta[m_2, f_0, f_1, f_2, f_3, s, m_0, m_1]_t$ |
|---|---|---|
| 0 | IDLE | [0, 0, 0, 0, 0, 0, 0, 0] |
| 1 | THRUST+ | [0, 0, 0, 0, 0, 0, 0.01, 0.01] |
| 2 | THRUST- | [0, 0, 0, 0, 0, 0, -0.01, -0.01] |
| 3 | NOSE_UP | [0, 0.025, 0.025, 0, 0, 0, 0, 0] |
| 4 | NOSE_DOWN | [0, -0.025, -0.025, 0, 0, 0, 0, 0] |
| 5 | NOSE_LEFT | [0.025, 0, 0, 0.025, 0.025, 0, 0, 0] |
| 6 | NOSE_RIGHT | [-0.025, 0, 0, -0.025, -0.025, 0, 0, 0] |

TABLE 1: discrete action space $\delta s_{a_t}$: (%). The notation $f_{0:3}$ correspond to the angle of left/right/top/bottom fins. Note that in this work thrust vectoring is disabled ($s, \delta s = 0$).

*C. Training Setup*

In this section, we describe the important factors that contribute to stabilize the training and increase the robustness of the trained policy. During training, the target position is sampled randomly within the range of 200 cubic meters w.r.t. the blimp. Random sampling is important to increase sample diversity and avoid overfitting to a specific track. We reset the task only after 200 seconds so that there is sufficient time for the blimp to reach any target, and use the spare time to learn to stay within the target range. During training, to increase the robustness of the policy, observations and actions are injected with 5% of noise and clip to the range $[-1, 1]$. Lastly, while the simulation step time is $0.1s$, the policy step is $0.5s$. Since the blimp has a relatively long response time, we found it important to increase the step time for the action to take effect.

We train the policy network with the QRDQN agent. The value-based method is in general more sample efficient compared to gradient-based methods and can therefore accelerate
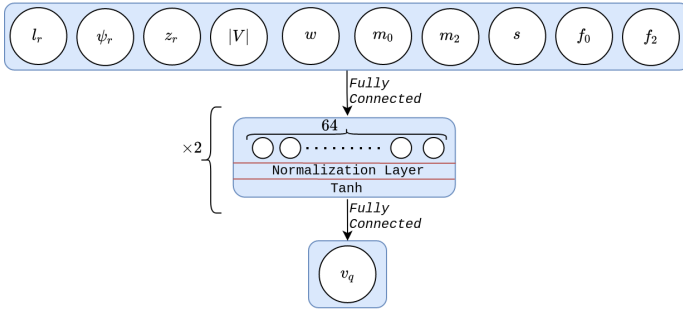
Fig. 3: The policy networks has the weights of 64 neurons by 2 layers. To prevent vanishing/exploding gradient, we add normalization layers to every linear layers. $v_q$ is a value array of size $K$ which evaluates each possible actions in $A$. Policy chooses action based on a greedy law, $a_t = argmax(v_q)$.

the training. QRDQN leverages a quantile network to estimate value function, which is important to stabilize training by alleviating chattering [29] effect and extreme value estimates. The architecture of our policy network is shown in Fig. 3. To reduce training time, we only apply less than 13 quantiles and sacrifice some estimation resolution.

## IV. Experiment

In the experiments, with the real world scenario in mind, we address the following questions. Is our DRL formulation with the compact representation able to solve the complex 3-D path-following task? In order to answer this question, we introduce the navigation and hovering tasks which are the building blocks for further complicated tasks. To evaluate the agent performance, a PID controller (from our previous work [1]) is considered as a benchmark, which is simple but well-known for its robustness. Finally, through various experiments (see sub-sec. IV-C) we evaluate if the RL agent is ready to be deployed in the real world. In other words, we evaluate the agent's robustness against unknown environmental changes.
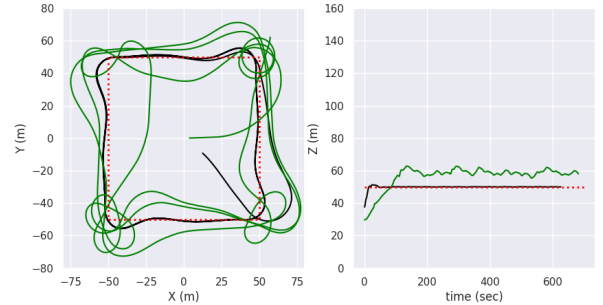
### A. Experiment Setup

