

# Emergency Pilot: Automated Flight Guidance After a Loss of Thrust Based on Deep Reinforcement Learning

Walter Laurito <sup>1,\*</sup>, Renè Titze <sup>2</sup> and Wolfram Schiffmann <sup>2</sup>

<sup>1</sup>*FZI Research Center of Information Technology, Karlsruhe, Germany*

<sup>2</sup>*Faculty of Mathematics and Computer Science, FernUniversität in Hagen, Germany*

Correspondence\*:  
Walter Laurito  
Laurito@fzi.de

## 2 ABSTRACT

3 The fatality rate for flights involving light, fixed-wing aircraft is relatively high, and most accidents  
4 occur after a loss of thrust in the lower airspace. In these emergency situations, with a lack  
5 of potential energy, an assistance or automated system could lead the aircraft safely to an  
6 appropriate landing field. Some solutions for path planning and guidance exist, however most of  
7 them rely on a simplified model of the environment and the aircraft's dynamics to generate a path.  
8 Thus, those solutions tend to be rigid and less reliable during an emergency; especially in the  
9 occurrence of wind. Moreover, many solutions don't consider the expected landing direction and  
10 the heading of the aircraft has to have, when reaching the landing field. In this work, we tackled  
11 these issues by focusing on the creation of a real-time guidance system for 3D trajectory planning  
12 after a loss of thrust based on deep reinforcement learning (DRL). In DRL, an agent learns  
13 through trial and error by interacting with an environment. DRL is especially useful in uncertain  
14 environments, where many parameters can't be calculated in advance, which is the case in an  
15 emergency. Therefore, to train the agent to guide a fixed-wing, engines-off aircraft to an arbitrary  
16 target position, we developed and implemented multiple simulation environments. Furthermore,  
17 we incorporated wind into one of these environments. The created software package of the  
18 environments can be found online <sup>1</sup>. Usually, complex calculations are needed to model the  
19 engines-off flight dynamics and to generate a 3D path (guidance) under wind. With DRL these  
20 calculations can be avoided. By using shaped reward functions, we trained a neural network to  
21 successfully select directions and glide angles to lead the aircraft to an arbitrary, chosen landing  
22 field in real-time while avoiding accidents. The success rate during our experiments was high:  
23 In most cases, the aircraft reaches the target position from the correct direction and with the  
24 expected heading.

25 **Keywords:** Deep Reinforcement Learning, Aviation, Flight Guidance, Emergency Landing, Trajectory Planning

## 1 INTRODUCTION

26 Compared to commercial aviation, the fatality rate for flights involving light, fixed-wing aircraft is still  
27 relatively high (*Air safety statistics in the EU - Statistics Explained 2020*), where many accidents occur  
28 after a loss of thrust (Dorr 2018). After a loss of thrust, the pilot has only a short amount of time to find a  
29 safe path to a nearby landing field. The pilot has to consider multiple things in parallel: The remaining

---

<sup>1</sup> Flight guidance env [https://github.com/lauritowal/guidance\\_flight\\_env](https://github.com/lauritowal/guidance_flight_env)

30 altitude to the ground, other passengers, the aircraft's velocity and probably even wind. Particularly  
31 inexperienced pilots might be overwhelmed by the fast-paced emergency situation during the flight. Instead  
32 of the inexperienced pilot, an automated system could lead the aircraft to a close landing field and ideally  
33 prevent fatal accidents.

34 Such an automated system could be structured as a Guidance, Navigation and Control (GNC) system.  
35 The guidance system, part of the GNC system, is responsible for generating a path to a specific target,  
36 a near landing field. The guidance system obtains the aircraft's location, velocity and attitude from a  
37 navigation system. From the obtained information, the guidance system generates the path and provides  
38 desired directions (headings) of the aircraft to a control system. Then, the control system leads the aircraft  
39 to follow these directions and eventually the aircraft moves along the path. This work focuses mostly on the  
40 creation of a real-time guidance system. The guidance system generates a 3D path that leads a fixed-wing  
41 aircraft after a loss of thrust to a near landing field.

