

---

# HARFANG3D DOG-FIGHT SANDBOX: A REINFORCEMENT LEARNING RESEARCH PLATFORM FOR THE CUSTOMIZED CONTROL TASKS OF FIGHTER AIRCRAFTS

---

**Muhammed Murat Özbek**  
Istanbul Technical University  
ozbekm17@itu.edu.tr

**Süleyman Yıldırım**  
Koç University  
suleymanyildirim22@ku.edu.tr

**Muhammet Aksoy**  
Istanbul Technical University  
aksoym15@itu.edu.tr

**Eric Kernin**  
HARFANG 3D NWNC  
eric.kernin@harfang3d.com

**Emre Koyuncu**  
Istanbul Technical University  
emre.koyuncu@itu.edu.tr

## ABSTRACT

*Abstract - The advent of deep learning (DL) gave rise to significant breakthroughs in Reinforcement Learning (RL) research. Deep Reinforcement Learning (DRL) algorithms have reached super-human level skills when applied to vision-based control problems as such in Atari 2600 games where environment states were extracted from pixel information. Unfortunately, these environments are far from being applicable to highly dynamic and complex real-world tasks as in autonomous control of a fighter aircraft since these environments only involve 2D representation of a visual world. Here, we present a semi-realistic flight simulation environment Harfang3D Dog-Fight Sandbox for fighter aircrafts. It is aimed to be a flexible toolbox for the investigation of main challenges in aviation studies using Reinforcement Learning. The program provides easy access to flight dynamics model, environment states, and aerodynamics of the plane enabling user to customize any specific task in order to build intelligent decision making (control) systems via RL. The software also allows deployment of bot aircrafts and development of multi-agent tasks. This way, multiple groups of aircrafts can be configured to be competitive or cooperative agents to perform complicated tasks including “Dog Fight”. During the experiments, we carried out training for two different scenarios: navigating to a designated location and within visual range (WVR) combat, shortly “Dog Fight”. Using Deep Reinforcement Learning techniques for both scenarios, we were able to train competent agents that exhibit human-like behaviours. Based on this results, it is confirmed that Harfang3D Dog-Fight Sandbox can be utilized as a 3D realistic RL research platform.*

## 1 Introduction

The training of machine learning models to make a sequence of decisions is known as Reinforcement Learning. In an uncertain and potentially complex environment, the RL agent meets a game-like simulation and learns to achieve a specified goal. In order to accomplish a given task, the agent uses trial and error method. Consequently, it is given either rewards or penalties for the actions that it takes. As a result, starting with an agent that behaves completely randomly, the algorithm is capable of developing sophisticated tactics without any prior knowledge about the environment or the game. The ability to adapt stochastic, complex, and dynamics environments to perform any task with superhuman skills makes Reinforcement Learning a powerful machine learning tool in the area of robotics and control. Therefore, researchers have been using RL in various control tasks [1], [2], [3], [4]: manipulation [5], [6], [7], [8], locomotion [9], [10], navigation [11], [12],[13], [14], flight [15], [16], interaction [17], [18], motion planning [19], [20] and more.

However, training of RL models differs from traditional machine learning algorithms due to its game-like nature. Unlike supervised and unsupervised learning, there is no labeled data nor a static dataset in Reinforcement Learning. This means that dataset changes dynamically as the agent interacts with its environment because the training samples

---

comes from agent’s previous experiences. This poses a number of problems when it comes to training an RL agent, which are high sample complexity, training time and realistic environment. Training agents in actual robotic systems is not a feasible option because of replicability, financial, time, and safety concerns. Hence, simulation environments customized for specific tasks are traditionally used for training and experimentation in Reinforcement Learning research [5], [9], [21]. On top of that, distributional training methods have been utilized to increase the sampling diversity and decrease the training time [22],[23], [24],[25]. A downside of carrying out experiments in a simulated environment is transferring the model into the real world. Researchers proposed several methods to overcome this domain problem, which is brought about by the lack of fidelity and stochasticity in the simulated environments, including calibration [6], [26], utilizing real-world data samples [13], [30], [31], and domain randomization [16], [27], [28], [29]. But still, how realistic a simulation environment is an important factor to be able to solve this problem.

Grounded on the facts explained above, it can be concluded that having access to a realistic, customizable and flexible simulation software that is specialized in the desired field is crucial to experiment with RL. To address the need of such platform in the area of autonomous control of fighter aircrafts, we propose the RL experimentation software Harfang3D Dog-Fight Sandbox. Harfang3D Dog-Fight Sandbox simulates an air-to-air confrontation scene with continuous sensory measurements received from the environment (e.g., Euler angles, coordinates, velocity vectors, and acceleration). Through the proposed platform, one can experiment with different simulation parameters, customize various tasks for fighter aircrafts, build multi agent systems, and run multiple parallel sessions during training by making use of the distributed structure. Moreover, the software can also be interfaced with the OpenAI Gym library [21]. For experimentation purposes, we trained our model using the Twin Delayed Deep Deterministic Policy Gradient (TD3) algorithm [32] without resorting to any real-world data, or expert knowledge.

## 2 Harfang3D Dog-Fight Sandbox: Reinforcement Learning Research Platform

Harfang3D Dog-Fight Sandbox is an easy-to-adapt, cross-platform, multi-language, powerful and optimized solution to integrate with embedded systems, into existing environments and combining features to meet both the general public and military challenges of real-time 3D imaging: providing strategic autonomy, reactivity, flexibility, interoperability and durability. It is entirely written in C++, with high level language API in Python, Lua and Golang.



Figure 1: Some examples of the realistic renders of the fighter jets. F-16 jet is on the left side, and TF-X is on the right side.

The platform provides features that facilitate adaptable, flexible, and inclusive research environment. The main features of Harfang3D Dog-Fight Sandbox include different control modes, custom scenarios, custom physics models, distributed run, and renderless mode eliminating the need of using a graphical interface.

### 2.1 Runtime

#### 2.1.1 Synchronization Modes

Harfang3D Dog-Fight Sandbox is capable of running above 1000 FPS on a modern PC if needed. But, this number may vary depending on the simulation run mode. We can investigate the simulation run modes in two parts as synchronous and asynchronous modes. In synchronous mode, simulation and DL model have to run at the same frequency because simulation is not updated unless the model gives a command as they run consecutively. Therefore, synchronous mode limits the frame rate of the simulation to approximately 50 Hz. However, simulation can run much faster in asynchronous mode because in this mode, both processes run independently at their own pace without having to wait

for each other's response. For instance, if the simulation runs at 1000 Hz and the DL model runs at 50 Hz, the model is going to make 1 inference at every 20 simulation steps since they can run simultaneously. In most cases, real-life scenarios can be modeled more accurately in asynchronous mode. Still, to avoid temporal problems, synchronous mode is preferred for experimentation (OpenAI games and applications run in synchronous mode) due to possible discrepancies between simulation frame rate and algorithm speed. Since simulation and algorithm run consecutively at each step, any temporal constraint is automatically eliminated.