We integrate our RL training environment in the ROS/Gazebo SITL simulation following the OpenAI-Gym framework. The QRDQN implementation is based on the StableBaseline3 [30]. The agent is trained for 7 days on a single computer (AMD Ryzen Threadripper 3960X, 24x 3.8GHz, NVIDIA GeForce RTX 2080 Ti, 11GB). Our simulation environment is designed based on our real robotic blimp (see Fig.1). The baseline PID controller is well-tuned to the simulation environment. Our previous work [1]) has shown that we could deploy it to the real world without further tuning, which implies a good quality of the simulation.
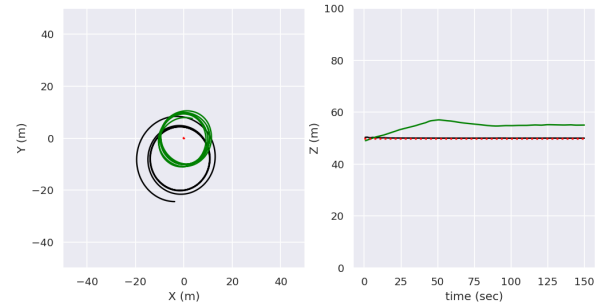
### B. Task Suite

To evaluate the performance of the agent, the navigation and hover tasks are introduced in the Sec.III-B. For convenience, we visualize the target waypoints in the world ENU frame.

*1) Navigation:* Four waypoints are created at an altitude of 50 m to form a square with sides of 100 m each. This has to be traversed in a counter-clockwise direction. A waypoint is registered when the blimp is within 15 m radius and then the

next waypoint is triggered. The early waypoint trigger allows less overshoot and achieves better performance. To make sure the comparison is fair, the track has to be performed 3 times to be marked as complete. The velocity for the PID controller is set to have a slow reference speed of $2m/s$ to prevent overshoot, while the agent is not subjected to any speed limit but $50\%$ maximum throttle. The results are shown in Fig. 4a. The PID controller has a stable performance during the whole task and remains a challenging baseline. On the other hand, although our trained RL policy can complete the navigation task successfully, it shows higher discrepancy from the reference path. It spends most of the time hovering above the waypoints and reduces the altitude until the next waypoint gets triggered.



(a) Navigation Task: Left – the planar trajectory of the blimp. Right – the altitude trajectory. Red: reference. Black: PID controller. Green: RL policy. The PID controller completes the task around 10% faster than the RL policy, which seems to favor an altitude 7 meters above the target altitude.



(b) Hover Task: Left: The planar trajectory of the blimp. Right: The altitude trajectory. Red: reference. Black: PID controller. Green: RL policy. The RL policy hovers with less radius around the target but it loiters 5 meters above the target altitude.

Fig. 4: Comparison of PID controller and RL policy in navigation and hover tasks.