42 *Reinforcement Learning* (RL) is a subfield of *Machine Learning* (ML), which is well suited for uncertain,  
43 changing environments. In RL, an agent learns directly by interacting with an environment. After learning,  
44 a successful agent is able to generalize (up to a certain point) to a newly presented situation by the  
45 environment. This means, the agent is capable of selecting the optimal action for a given observation at  
46 each time step. The combination of deep neural networks (Goodfellow, Bengio, and Courville 2016) with  
47 reinforcement learning led to the term *Deep Reinforcement Learning* (DRL). In the past few years, DRL  
48 was applied to a variety of complex problems and demonstrated remarkable achievements, for example, in  
49 the field of robotics (Akkaya et al. 2019; Gu et al. 2016), digital games (Mnih et al. 2013; Berner et al.  
50 2019 and board games (Schrittwieser et al. 2020; Vinyals et al. 2019). Using DRL for flight guidance has  
51 the following advantages:

- 52 •Although the training phase of an agent can be relatively long, after the training phase is completed, the  
53 agent can be used in real-time. The trained agent does not need to perform computationally expensive  
54 operations to generate a path. During the flight, decisions for the aircraft's next direction can be made  
55 nearly instantly.
- 56 •A trained DRL-agent could probably perform better than humans and even reach superhuman perfor-  
57 mance, as shown for example in (Mnih et al. 2013) and (Silver et al. 2017). Therefore, when combined  
58 with a conventional control system, DRL could be ideal for guidance in an emergency, since a trained  
59 agent could reach a landing field where even a human expert pilot would fail.
- 60 •There is no need to have a model of the aircraft, neither of the surroundings. Therefore, the same approach  
61 must be used to train an agent on different aircraft and environments.

62 There are works applying DRL to path planning in aviation. However, to our knowledge, no work exists  
63 that uses DRL to specifically address the presented problem above: Generating a 3D path to a landing field,  
64 for a fixed-wing aircraft in case of total loss of thrust. The exceptional achievements in the field of DRL  
65 and the described advantages to flight guidance, led to the decision of the authors of this work, to study the  
66 potential of DRL applied to the described problem above.

67 In this work, a trained DRL agent acts as the guidance system. To train the agent, three custom OpenAI  
68 Gym (Brockman et al. 2016) environments were developed. Those environments were created with a  
69 typical emergency situation in mind: Low remaining altitude, loss of thrust, and additional wind. After the  
70 training, the agent was evaluated in different experiments. Overall, in many cases, the agent guided the  
71 aircraft successfully to the target. In one of the experiments, it achieved a high mean success rate of 73 %.

72 For the training, the light Cessna172P aircraft was selected inside the simulation. The reason for this  
73 selection is stated at the beginning of this section: The fatality rate for flights involving light, fixed-wing  
74 aircraft is still relatively high (*Air safety statistics in the EU - Statistics Explained* 2020). Yet, since  
75 model-free DRL is used, it is not necessary for the DRL-agent to receive information about the specific  
76 aircraft in advance. The agent learns about the aircraft dynamics from experience during the training phase.  
77 This has the advantage that the methods developed in this work can easily be extended and applied to other  
78 types of aircraft, by simply using a different aircraft model during training.