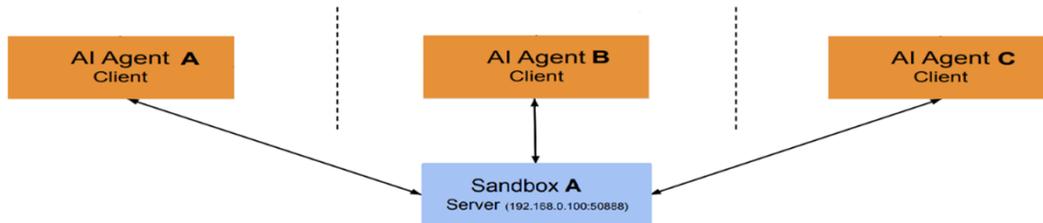
### 2.1.2 Distributed Run

One of the best properties of Harfang3D Dog-Fight Sandbox is that agents can independently run in different environments in parallel at the same time. Additionally, multiple agents can also run in one environment as cooperative or competitive groups. This will give us the opportunity to decrease training time and to boost exploration as well as to create highly complex scenarios and complicated tasks.

For example, this property may be exploited to speed up the training process by running one agent in multiple sandboxes. This way, one agent can take advantage of numerous training samples coming from variety of sources. As this may multiply the computational cost, it will also immensely accelerate the training while enforcing exploration.

Another example is to train multiple agents in multiple environments in parallel. With this setup, different hyperparameter configurations can be experimented at the same time in parallel when training your agent. After that, the best agent can be trained optimally in multiple environments as discussed in the first example.

#### Variant 1:



#### Variant 2:

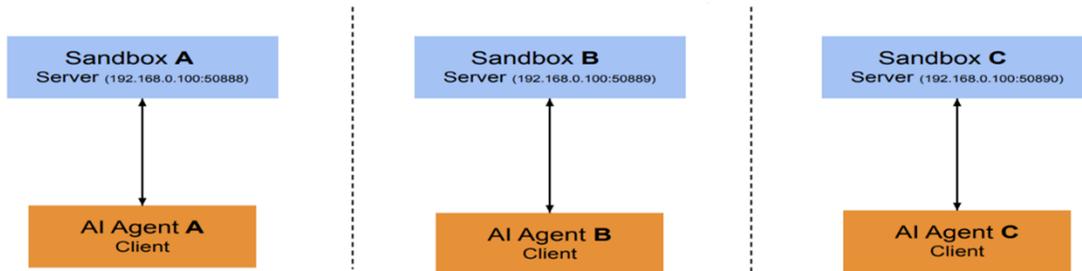


Figure 2: Distributional architecture of Harfang3D Dog-Fight Sandbox and possible parallel running scenarios.

## 2.2 Simulation Components

Harfang has mainly two different trainable simulation components. These are aircrafts and missiles. There are five models of fighter jets including world-renowned F-16 and recently being developed TFX. Missiles, on the other hand, consist of two categories (SAM and AAM) and six models in total. Each of these units have unique specialties due to distinct characteristics.

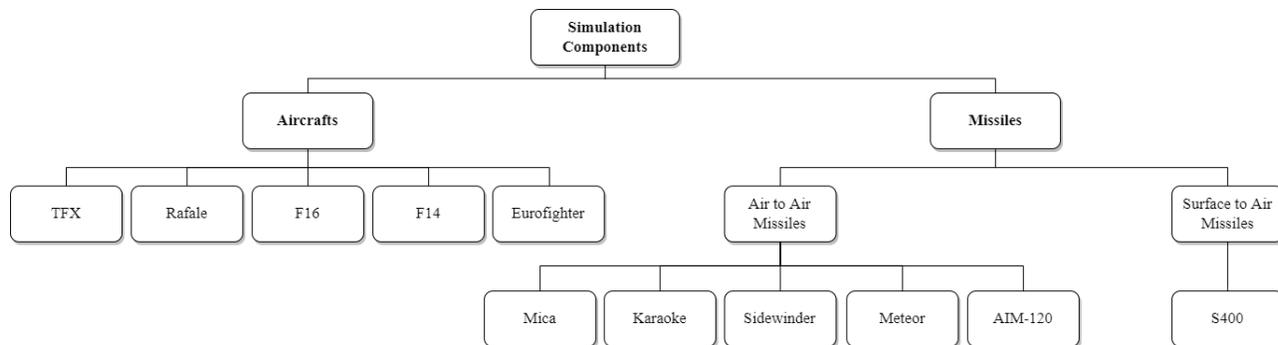


Figure 3: A general diagram for the simulation components and their sub categories.

### 2.2.1 Aircrafts

There currently are five different fighter aircrafts in the simulation such as Turkish Fighter X (TF-X) by Turkish Aerospace Industries (TAI), two variants of F-14 by Northrop Grumman (USA), F-16 by General Dynamics (USA), and Rafale by Dassault (France). Model parameters of the aircrafts are listed in Table 1 and they can be altered as desired. Due to the fact that each aircraft has distinct flight dynamics model, each of them has advantageous and disadvantageous aspects accordingly. Therefore, during training, AI agents experience different aircraft with varying dynamics and ammunition configurations. Consequently, agents are simulated against various opponents and get to learn complex tactics. There are two types of weapons: missiles and machine guns. Different missile configurations of each aircraft can also be found in Table 1. Prior to the launching of a missile, the target must be locked. The ammunition of the machine gun is not limited.

	TFX	Rafale	F16	F14	Eurofighter
Thrust Force	20	15	15	10	13
Post Combustion Force	20	7.5	15	5	9
Angular Frictions	0.000175	0.000165	0.000175	0.000175	0.00019
Speed Ceiling Force	2500	2200	1750	1750	2500
Max Safe Altitude	25000	15240	15700	15700	16800
Max Altitude	30000	25240	25700	25700	26800
Missile Number	4	6	12	4	6
Missile Config	4xAIM-120	2xMica 4xMeteor	8xKaraoke 2xAIM-120 2xCFT	4xSidewinder	2x Meteor 4xMica

Table 1: This table demonstrates the different parameters of five different fighter aircrafts and their configurations.

### 2.2.2 Missiles

In this simulation, missiles can be gathered under two categories. These are air-to-air missiles (AAMs) and surface-to-air missiles (SAMs). Mica, Karaoke, Sidewinder, Meteor, and AIM-120 are currently available air-to-air missiles in Harfang. For the surface-to-air missile category, S-400 is a popular option that is provided in this simulation. Properties of each missile differ from each other to enrich customizability.

**Air to Air Missiles** AAMs can be differentiated according to four main properties: thrust, endurance, damage, angular friction. Mica has the most thrust power, which makes it the fastest missile among four others listed in Table 2, but it also has the least endurance. Karaoke gives the largest amount of damage even-though it is the slowest missile, while Meteor has the longest endurance. Table 2 demonstrates the properties of five different missile variants. Please note that these values can be easily varied.