*2) Hovering:* The hovering task requires the blimp to stay as close to the target as possible without spending excessive amount of energy. The target waypoint is positioned at $(0, 0, 50)$ in the world ENU frame. The blimp is spawned at the target position. The result is in Fig. 4b. The PID controller requires a larger radius compare to the RL policy. On the other hand, similar to the navigation task, the RL policy tends to hover 5 meters above the target altitude. Our initial reasoning for this behavior was lack of training. However, after continuing training the same policy, the results become worse. When the waypoint is an arbitrary point in space and far from the origin, the agent accelerates towards the target at first, then hovers close to it,
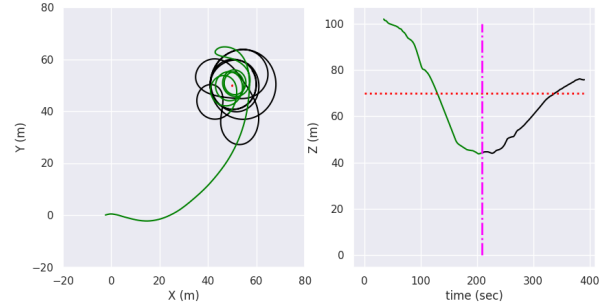
and finally makes a slow approach towards it. We argue that hovering close to the target altitude gives long-term advantages to the agent, analyzed as follows. First, the agent receives more action rewards ($r_t^{act}$) as it does not need to command anymore but only needs to wait until it slowly approaches the target altitude. Second, during this time, the distance is short enough to receive a good amount of hover reward ($r_t^{hover}$); and if the agent would rush to the target, it is most likely to overshoot and spend an excessive amount of energy to come back to the hover position, and subsequently overshoot again. Third, since the speed of the blimp becomes very small when approaching the target, the agent can easily stay longer within the target range and continuously receive abundant success rewards ($r_t^{success}$).

We are able to reproduce this behavior as shown in Fig. 5a. The blimp is spawned at $(0, 0, 100)$, which is above the target altitude $(50, 50, 70)$. We first observe that the blimp approaches the target, then stays close to it with a low speed. During this time, it still receives a good amount of tracking and success reward as shown in Fig. 5b. The blimp then continues to sink 25m below the target, after which the policy brings it back and raises it above the target altitude. This overall behavior required $\sim 200$s. This is followed by the hovering behavior as observed in Fig. 4b. Such a behavior comes from the fact that the total reward is dominated by the success reward and the loss of altitude does not result in significant punishment in tracking reward. Therefore, manipulating the reward function and increase altitude weight could help get rid of this behavior.
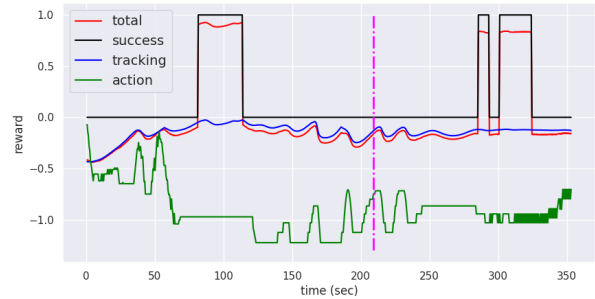
### C. Robustness Study

To show the robustness of the agent, we test our agent with i) a fixed wind field, ii) changes in the blimp buoyancy, and ii) changes in the weight distribution along the gondola. In Fig. 6, the agent is able to handle small wind disturbance at $2m/s$ but fails at $4m/s$. Under the wind condition, it takes a significantly large amount of time to finish the task. Notice that when the wind is at $2m/s$, the agent trajectory seems to be smoother. This is because wind slows down the agent and prevents overshooting the target. At the wind speed of $4m/s$, the agent tends to slow down when it approaches the target as in Fig. 7. This is the side effect from training navigation and hovering task together as the agent tries to slow down to stay within the range of success reward. It starts to slow down around 40 m to the target. As a result, although the agent has enough thrust power to overcome $4m/s$ wind, it gradually reduces both motors to zero speed and then is blown away by the wind. A naïve workaround is to toggle the target switch when the blimp is 40 m from the target. But a toggle with such a huge radius is not realistic. Notice that in Fig. 7, the motors always have smooth transition due to the discrete action design which only allows $1\%$ motor speed change every 0.5 seconds.