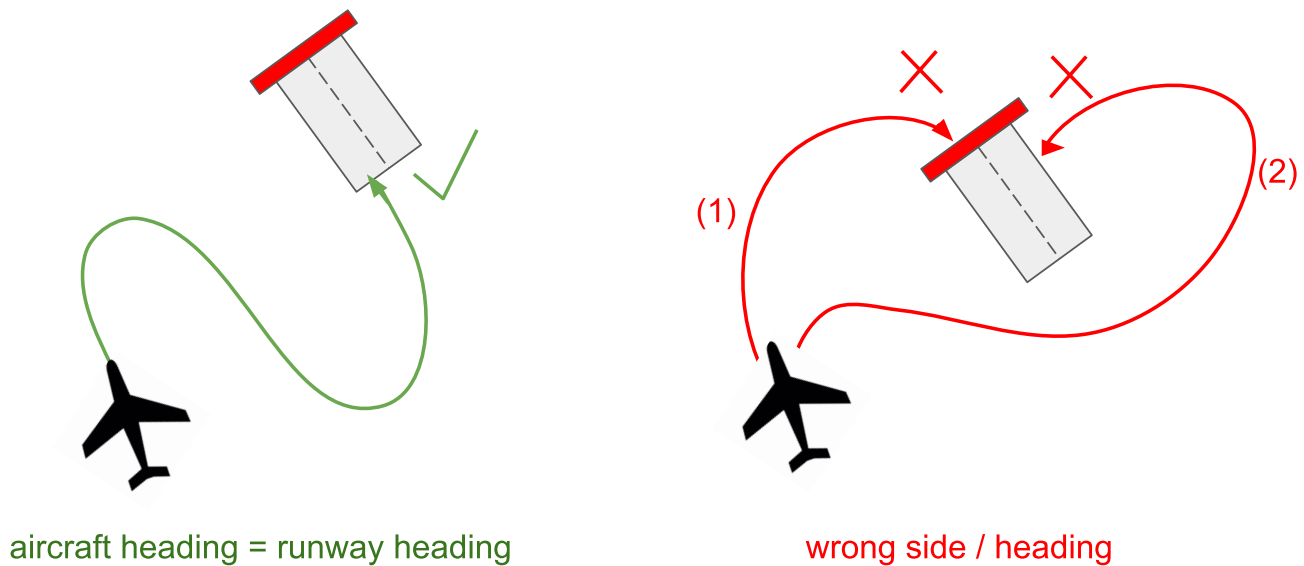
79 The main contribution of this work is a first insight into the applicability of DRL to create a flight  
80 emergency guidance system for fixed-wing aircraft. This guidance system is responsible for generating a  
81 3D path to a landing field after a loss of thrust. Furthermore, the developed and implemented environments  
82 used to train and evaluate the DRL agent can be used for future research by others.

83 The remainder of this paper is structured as follows: Section 2 states the problem in more detail. Section 3  
84 discusses related works. Section 4 introduces fundamentals and background knowledge. Section 5 presents  
85 the methods. Section 6 contains the experiments with the setup. Section 7 presents the results, and Section  
86 8 the discussion. Finally, the last Section 9 provides the conclusion and an outlook for further work.

## 2 PROBLEM STATEMENT

87 In this work, the goal is to study the potential of deep reinforcement learning applied to the creation of a  
88 flight emergency guidance system. A DRL agent acts as the guidance system and shall have the following  
89 capabilities:

- 90 •Generating a 3D path from the aircraft position to an arbitrary target position. This is done by the guidance  
91 system communicating the desired direction (heading) to the control system in each time step.
- 92 •Considering that the initial configuration of the aircraft is arbitrary. This includes the heading, altitude,  
93 position, etc.
- 94 •Guiding the engines-off aircraft (loss of thrust) to reach the target with the expected heading. During the  
95 last part of the landing procedure, the final approach, the aircraft needs to turn into line with the runway  
96 to be able to proceed to the round-out stage smoothly (Crocker 2007). The round-out stage, at around 15  
97 ft above the runway, is the last phase just before touchdown (United States Department of Transportation  
98 2016; Allerton 2009, p. 8–6.). In a normal situation, during the round-out stage, the nose of the aircraft is  
99 raised for some time to reduce the rate of descent before touching the ground (Allerton 2009, p. 194).  
100 Therefore, to be able to proceed to the round-out stage smoothly, a major weight shall be given to the  
101 following: The aircraft must end the final approach with the expected heading, which includes the fact  
102 that the aircraft needs to reach the landing field from the correct side, as can be seen in Figure 1. Reaching  
103 the landing field from the correct side is especially important for emergency situations, since obstacles  
104 could be found on the other side of the landing field. Moreover, the influence of existing wind on the  
105 aircraft affect the choice of the particular direction of the landing field ((ASA) 2017, Chapter 8, p. 14).
- 106 •Leading the aircraft also under the influence of constant wind, since wind can affect the flight path  
107 significantly. Wind speed in the range of 0 to 3000 ft/min (Beaufort scale 1-7) shall be considered.