	Mica	Karaoke	Sidewinder	Meteor	AIM-120	S-400
Thrust Force	150	70	100	80	120	200
Endurance	15	35	20	40	20	210
Damage	20-30	50-70	30-40	40-60	25-35	100
Angular Frictions	0.00014	0.00005	0.00008	0.00005	0.00008	0.000025

Table 2: This table demonstrates the general properties of each missile. Damage values are given as an interval rather than a strict value because the magnitude of the damage is chosen randomly from the given value range. S-400 missile is considered as the most powerful missile in this simulation; thus, its damage is directly set to 100 and it can destroy an airplane with a single hit.

**Surface to Air Missiles** One of the most popular surface-to-air missile system is S-400. This system has the longest endurance and the most thrust power when compared to its counterparts. Therefore, S-400 is used in the Harfang 3D Dog Fight sandbox. Table 2 shows physical model parameters of S-400.

### 2.3 Scenario Options

Environment customizability is a key property in RL research. For this purpose, Harfang3D Dog-Fight Sandbox possesses an highly adaptable architecture which enables the creation of diverse tasks and scenarios from many perspectives. This includes being able to deploy different number of players, appoint other bot, AI, or human players as enemy or ally, program custom physical models, modify physical aspects (e.g., wing area) or other features of an aircraft. When combined with its diverse, large, and continuous observation space, it makes Harfang3D Dog-Fight Sandbox a powerful tool capable of creating numerous distinct scenarios. The concept is straightforward: In a three-dimensional space, aircrafts will fight each other. The starting point is always an aircraft carrier. Different protagonists can be piloted by a human, by an autopilot or by an AI that has been developed outside the Sandbox.

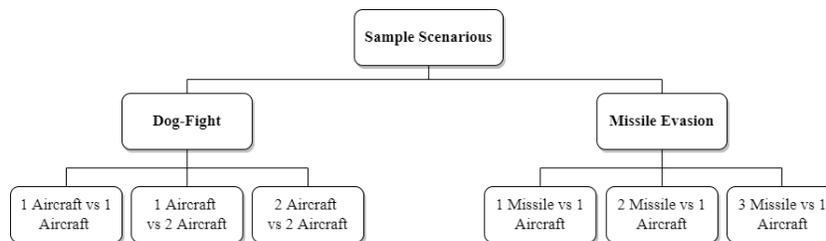


Figure 4: This diagram illustrates sample customized scenarios that can be implemented in Harfang.

We have established two main use cases for Harfang: missile evasion and dog fight. Dog fight can be performed in 1 versus 1, 1 versus 2, and 2 versus 2 format. Multiple agents from the same team can cooperate while communicating with each other to defeat the opponent team. On the other hand, missile evasion can be performed in 1 versus 1, 1 versus 2, and 1 versus 3 format. For now, simulation only supports a maximum of three missiles at once for an aircraft.

#### 2.3.1 Autopilot Mode

Prior to handing over the controls to an AI, it is important that the aircraft can be steered by an autopilot. The latter simply operates the heading and altitude of the plane, as well as the thrust. Once the autopilot is capable of steering the aircraft, the interfacing with an external AI is no longer an issue. In the absence of an external AI control, the algorithm provided with the Sandbox to simulate combat scenarios is fairly elementary:

1. Head to the target.
2. Determine the altitude of the target.
3. If the altitude is below a certain limit, climb to a certain altitude.
4. If the target is in range and aligned, fire the machine gun.
5. If the target is locked, fire a missile.
6. Wait between two missile shots (to avoid unloading everything at once).

## 2.4 Physics Model

In an air combat scenario, the physics of the aircraft is critical. The model used in the Sandbox is a compromise between simplicity and efficiency.

### 2.4.1 Air Density

Air density is an important parameter of the physical model because thrust power that is generated by the turbines is strongly dependant on air density. Air density is represented as  $\rho$ . For calculations, ideal gas constant R, molar mass of dry air M, and gravitational acceleration G are accepted as  $8.3144621 J/(mol.k)$ ,  $0.0289652 kg/mol$ ,  $9.80665 m/s^2$ , respectively. When the necessary calculations are done, it can be seen that the air density decreases as the altitude increases. Temperature is accepted absolute temperature. Note that temperature is kept constant throughout the simulation but it can also be varied.

$$\rho = \frac{p}{RT} \quad (1)$$

### 2.4.2 Velocity and Dynamic Pressure

The aircraft responses depend essentially on its velocity, i.e. on the pressure applied by the atmosphere on the wing (and other control surfaces). In the aerodynamic model of the Sandbox, the dynamic pressure is evaluated first in order to obtain a coefficient which will be used in several other parameters calculations.

$$q = \frac{\rho v^2}{2} \quad (2)$$

Where q is dynamic pressure.  $\rho$  is density of fluid and v is velocity. For the dynamic pressure in 3 axes, x, y and z components are calculated separately.

### 2.4.3 Forces

Several forces come into play in the flight simulation: gravity, lift, drag, turbine thrust. Lift and Drag will vary with dynamic pressure. The wind force is not taken into account.

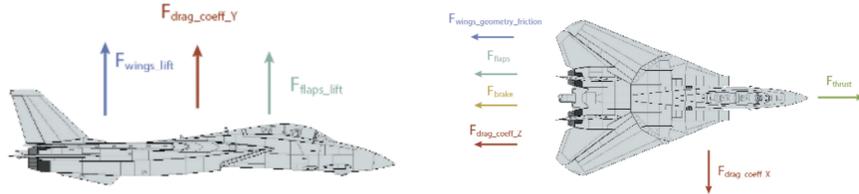


Figure 5: Aerodynamic forces acting on the aircraft model.

$$\vec{F}_{lift} = q_z \cdot (\vec{F}_{wings\_lift} + \vec{F}_{flaps\_lift}) \quad (3)$$

$$\vec{F}_{drag} = \frac{\vec{V}}{\|\vec{V}\|} \cdot q_z \cdot (\vec{F}_{drag\_coeff} + \vec{F}_{flaps} + \vec{F}_{wings\_geometry\_friction}) \quad (4)$$

$$\vec{F}_{move} = \vec{F}_{thrust} + \vec{F}_{lift} - \vec{F}_{drag} + \vec{F}_{gravity} \quad (5)$$

### 2.4.4 Moments

The pilot, by operating the controls, acts on the moments (angular forces) of the aircraft: pitch, yaw and roll. The rotational movement of the aircraft is the sum of these three moments.

$$\vec{M}_{pitch} = \vec{X}_{axis} \cdot q_z \cdot (pitch\_Level) \cdot (pitch\_friction) \quad (6)$$

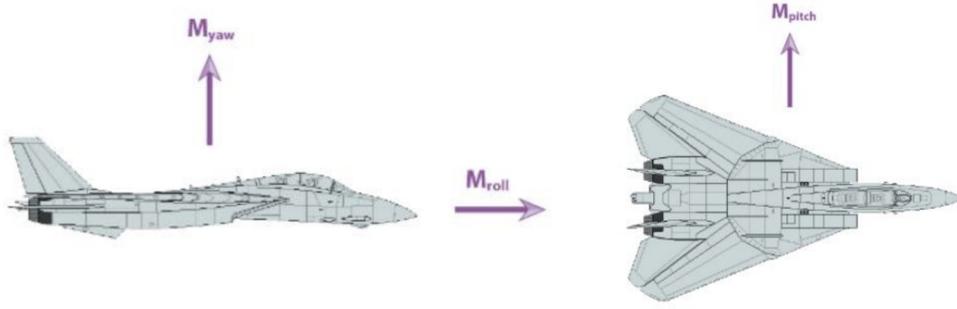


Figure 6: Aerodynamic moments acting on the aircraft model.