Another common scenario for the blimp in real world is the buoyancy change. Depending on the weather, the buoyancy of the blimp can change significantly. We test the performance of the agent with a decreasing amount of buoyancy w.r.t. the original state. Not to our surprise, the result in Tab. 2a shows that the decreasing amount of buoyancy does lead to worse



(a) What happens if the agent is spawned far from the target? Left: the planar view of the blimp approaching the target at (50,50,70). Right: the altitude trajectory of the blimp. Red: waypoint. The blimp significantly loses its altitude near the target position.



(b) Reward exploitation by the agent. Red: total reward. Black: Success reward. Blue: tracking reward. Green: action reward. The agent still receives a good amount of reward despite the altitude loss.

Fig. 5: Analysis of the hovering behavior. Magenta: the time agent responds to altitude loss. The total reward is calculated based on the (6). The tracking reward does not penalize altitude discrepancy enough and causes the strange behavior.
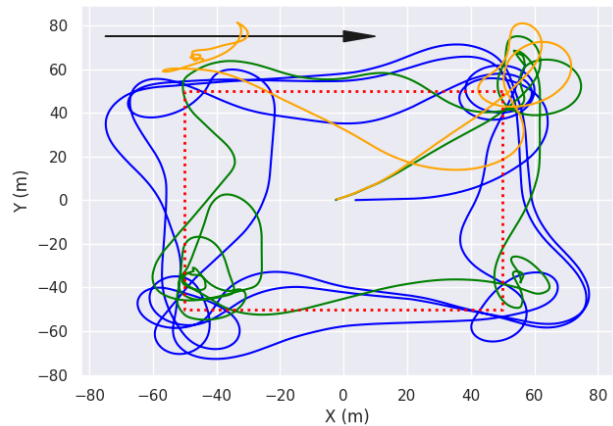


Fig. 6: Effect of wind. Blue: no wind. Green: $2m/s$ wind. Orange: $4m/s$ wind. The dark arrow indicates the direction of the wind field. This experiment was conducted for 3 laps for the 'no wind' case but for only 1 lap for the other two cases. When the wind speed is $4m/s$ the agent is not able to complete the task.

performance. With $95\%$ buoyancy, policy performance suffers significantly and takes much longer to finish the task. Lower than that the agent is not able to control the blimp at all. The effect of weight distribution (also commonly affected in real
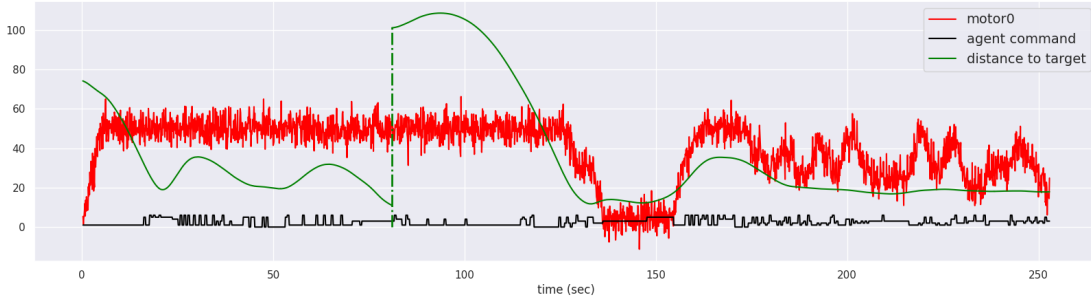
Fig. 7: Behavior of the DRL policy when the blimp is subject to $4m/s$ wind: Motor0 is the motor output ($\%$). Distance to target is calculated by euclidean distance ($m$). The left side of the green dotted line is the first waypoint, and right side is the second. The second waypoint is harder to reach since the blimp has to move against the wind. Whenever it is close to the target, it starts to transition into hover mode. This causes the blimp to be blown away by the wind, from which it is never able to recover.

| Buoyancy | Avg. Time (sec) |
|----------|-----------------|
| 100% | 238 |
| 95% | 545 |
| 90% | NA |
| 85% | NA |
| 80% | NA |

| Added mass (g) | Avg. Time (sec) |
|----------------|-----------------|
| 0 | 238 |
| -100 | 237 |
| -250 | NA |
| 100 | 328 |
| 250 | NA |

(a) Effect of buoyancy on average time to complete a square.