**Figure 1.** The image on the left shows a correct path to reach the runway with the expected heading. The image on the right shows two incorrect paths, where the aircraft would reach the runway with a high heading error. Here, the red rectangle at the top of the runway indicates the wrong side for landing.

108 For the agent to obtain the above capabilities, it needs to be trained. Therefore, multiple environments  
109 need to be created to train the agent in simulation.

110 In addition to the guidance system, a control system is essential, which receives the guidance system's  
111 instructions (desired heading and pitch angle). Based on those instructions, the control system is responsible  
112 to output necessary commands to correct the aircraft's orientation. Furthermore, the control system is used  
113 for maintaining the *best glide speed*. The best glide speed, is the speed that allows the aircraft to travel the  
114 greatest forward distance for each increment of altitude lost (Administration 2004). For the Cessna172P,  
115 the best glide speed is at around 65 Knots-Indicated Air Speed (KIAS) (*Best Glide Speed and Distance*  
116 2018). Only a minimal control system shall be created, since the control system is not the focus of this  
117 work.

118 A navigation system is assumed to be available, which is part of many aircraft, and also available in the  
119 simulation in this work. Usually, the navigation system obtains information about the aircraft's location,  
120 velocity and attitude from the aircraft's sensors and provides this information to the guidance system  
121 Allerton 2009, p. 247.

122 Success is defined by the aircraft's distance to the target and the heading error at the end of the episode.  
123 Furthermore, the difference in altitude of the aircraft to the target position needs to be just above ground  
124 level.

125 The landing field can be a conventional runway or any other quite flat area which can be used for an  
126 emergency landing (for example, a grass field).

127 Guiding the aircraft through the round-out stage and touchdown phase is not part of this work. Moreover,  
128 considering obstacles on the flight path is left to future research.

### 3 RELATED WORKS

129 Klein et al. have proposed the Emergency Landing Assistant (ELA) for an aircraft after a loss of thrust  
130 (Klein, Klos, Lenhardt, et al. 2018; Klein, Klos, and Schiffmann 2020). There, Dubins paths were generated  
131 to obtain a flight path from an arbitrary aircraft's configuration (position, heading, etc.) to an emergency  
132 landing field, while considering constant wind. The main idea of their approach was the following: Moving  
133 the landing field's threshold opposite to the wind direction by a length of the wind speed times the glide  
134 time. This created a new virtual target position. Next, by following the path to this virtual threshold instead,  
135 the aircraft would be compensating for the displacement caused by the constant wind. Eventually, the  
136 aircraft would then reach the threshold of the real landing field on the ground. Klein et al. have inspired the  
137 author of this thesis to create an alternative solution to ELA, based on DRL. Stephan et al. also used Dubins  
138 curves successfully to generate 3D paths to an emergency landing field, however they did not consider  
139 constant wind (Stephan and Fichter 2016).