$$\vec{M}_{yaw} = \vec{Y}_{axis} \cdot q_z \cdot (yaw\_Level) \cdot (yaw\_friction) \quad (7)$$

$$\vec{M}_{roll} = \vec{Z}_{axis} \cdot q_z \cdot (roll\_Level) \cdot (roll\_friction) \quad (8)$$

Total moment can be found from;

$$\vec{M} = \vec{M}_{pitch} + \vec{M}_{roll} + \vec{M}_{yaw} \quad (9)$$

Where  $q_z$  is the local Z component of dynamic pressure vector. The Z axis is the front component of plane displacement. If  $q_z = 0$ , that means the plane only moves along its vertical (Y) and/or lateral (X) axis, the aircraft is in a stall and can no longer be controlled. If  $q_z = 1$ , that means the plane is moving along its frontal (Z) axis, the aircraft is in stable flight.

### 2.4.5 Customized Physics

$$\begin{bmatrix} 1 & \tan(\theta)\sin(\phi) & \cos(\phi)\tan(\theta) \\ 0 & \cos(\phi) & -\sin(\phi) \\ 0 & \sec(\theta)\sin(\phi) & \sec(\theta)\cos(\phi) \\ x & y & z \end{bmatrix}$$

The sandbox currently runs on basic physics dynamics; nevertheless, users can easily configure customized dynamics for either aircrafts or missiles. Using the function `set_machine_custom_physics_mode`, user can choose which aircraft or missile is to be modified. Then, the aerodynamics matrix demonstrated above can be utilized to transfer custom dynamics model into the simulation. Here,  $x$ ,  $y$  and  $z$  represent the Cartesian coordinates of the plane while  $\theta$  and  $\phi$  angles represent the data which come from the dynamics. After that, with the help of `update_machine_kinetics` function, position data from the dynamics function can be sent to sandbox for update at each time frame.

## 2.5 Graphics and Environment

The visual realism of the simulation is a key factor in creating an immersive experience for the Sandbox user, especially when using a virtual reality headset. The purpose of the Sandbox is multifaceted, but since one of the objectives is to confront human pilots with algorithm-driven aircraft, the quality of the immersion is central in the acceptance of the simulation.

### 2.5.1 Rendering Pipeline

The Sandbox uses a rendering pipeline that aims to provide a rich and complex visualisation without compromising the performance of an application that is implemented in the Python language.

The rendering is split into several passes that are calculated sequentially as follows, for each frame:

1. Rendering of water reflections
2. Rendering of the ocean and sky
3. Rendering of the 3D scene's main elements

On top of these 3 passes, two layers overlay 2D information about the status of the simulation as well as 3D information regarding the debugging process.

**Reflection Rendering Pass** The scene is rendered from a camera that is mirrored to the main view angle, relative to the plane of the ocean. All solid 3D objects are included in this pass, be they static or dynamic (terrain, ships, aircraft), okyanus uzerindeki yansimalar 7 gorselinde detayli gorunebilir.

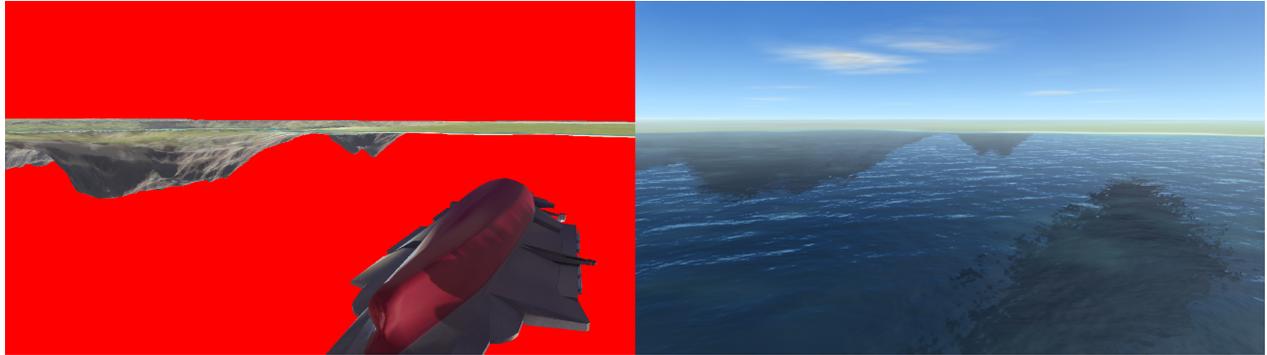


Figure 7: These images depict how the reflections on the ocean is rendered. Left image shows the scene without the ocean and right image shows the scene after the ocean is added.

**Ocean and Sky Rendering Pass** The ocean and the sky are rendered by a fairly simple raytracing routine included in a GPU shader. The framebuffer of the reflection in the water, obtained from the previous pass, is used for that purpose. The surface of the ocean is a 3D plane, which can intersect with 3D models of the scene (ship hulls, aircraft grounded in the water, emerged terrain). This is why the DepthBuffer is also generated in this raytracing calculation.

Our commitment to provide a solution that can be accessed on standard hardware specifications led us to avoid using RTX extensions.

**3D Scene Rendering Pass** The 3D scene that is the main part of the simulation (aircraft, vehicles, weapons, terrain and particles) is rendered through the primary camera ( subjective view camera, or any other external view of the aircraft). This scene is rendered over the ocean and the sky. As the ocean DepthBuffer has been rendered, the parts of the 3D models below the water surface are not drawn. Passes 2 and 3 thereby combine seamlessly to create a coherent image.

**3D Overlays Pass** The 3D texts as well as the debug displays (3D vectors, flight paths, collision boxes) are rendered in this pass. You can see the generation of the collision boxes as an example in Figure 8.

**2D Overlays Pass** The HUD and other instructions which describe the various phases of the simulation are shown in 2D, within the screen plane. For instance, HUD feature can be seen in Figure 9.

**Post-Processing** Once the 3D renders and overlays are displayed, a final pass is made which has no other purpose than aesthetics. At this point, it is mainly the fade in/fade out effects that are being applied to provide harmonious transitions between each phase.