(b) Effect of trim weight on average time to complete a square

TABLE 2: Buoyancy and weights change cause significant impact on the agent

world) also can not be ignored as it could introduce unnecessary vibrations if not balanced. To this end we perform another experiment, where we add and remove ballast to the front end of the blimp to break this balance. Results in Tab. 2b suggest that 100g of mass change does not affect the performance, but larger than that would impair the policy. These two experiments have shown that the RL agent is currently sensitive to the environment changes. We expect to improve the performance of the agent by increasing the penalty for the altitude loss.

## V. DISCUSSION

In this work, we integrated the ROS/Gazebo SITL blimp simulation together with the RL training environment. We have derived a compact representation of the state space and action space which allows less training time and guarantee the actuator continuity. We have shown that such a setting is able to successfully complete the task. The trained policy network has a certain degree of robustness against wind and parameter uncertainty.

On the other hand, we have observed and analysed how the agent exploits the reward function. The altitude loss is unacceptably large for this agent to be deployed to the real world. Increasing the altitude reward weights and punishing the altitude loss could potentially address this issue. However, further experiments are needed to verify this hypothesis. In this work, the reverse thrust and thrust vectoring were not enabled. Given a more diverse action space, the agent is more likely to gain more rewards by staying closer to the target. Another problem is that when training navigation and hover in the same time the agent learns the conservative behavior when approaching a waypoint. This could be potentially eliminated

by including disturbance to the training and making it harder to exploit the weakness of this approach. A more promising solution would be multi-task learning which trains navigation and hovering task independently.

There are many other open issues not been addressed in this work so far. In real world experiments, not presented in this paper, we have encountered several difficulties even when flying with a PID controller. For example, in this work we assumed the lateral movement can be neglected and longitude velocity is always positive. To our observation in real world, this is a dangerous assumption as it does not hold in the presence of moderate to strong wind. When the wind speed is larger than the vehicle's, the speed can become negative and lateral movement can be created if the wind is blowing from the side. This can be dangerous and cause undesired behavior for the policy network.

Finally, blimp control has not received enough attention and still remains an underdeveloped field. RL-based methods do not provide any stability guarantee but provide the potential to learn continuously from data and improve its own performance. Conversely, the nonlinear controllers are robust against parameter uncertainty and disturbance at the expense of control performance. How to leverage these two approaches is the key to the success of future blimp control methods. Since blimp dynamic is heavily dependent on the environment, it serves as a perfect robotic platform to study adaptive learning control. Secondly, the modern DRL algorithms are still not sample efficient enough. Our next step is to leverage parallel training to accelerate gathering experience. This also allows us to increase the diversity in the training and offer the potential to leverage multi-tasking learning as mentioned in [46]. Lastly, for the agent to counter partial observations such as wind disturbances, it is important to include past experiences in the decision-making process. For example, a recurrent network architecture might be a possible solution.

## REFERENCES

[1] E. Price, Y. T. Liu, M. J. Black, and A. Ahmad, "Simulation and control of deformable autonomous airships in turbulent wind," *arXiv preprint arXiv:2012.15684*, 2020.

[2] L. F. Gonzalez, G. A. Montes, E. Puig, S. Johnson, K. Mengersen, and K. J. Gaston, "Unmanned aerial vehicles (uavs) and artificial intelligence revolutionizing wildlife monitoring and conservation," *Sensors*, vol. 16, no. 1, p. 97, 2016.

[3] L. Wang, F. Chen, and H. Yin, "Detecting and tracking vehicles in traffic by unmanned aerial vehicles," *Automation in construction*, vol. 72, pp. 294–308, 2016.