140 The presented works demonstrated solutions for generating paths based on simple Dubins curves, however  
141 some shortcomings still exist. First, the radius  $r$  of the turns (L, R) needs to be provided in advance. The  
142 turn radius depends on the aircraft's properties and current flight dynamics (forward speed, turn speed,  
143 roll angle, etc.). For this reason, a model of the aircraft is needed to select the correct radius. Hence, the  
144 aircraft's model has to be quite realistic to ensure that the aircraft in emergency can actually carry out the  
145 required maneuvers needed to turn and follow the path. This also implies that, inconveniently, for each  
146 aircraft a different model is required. Furthermore, the necessity to provide the radius in advance for the  
147 path generation also reduces available maneuvers, because the radius can't be changed during an aircraft's  
148 ongoing turn operation. Moreover, having a fixed radius restricts the aircraft to turn with the same roll  
149 angle (bank angle) and velocity. This means, that unexpected perturbations (e.g., obstacles, changes in  
150 wind speed or direction, pilot errors, etc.) could render a generated path useless. Klein et al. tried to solve  
151 this by recalculating new paths during the flight (Klein, Klos, and Schiffmann 2020), however it could still  
152 be impossible following these generated paths. Furthermore, the bank angle and speed affects the aircraft's  
153 available remaining altitude and time before reaching the ground. Therefore, during an emergency, it is  
154 of major importance to select an adequate radius. Using a fixed radius for all turns however could mostly  
155 result in a compromise and not in the optimal solution for a specific situation. Another shortcoming in  
156 (Klein, Klos, Lenhardt, et al. 2018) is that in some cases the position of the virtual landing field can't be  
157 calculated directly and needs to be approximated. This could lead to the aircraft not terminating exactly at  
158 the correct position of the actual landing field.

159 As described above, by creating a solution based on deep reinforcement learning no model in advance  
160 is needed, which eliminates the problem of creating a model and specifying its parameters. Therefore,  
161 the same approach based on DRL could easily be applied for training models using different aircraft.  
162 Furthermore, setting a turn radius in advance is not needed, since a trained DRL agent can decide to adjust  
163 the turn radius as needed during the flight. The path is generated in real-time and not beforehand, which  
164 makes the DRL solution more flexible. The DRL agent would be trained on different scenarios to be able  
165 to generate an optimal path for a specific situation.

166 Apart from approaches based on Dubins paths, there are other solutions. One of them is (Váňa et al. 2018).  
167 Váňa et al. proposed a solution based on Rapidly-exploring Random Trees (RRT\*). There, trajectories to  
168 potential emergency landing fields to the aircraft's location were generated and collected at any time during  
169 the flight. In case of a loss of thrust, a path to a landing site was then already available and thus needed not  
170 to be generated via the computational expensive RRT\* algorithm. However, there were cases where not  
171 enough trajectories could be generated in time, for example in the case of a loss of thrust directly after take

---

172 off. Furthermore, small changes occurring during the flight were not taken into account. The authors did  
173 not consider wind, for example, which can have a substantial influence on the quality of a selected path.

174 Most literature on using deep reinforcement learning for flight guidance focused on path planning in  
175 situations where the aircraft's engines were fully functional. As described before, no previous work was  
176 found that used DRL to specifically address the presented problem above: Generating a 3D path to a  
177 landing field for a fixed-wing aircraft after a loss of thrust. Nonetheless, there are of course works which  
178 addressing the problem of path generation for aircraft in general. Those used a solution based on DRL.  
179 Some are shortly introduced.

180 In (Z Wang et al. 2018), a simple Deep Q-Network (DQN) was used to train an agent model for guidance,  
181 leading the fixed-wing aircraft to a landing field. During the training, the agent was rewarded for landing  
182 successfully and penalized for leaving a specific sector or if running out of fuel. However, the approach  
183 was simplified to the 2D space. Moreover, the occurrence of wind was not investigated.

184 In (Yan, Xiang, and C Wang 2019), the authors used a Dueling Deep Q Network (D3QN) algorithm  
185 in a dynamic environment to train an Unmanned Aerial Vehicle (UAV) for path planning while avoiding  
186 static and moving obstacles. The paper did not provide a solution in a situation with wind. Furthermore,  
187 the simplification of the environment to a 2D image map makes it difficult to apply the solution during a  
188 real-world emergency landing.