## 2.5.2 Earth and Space Geometry

To increase visual realism, the curvature of the horizon and the atmospheric hue change according to the altitude of the aircraft are also simulated. In first person view, in exterior view, or in VR, the immersion is thus reinforced. The essence of the simulation enables us to restrict the altitude that the aircraft can reach thereby avoiding the complexity of

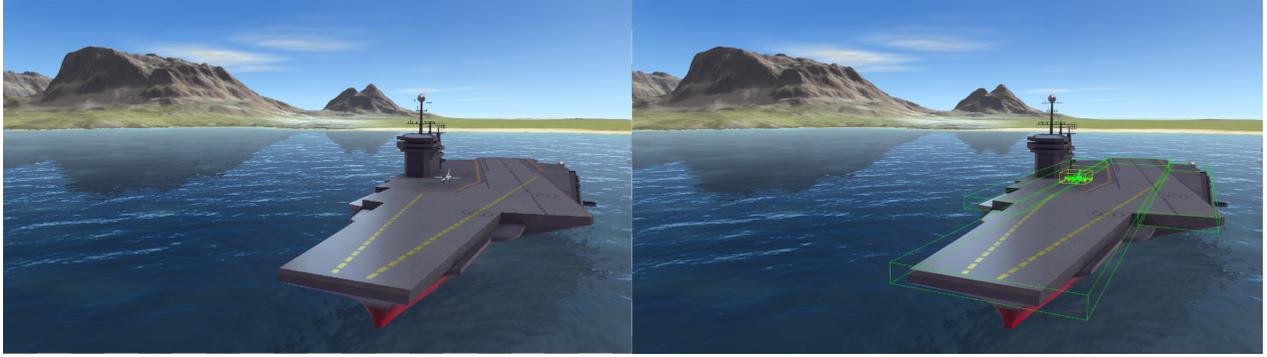


Figure 8: Three dimensional objects are created independently before being integrated. You can see the example of this aircraft carrier.



Figure 9: HUD feature provides useful indicators about the flight (altitude, pitch, roll, etc.).

a planetary visualization when viewed from orbit. Again, the suggested realism strengthens the acceptance factor of the simulation by the users, without increasing the complexity of the implementation whose primary goal is to be used as a sandbox.

**Intersection Equation** We assume that the distance between the camera and the Earth's surface is still very short compared to the planet's radius. We therefore require an equation that gives more accurate results than the standard sphere/ray intersection, by solving a second order equation. You can check the geometric illustrations for this calculations in Figure 10.

h: camera altitude

r: Earth's radius

$\alpha$ : angle between the ray and the surface normal

d: distance on the ray between the camera and the surface

Calculation of d:

$$d = \cos(\alpha) \times (r + h) - x$$

$$x = r \times \sin(\beta)$$

Calculation of  $\beta$ :

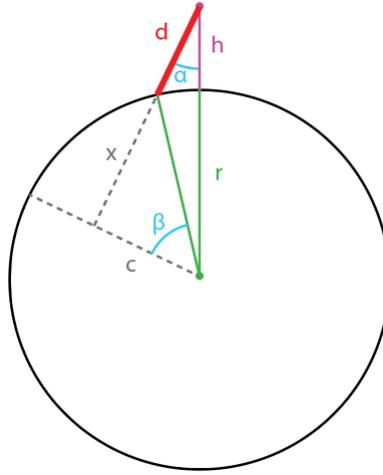


Figure 10: Geometric illustrations of the equations above.

$$c = \sin(\alpha) \times (r + h) = r \times \cos(\beta)$$

$$\cos(\beta) = \frac{\sin(\alpha) \times (r + h)}{r}$$

$$\beta = \text{acos} \left( \frac{\sin(\alpha) \times (r + h)}{r} \right)$$

We then have our distance d along the ray:

$$d = \cos(\alpha) \times (r + h) - r \times \sin \left( \text{acos} \left( \frac{\sin(\alpha) \times (r + h)}{r} \right) \right)$$

**Altitude Accuracy Issue** The method described in this document suffers from one limitation: at low altitudes, the numerical accuracy is not sufficient. If h is about 100 m, the radius of the Earth r is about 6,000,000 meters. The precision of 32-bit floating point numbers doesn't allow such a large range:

Low altitude  $h \ll r$ :

$$d = \frac{h}{\cos(\alpha)}$$

The answer is to mix the intersection with a plane and the intersection with the sphere. The planar intersection occurs when the magnitude of h is insignificant versus r, whereas the spherical intersection occurs when h grows big enough compared to r. For transitory levels of h relative to r, a fading applies between the two distances. The tests showed that the transition is not noticeable to the viewer. Figure 11 exemplifies a low and a high altitude sky view.

High altitude  $h \approx r$ :

$$d = \cos(\alpha) \times (r + h) - r \times \sin \left( \text{acos} \left( \frac{\sin(\alpha) \times (r + h)}{r} \right) \right)$$

$$d = \frac{h}{\cos(\alpha)}$$



Figure 11: Left image shows the sky view at low altitudes and right image shows the sky view at high altitudes.

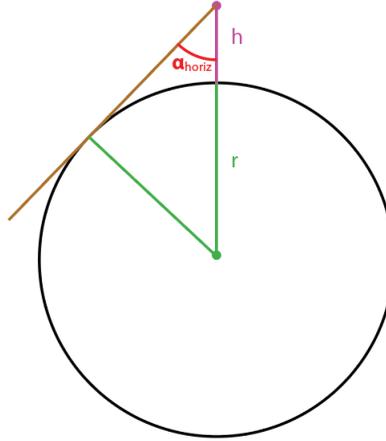


Figure 12: Geometric illustration of the ratio of the camera's altitude to the Earth's radius.

**Interpolation Between Planar and Spherical Distances** To compute the interpolation between the planar distance and the spherical distance, we need to figure out the angle between the normal of the surface of the sphere and the horizon line. Geometric illustration of the calculation is shown in Figure 12.

$$\alpha_{horiz} = \frac{\pi}{2} - \text{atan} \left( \frac{\sqrt{h \times (h + 2 \times r)}}{r} \right)$$

Then we compute the ratio of the camera's altitude to the Earth's radius:

$$ratio = \frac{h}{r}$$

1. If ratio < 0.001, then only the planar distance should be used.
2. If ratio > 0.01, then only the spherical distance is used.
3. If 0.001 < ratio < 0.01, an interpolation applies between the planar distance and the spherical distance, following the angle of the radius.

**Generating The Atmosphere** We define 3 colours:

1. Colour of the lower atmosphere (a tone of light blue).
2. Colour of the upper atmosphere (blue).

3. Colour of the outer space beyond the atmosphere (black).

Then we will need the angle between the line of the horizon and the edge of the atmosphere  $\alpha_{atm}$ . Two cases are possible:

1. The camera is within the atmosphere
2. The camera is above the atmosphere

To determine the colour of the atmosphere, we may choose between a "physical" and an "aesthetic" model.

The physical model is to estimate the amount of light diffracted and absorbed by the atmosphere. For this purpose, we would need to find out the thickness of the atmosphere crossed by the ray. However, this quickly leads to a complexity that is not necessarily relevant in this case.

For the Sandbox, the aesthetic approach was preferred, providing greater control over the atmospheric rendering as a function of the aircraft's altitude.