[4] K. Cesare, R. Skeele, S.-H. Yoo, Y. Zhang, and G. Hollinger, "Multi-uav exploration with limited communication and battery," in *2015 IEEE International Conference on Robotics and Automation (ICRA)*, 2015, pp. 2230–2235.

[5] S. Rathinam, P. Almeida, Z. Kim, S. Jackson, A. Tinka, W. Grossman, and R. Sengupta, "Autonomous searching and tracking of a river using an uav," in *2007 American Control Conference*, 2007, pp. 359–364.

[6] A. Ahmad, E. Price, R. Tallamraju, N. Saini, G. Lawless, R. Ludwig, I. Martinovic, H. H. Bülthoff, and M. J. Black, "Aircap – aerial outdoor motion capture," Nov. 2019.

[7] I. Mademlis, V. Mygdalis, N. Nikolaidis, M. Montagnuolo, F. Negro, A. Messina, and I. Pitas, "High-level multiple-uav cinematography tools for covering outdoor events," *IEEE Transactions on Broadcasting*, vol. 65, no. 3, pp. 627–635, 2019.

[8] N. Saini, E. Price, R. Tallamraju, R. Ludwig, R. Enficiaud, A. Ahmad, and M. Black, "Markerless outdoor human motion capture using multiple autonomous micro aerial vehicles," 08 2019.

[9] Y. Li and M. Nahon, "Modeling and simulation of airship dynamics," *Journal of Guidance, Control, and Dynamics*, vol. 30, no. 6, pp. 1691–1700, 2007.

[10] J. Azinheira, P. Rives, J. Carvalho, G. Silveira, E. de Paiva, and S. Bueno, "Visual servo control for the hovering of all outdoor robotic airship," in *Proceedings 2002 IEEE International Conference on Robotics and Automation (Cat. No.02CH37292)*, vol. 3, 2002, pp. 2787–2792 vol.3.

[11] Hong Zhang and J. P. Ostrowski, "Visual servoing with dynamics: control of an unmanned blimp," in *Proceedings 1999 IEEE International Conference on Robotics and Automation (Cat. No.99CH36288C)*, vol. 1, May 1999, pp. 618–623 vol.1.

[12] S. van der Zwaan, A. Bernardino, and J. Santos-Victor, "Vision based station keeping and docking for an aerial blimp," in *Proceedings. 2000 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS 2000) (Cat. No.00CH37113)*, vol. 1, Oct 2000, pp. 614–619 vol.1.

[13] T. Takaya, H. Kawamura, Y. Minagawa, M. Yamamoto, and A. Ohuchi, "Pid landing orbit motion controller for an indoor blimp robot," *Artificial Life and Robotics*, vol. 10, no. 2, pp. 177–184, 2006.

[14] Jinjun Rao, Zhenbang Gong, Jun Luo, and Shaorong Xie, "A flight control and navigation system of a small size unmanned airship," in *IEEE International Conference Mechatronics and Automation, 2005*, vol. 3, July 2005, pp. 1491–1496 Vol. 3.

[15] T. Fukao, T. Kanzawa, and K. Osuka, "Inverse optimal tracking control of an aerial blimp robot," in *Proceedings of the Fifth International Workshop on Robot Motion and Control, 2005. RoMoCo '05.*, 2005, pp. 193–198.

[16] A. Moutinho and J. Azinheira, "Stability and robustness analysis of the aurora airship control system using dynamic inversion," in *Proceedings of the 2005 IEEE International Conference on Robotics and Automation*, 2005, pp. 2265–2270.

[17] J. R. Azinheira and A. Moutinho, "Hover control of an uav with backstepping design including input saturations," *IEEE Transactions on Control Systems Technology*, vol. 16, no. 3, pp. 517–526, 2008.