When the camera is inside the atmosphere, we can compute  $\alpha_{atm}$  like this:

$$\alpha_{atm} = \pi - \alpha_{horiz}$$

However, this poses a special difficulty when the camera is just above the atmosphere:

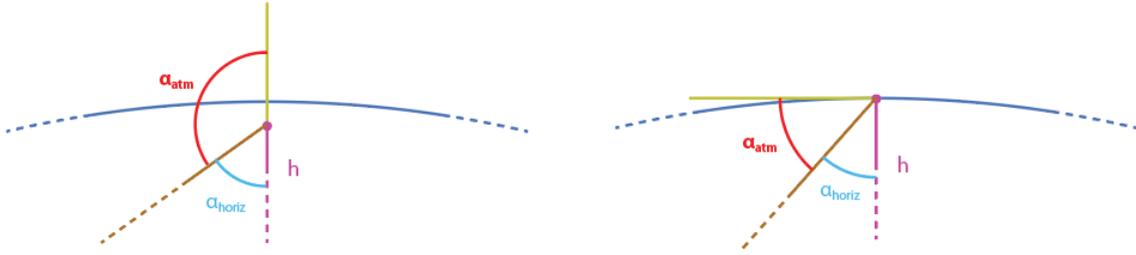


Figure 13: Change in  $\alpha_{atm}$  according to  $h$ .

As seen from figure 13 when  $h$  goes above the atmosphere,

$$\alpha_{atm} = \frac{\pi}{2} - \alpha_{horiz}$$

To circumvent this sharp edge, we can interpolate the angles along the thickness of the atmosphere. To do this, we determine the parametric altitude  $F$  of the camera within the atmosphere:

$$F = \frac{h}{A_t}$$

#### **F <= 1 : Camera Within The Atmosphere**

$$\alpha_{atm} = \pi \times (1 - F) + \left(\frac{\pi}{2} \times F\right) - \alpha_{horiz}$$

An interpolation applies to the colour beyond the atmosphere:

1. F=0 : SpaceColor = colour of the upper atmosphere
2. F=1 : SpaceColor = colour of space beyond the atmosphere

**F > 1 : Camera Above The Atmosphere** In this case, the equation is a variation of the calculation of the angle between the normal of the sphere and the horizon line.

$$\alpha_{atm} = \frac{\pi}{2} - \text{atan} \left( \frac{\sqrt{(h - A_t) \times ((h - A_t) + 2 \times (r + A_t))}}{r + A_t} \right) - \alpha_{horiz}$$

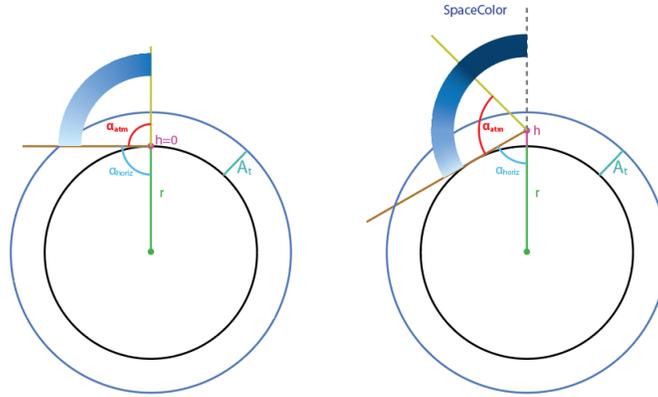


Figure 14: Sky/Space color and viewing angles when the camera is within the atmosphere.

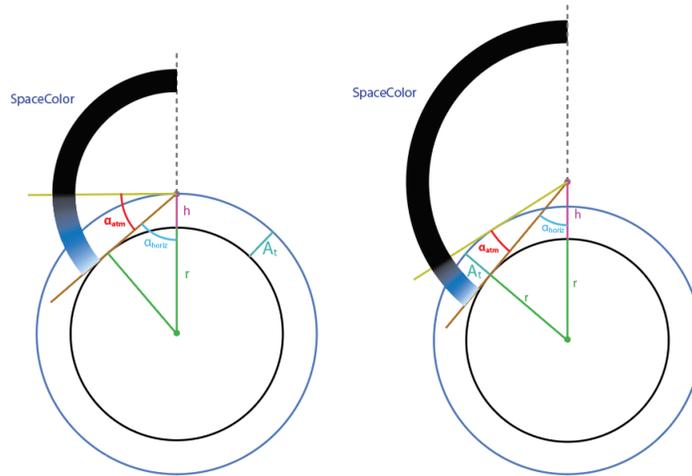


Figure 15: Sky/Space color and viewing angles when the camera is above the atmosphere.

### 2.5.3 Pursuit Camera Model

The camera is aimed at a position slightly delayed in time from the aircraft. To render the movement more realistically, a subtle Perlin noise is added to the orientation. The intensity of this noise varies according to three parameters: the velocity of the aircraft, the afterburner activation, the acceleration intensity.

## 2.6 Renderless Mode

It is known that 3D rendering with high graphics, puts a lot of stress on the processing units and inhibits the performance of the computer. It also unnecessarily raises the system requirements for training, making it impossible to run on many device. Once the renderless mode is turned on, rendering is deactivated while the algorithm continues learning in the background. Thus, excessive performance constraint is disposed of. Renderless mode also makes running the software on server side possible since the program does not impose any graphical (visual) output.

## 3 Experiment

So far, we have discussed diverse scenarios that can be implemented in Harfang. In this part, for the purpose of validating the usability of our simulation environment, a simple navigation task is experimented. The details of this experiment are discussed below.

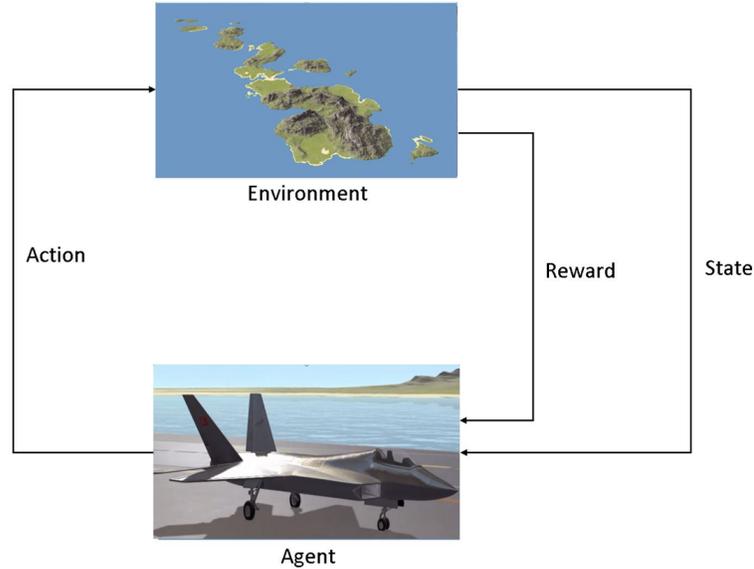


Figure 16: Interaction of a reinforcement learning agent with the environment.

### 3.1 Scenario



Figure 17: Playground of the agent. The red dot in the center represents the goal and the two black squares are the borders of starting region.