[18] L. Beji, A. Abichou, and Y. Bestaoui, "Stabilization of a nonlinear underactuated autonomous airship-a combined averaging and backstepping approach," in *Proceedings of the Third International Workshop on Robot Motion and Control, 2002. RoMoCo '02.*, 2002, pp. 223–229.

[19] S. Q. Liu, Y. J. Sang, and J. F. Whidborne, "Adaptive sliding-mode-backstepping trajectory tracking control of underactuated airships," *Aerospace Science and Technology*, vol. 97, p. 105610, 2020. [Online]. Available: http://www.sciencedirect.com/science/article/pii/S1270963819310181

[20] L. Cheng, Z. Zuo, J. Song, and X. Liang, "Robust three-dimensional path-following control for an under-actuated stratospheric airship," *Advances in Space Research*, 09 2018.

[21] Hiroaki Fukushima, Ryosuke Saito, Fumitoshi Matsuno, Yasushi Hada, Kuniaki Kawabata, and Hajime Asama, "Model predictive control of an autonomous blimp with input and output constraints," in *2006 IEEE Conference on Computer Aided Control System Design, 2006 IEEE International Conference on Control Applications, 2006 IEEE International Symposium on Intelligent Control*, Oct 2006, pp. 2184–2189.

[22] S. Liu, Y. Sang, and H. Jin, "Robust model predictive control for stratospheric airships using lpv design," *Control Engineering Practice*, vol. 81, pp. 231 – 243, 2018. [Online]. Available: http://www.sciencedirect.com/science/article/pii/S0967066118305100

[23] H. J. Kim, M. I. Jordan, S. Sastry, and A. Y. Ng, "Autonomous helicopter flight via reinforcement learning," in *Advances in Neural Information Processing Systems 16*, S. Thrun, L. K. Saul, and B. Schölkopf, Eds. MIT Press, 2004, pp. 799–806. [Online]. Available: http://papers.nips.cc/paper/2455-autonomous-helicopter-flight-via-reinforcement-learning.pdf

[24] J. Hwangbo, I. Sa, R. Siegwart, and M. Hutter, "Control of a quadrotor with reinforcement learning," *IEEE Robotics and Automation Letters*, vol. 2, no. 4, p. 2096–2103, Oct 2017. [Online]. Available: http://dx.doi.org/10.1109/LRA.2017.2720851

[25] H. A. M. C. e. a. Silver, D., "Mastering the game of go with deep neural networks and tree search," *Nature529*, 2016.

[26] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. A. Riedmiller, "Playing atari with deep reinforcement learning," *ArXiv*, vol. abs/1312.5602, 2013.

[27] H. Zhu, A. Gupta, A. Rajeswaran, S. Levine, and V. Kumar, "Dexterous manipulation with deep reinforcement learning: Efficient, general, and low-cost," 2018.

[28] M. G. Bellemare, S. Candido, P. S. Castro, J. Gong, M. C. Machado, S. Moitra, S. S. Ponda, and Z. Wang, "Autonomous navigation of stratospheric balloons using reinforcement learning," *Nature*, vol. 588, no. 7836, pp. 77–82, 2020.

[29] B. Mavrin, H. Yao, L. Kong, K. Wu, and Y. Yu, "Distributional reinforcement learning for efficient exploration," in *International conference on machine learning*. PMLR, 2019, pp. 4424–4434.

[30] A. Raffin, A. Hill, M. Ernestus, A. Gleave, A. Kanervisto, and N. Dormann, "Stable baselines3," https://github.com/DLR-RM/stable-baselines3, 2019.

[31] E. Liang, R. Liaw, R. Nishihara, P. Moritz, R. Fox, K. Goldberg, J. E. Gonzalez, M. I. Jordan, and I. Stoica, "RLlib: Abstractions for distributed reinforcement learning," in *International Conference on Machine Learning (ICML)*, 2018.