In this experiment, agents are initialized at a random position and altitude within the region lies between the two black square shown in Figure 17. The task is to reach the destination point which is marked as red regardless of the starting position.

### 3.2 Method

TD3 + HER are used in this project, TD3 [32] was published in 2018. The main difference of this method when compared to other methods is that even if it uses actor-critic as a base, we cannot deny that classical methods like DDPG [33] have some problem. Those problems are mainly an overestimation, it basically can be explained as giving actions to agents in places it did not present, so this is the main reason why TD3 + HER come up with better results when compared to DDPG and classical actor-critic methods. We tested DDPG at the beginning, but we did not able to perform even basic actions in DDPG, that is why we have moved to TD3 instead.

---

### 3.3 Action Space

For this problem, we adopted a low-level control approach. This means that the agent directly manipulates the control surfaces of the plane. Therefore, there were 3 actions outputted by the algorithm: rudder, elevator, aileron. The reason why we choose a low-level control strategy is because we intended to see what the agent is capable of learning without limiting its abilities.

### 3.4 State Space

The agent had access to 13 state signals in total: (x, y, z) coordinate differences between the goal and the agent,  $(\phi, \theta, \psi)$  Euler angles, horizontal and vertical velocity, heading angle, pitch attitude, and acceleration in (x, y, z) directions. Although there were other signals present in the simulation environment, we decided that these signals were adequate in order for agent to accomplish the task.

### 3.5 Reward

We tried to keep the reward as simple as possible without hindering the agents ability to learn how to reach the goal. Hence, the reward structured according to the distance of the plane from the goal position as depicted in the formula below. This reward formula aims to promote the agent going towards the goal as it receives higher reward with the decreasing distance because the reward function outputs a negative value. In addition to that, the agent receives +100 reward if it reaches the goal, and -100 if it crashes or cannot reach the goal within the duration of an episode.

$$R = -10^{-5} \sqrt{(UAV_x - Goal_x)^2 + (UAV_y - Goal_y)^2 + (UAV_z - Goal_z)^2} \quad (10)$$

Table 3: Hyperparameter setting used in the experiment.

Hyper-Parameter	Values
Critic Learning Rate	$10^{-3}$
Actor Learning Rate	$10^{-3}$
Optimizer	Adam
Target Update Rate ( $\tau$ )	$5.10^{-3}$
Batch Size	100
Iterations per time step	1
Discount Factor	0.99
Normalized Observations	True
Gradient Clipping	False
Actor First Hidden Layer Neuron number	400
Actor Second Hidden Layer Neuron number	300
Critic First Hidden Layer Neuron number	400
Critic Second Hidden Layer Neuron number	300

## 4 Result

After a training of about 1400 episodes, our agent was capable of navigating to the desired destination. It can be seen from Figure 6 that the learning slowed down at about 800 episodes and came to a convergence at 1000 episodes. Even though there were some visible oscillations in the reward before 1000 episodes, at the end, our agent managed to exploit the correct behaviours and successfully learned to perform this basic navigation task. This results proves that Harfang3D Dog-Fight Sandbox is an eligible platform for reinforcement learning research offering a broad array of scenarios, high accessibility, easy implementation, and extremely flexible training opportunities with renderless mode and distributional structure.

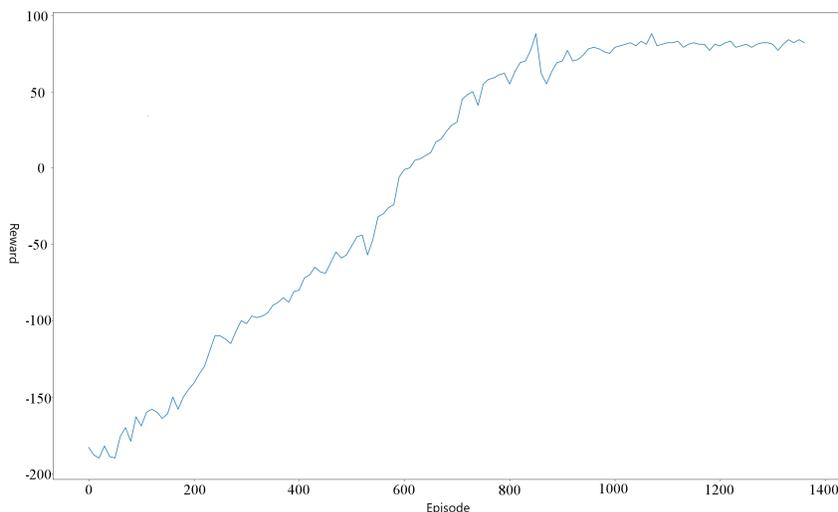


Figure 18: Episodic reward plot of the training phase.

## Acknowledgements

We would like to thank Thomas Simonnet, François Gutherz, Philippe Herber, Serge Bidault and the other members of the HARFANG team for the technical support.

## References

- [1] M. J. Matari ´c, “Reinforcement learning in the multi-robot domain,” in *Robot colonies*. Springer, 1997, pp. 73–83.
- [2] M. Asada, S. Noda, S. Tawaratsumida, and K. Hosoda, “Purposive behavior acquisition for a real robot by vision-based reinforcement learning,” *Machine learning*, vol. 23, no. 2-3, pp. 279–303, 1996.
- [3] V. Gullapalli, J. A. Franklin, and H. Benbrahim, “Acquiring robot skills via reinforcement learning,” *IEEE Control Systems Magazine*, vol. 14, no. 1, pp. 13–24, 1994.
- [4] S. Mahadevan and J. Connell, “Automatic programming of behaviorbased robots using reinforcement learning,” *Artificial intelligence*, vol. 55, no. 2-3, pp. 311–365, 1992.
- [5] M. Andrychowicz, F. Wolski, A. Ray, J. Schneider, R. Fong, P. Welinder, B. McGrew, J. Tobin, O. P. Abbeel, and W. Zaremba, “Hindsight experience replay,” in *Advances in Neural Information Processing Systems*, 2017, pp. 5048–5058.
- [6] M. Andrychowicz, B. Baker, M. Chociej, R. Jozefowicz, B. McGrew, J. Pachocki, A. Petron, M. Plappert, G. Powell, A. Ray et al., “Learning dexterous in-hand manipulation,” *arXiv preprint arXiv:1808.00177*, 2018.
- [7] S. Gu, E. Holly, T. Lillicrap, and S. Levine, “Deep reinforcement learning for robotic manipulation with asynchronous off-policy updates,” in *2017 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2017, pp. 3389–3396.
- [8] A. A. Rusu, M. Veerk, T. Rothrl, N. Heess, R. Pascanu, and R. Hadsell, “Sim-to-real robot learning from pixels with progressive nets,” in *Proceedings of the 1st Annual Conference on Robot Learning*, ser. *Proceedings of*