[32] S. Guadarrama, A. Korattikara, O. Ramirez, P. Castro, E. Holly, S. Fishman, K. Wang, E. Gonina, N. Wu, E. Kokiopoulou, L. Sbaiz, J. Smith, G. Bartók, J. Berent, C. Harris, V. Vanhoucke, and E. Brevdo, "TF-Agents: A library for reinforcement learning in tensorflow," https://github.com/tensorflow/agents, 2018, [Online; accessed 25-June-2019]. [Online]. Available: https://github.com/tensorflow/agents

[33] M. Hoffman, B. Shahriari, J. Aslanides, G. Barth-Maron, F. Behbahani, T. Norman, A. Abdolmaleki, A. Cassirer, F. Yang, K. Baumli, S. Henderson, A. Novikov, S. G. Colmenarejo, S. Cabi, C. Gulcehre, T. L. Paine, A. Cowie, Z. Wang, B. Piot, and N. de Freitas, "Acme: A research framework for distributed reinforcement learning," *arXiv preprint arXiv:2006.00979*, 2020. [Online]. Available: https://arxiv.org/abs/2006.00979

[34] Y. Liu, Z. Pan, D. Stirling, and F. Naghdy, "Control of autonomous airship," in *2009 IEEE International Conference on Robotics and Biomimetics (ROBIO)*, 2009, pp. 2457–2462.

[35] B. Singh, R. Kumar, and V. P. Singh, "Reinforcement learning in robotic applications: a comprehensive survey," *Artificial Intelligence Review*, pp. 1–46, 2021.

[36] J. Ko, D. J. Klein, D. Fox, and D. Haehnel, "Gaussian processes and reinforcement learning for identification and control of an autonomous blimp," in *Proceedings 2007 IEEE International Conference on Robotics and Automation*, April 2007, pp. 742–747.

[37] A. Rottmann and W. Burgard, "Adaptive autonomous control using online value iteration with gaussian processes," in *2009 IEEE International Conference on Robotics and Automation*, May 2009, pp. 2106–2111.

[38] A. Rottmann, C. Plagemann, P. Hilgers, and W. Burgard, "Autonomous blimp control using model-free reinforcement learning in a continuous state and action space," in *2007 IEEE/RSJ International Conference on Intelligent Robots and Systems*, Oct 2007, pp. 1895–1900.

[39] C. Nie, Z. Zheng, and M. Zhu, "Three-dimensional path-following control of a robotic airship with reinforcement learning," *International Journal of Aerospace Engineering*, vol. 2019.

[40] I. Carlucho, M. De Paula, S. Wang, B. V. Menna, Y. R. Petillot, and G. G. Acosta, "Auv position tracking control using end-to-end deep reinforcement learning," in *OCEANS 2018 MTS/IEEE Charleston*, 2018, pp. 1–8.

[41] Y. Zhang, T. Liu, M. Long, and M. Jordan, "Bridging theory and algorithm for domain adaptation," in *International Conference on Machine Learning*. PMLR, 2019, pp. 7404–7413.

[42] S. Mysore, B. Mabsout, R. Mancuso, and K. Saenko, "Regularizing action policies for smooth control with reinforcement learning," 2021.

[43] J. Parras and S. Zazo, "Robust deep reinforcement learning for underwater navigation with unknown disturbances," in *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021, pp. 3440–3444.

[44] H. Saiki, T. Fukao, T. Urakubo, and T. Kohno, "Hovering control of outdoor blimp robots based on path following," in *2010 IEEE International Conference on Control Applications*, Sep. 2010, pp. 2124–2129.

[45] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*, 2018.

[46] L. Espeholt, H. Soyer, R. Munos, K. Simonyan, V. Mnih, T. Ward, Y. Doron, V. Firoiu, T. Harley, I. Dunning *et al.*, "Impala: Scalable distributed deep-rl with importance weighted actor-learner architectures," in *International Conference on Machine Learning*. PMLR, 2018, pp. 1407–1416.