- 
- Machine Learning Research, S. Levine, V. Vanhoucke, and K. Goldberg, Eds., vol. 78. PMLR, 13–15 Nov 2017, pp. 262–270. [Online]. Available: <http://proceedings.mlr.press/v78/rusu17a.html>
- [9] J. Hwangbo, J. Lee, A. Dosovitskiy, D. Bellicoso, V. Tsounis, V. Koltun, and M. Hutter, “Learning agile and dynamic motor skills for legged robots,” *Science Robotics*, vol. 4, no. 26, 2019. [Online]. Available: <https://robotics.sciencemag.org/content/4/26/eaau5872>
- [10] Z. Xie, G. Berseth, P. Clary, J. Hurst, and M. van de Panne, “Feedback control for cassie with deep reinforcement learning,” in 2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). IEEE, 2018, pp. 1241–1246.
- [11] S.-H. Hsu, S.-H. Chan, P.-T. Wu, K. Xiao, and L.-C. Fu, “Distributed deep reinforcement learning based indoor visual navigation,” in 2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). IEEE, 2018, pp. 2532–2537.
- [12] J. Choi, K. Park, M. Kim, and S. Seok, “Deep reinforcement learning of navigation in a complex and crowded environment with a limited field of view,” in 2019 International Conference on Robotics and Automation (ICRA). IEEE, 2019, pp. 5993–6000.
- [13] Y. Zhu, R. Mottaghi, E. Kolve, J. J. Lim, A. Gupta, L. Fei-Fei, and A. Farhadi, “Target-driven visual navigation in indoor scenes using deep reinforcement learning,” in 2017 IEEE international conference on robotics and automation (ICRA). IEEE, 2017, pp. 3357–3364.
- [14] G. Kahn, A. Villafior, B. Ding, P. Abbeel, and S. Levine, “Self supervised deep reinforcement learning with generalized computation graphs for robot navigation,” in 2018 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2018, pp. 1–8.
- [15] 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*, 2004, pp. 799–806.
- [16] F. Sadeghi and S. Levine, “CAD2RL: Real single-image flight without a single real image,” arXiv preprint [arXiv:1611.04201](https://arxiv.org/abs/1611.04201), 2016.
- [17] C. Chen, Y. Liu, S. Kreiss, and A. Alahi, “Crowd-robot interaction: Crowd-aware robot navigation with attention-based deep reinforcement learning,” in 2019 International Conference on Robotics and Automation (ICRA). IEEE, 2019, pp. 6015–6022.
- [18] S. Christen, S. Stevsic, and O. Hilliges, “Guided deep reinforcement learning of control policies for dexterous human-robot interaction,” arXiv preprint [arXiv:1906.11695](https://arxiv.org/abs/1906.11695), 2019.
- [19] M. Everett, Y. F. Chen, and J. P. How, “Motion planning among dynamic, decision-making agents with deep reinforcement learning,” in 2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). IEEE, 2018, pp. 3052–3059.
- [20] G. Sartoretti, J. Kerr, Y. Shi, G. Wagner, T. S. Kumar, S. Koenig, and H. Choset, “Primal: Pathfinding via reinforcement and imitation multi-agent learning,” *IEEE Robotics and Automation Letters*, vol. 4, no. 3, pp. 2378–2385, 2019.
- [21] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba, “OpenAI Gym,” arXiv preprint [arXiv:1606.01540](https://arxiv.org/abs/1606.01540), 2016.
- [22] L. Fan, Y. Zhu, J. Zhu, Z. Liu, O. Zeng, A. Gupta, J. Creus-Costa, S. Savarese, and L. Fei-Fei, “Surreal: Open-source reinforcement learning framework and robot manipulation benchmark,” in *Conference on Robot Learning*, 2018, pp. 767–782.
- [23] J. Liang, V. Makoviychuk, A. Handa, N. Chentanez, M. Macklin, and D. Fox, “Gpu-accelerated robotic simulation for distributed reinforcement learning,” in *Proceedings of The 2nd Conference on Robot Learning*, ser. *Proceedings of Machine Learning Research*, A. Billard, A. Dragan, J. Peters, and J. Morimoto, Eds., vol. 87. PMLR, 29–31 Oct 2018, pp. 270–282. [Online]. Available: <http://proceedings.mlr.press/v87/liang18a.html>
- [24] L. Espeholt, H. Soyer, R. Munos, K. Simonyan, V. Mnih, T. Ward, Y. Doron, V. Firoiu, T. Harley, I. Dunning, S. Legg, and K. Kavukcuoglu, “IMPALA: Scalable distributed deep-RL with importance weighted actor-learner architectures,” in *Proceedings of the 35th International Conference on Machine Learning*, ser. *Proceedings of Machine Learning Research*, J. Dy and A. Krause, Eds., vol. 80. Stockholmsmssan, Stockholm Sweden: PMLR, 10–15 Jul 2018, pp. 1407–1416. [Online]. Available: <http://proceedings.mlr.press/v80/espeholt18a.html>
- [25] E. Liang, R. Liaw, R. Nishihara, P. Moritz, R. Fox, K. Goldberg, J. Gonzalez, M. Jordan, and I. Stoica, “RLlib: Abstractions for distributed reinforcement learning,” in *Proceedings of the 35th International Conference on Machine Learning*, ser. *Proceedings of Machine Learning Research*, J. Dy and A. Krause, Eds., vol. 80. Stockholmsmssan, Stockholm Sweden: PMLR, 10–15 Jul 2018, pp. 3053–3062. [Online]. Available: <http://proceedings.mlr.press/v80/liang18b.html>

- 
- [26] J. Tan, T. Zhang, E. Coumans, A. Iscen, Y. Bai, D. Hafner, S. Bohez, and V. Vanhoucke, "Sim-to-real: Learning agile locomotion for quadruped robots," in *Proceedings of Robotics: Science and Systems*, Pittsburgh, Pennsylvania, June 2018.
- [27] I. Higgins, A. Pal, A. Rusu, L. Matthey, C. Burgess, A. Pritzel, M. Botvinick, C. Blundell, and A. Lerchner, "Darla: Improving zeroshot transfer in reinforcement learning," in *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. JMLR. org, 2017, pp. 1480–1490.
- [28] H. Bharadhwaj, Z. Wang, Y. Bengio, and L. Paull, "A data-efficient framework for training and sim-to-real transfer of navigation policies," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 782–788.
- [29] X. B. Peng, M. Andrychowicz, W. Zaremba, and P. Abbeel, "Sim-to-real transfer of robotic control with dynamics randomization," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2018, pp. 1–8.
- [30] F. Muratore, F. Treede, M. Gienger, and J. Peters, "Domain randomization for simulation-based policy optimization with transferability assessment," in *Conference on Robot Learning*, 2018, pp. 700–713.
- [31] A. Mandlekar, Y. Zhu, A. Garg, L. Fei-Fei, and S. Savarese, "Adversarially robust policy learning: Active construction of physicallyplausible perturbations," in *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2017, pp. 3932– 3939.
- [32] S. Fujimoto, V. Hoof, H., and D. Meger, "Addressing function approximation error in actor-critic methods", arXiv preprint arXiv:1802.09477, 2018.
- [33] Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D. and Wierstra, D., "Continuous control with deep reinforcement learning", *International Conference on Learning Representations*, 2016.