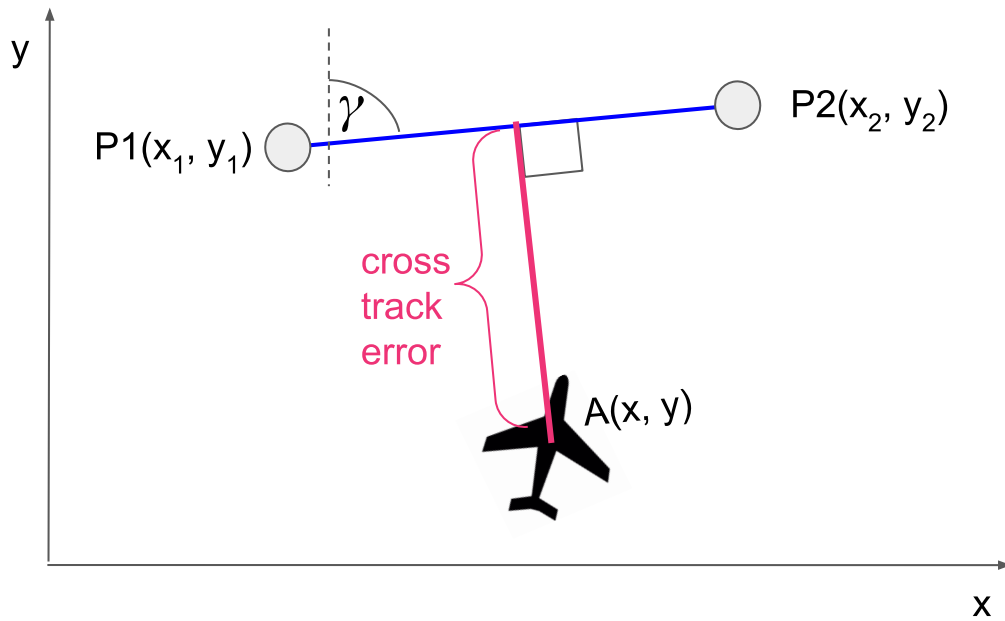
189 In (Xi and Liu 2020), the authors addressed the problem of path planning to a specific target for an UAV  
190 in dynamically changing environments with obstacles. They proposed a training scheme for reinforcement  
191 learning based on the Deep Deterministic Policy Gradient (DDPG) algorithm. The purpose of the scheme  
192 was to train the agent such that the UAV behaved in accordance with the real human intent. However, the  
193 presented solution was only based on a 2D environment and did not consider wind.

194 Bouhamed et al. (Bouhamed et al. 2020) used the technique of transfer learning to train an agent model  
195 with DDPG to lead an UAV to a static or moving target. As in the case of many other methods mentioned  
196 before, no wind was considered in the proposed solution.

197 Apart from works directed to aviation, two other works are shortly presented, which inspired the use of  
198 the cross track error in this work. Both following works applied DRL to the *path-following problem*, where  
199 a vehicle needs to reach and follow a *predefined* path. To measure if a vehicle is following the path, the  
200 cross track error is calculated. The cross track error is the normal from a vehicle to a specific path. The  
201 objective is to reduce the cross track error to zero. By doing this, the vehicle approaches the path and stays  
202 on it:

203 In (Martinsen and Lekkas 2018b) DRL was applied to the path-following problem for a mariner-class  
204 vessel. They showed that the trained agent was able to minimize the cross track error to a straight-line path  
205 also under the presence of ocean currents.

206 In (Havenstrøm, Rasheed, and San 2021), DRL was applied to the path-following problem for an  
207 underwater vehicle. Apart from the horizontal cross track error, they calculated the along-track error and  
208 an additional error to the path, which they called vertical track error. The additional vertical track error  
209 allowed the agent to follow the path also in the three-dimensional space.



**Figure 2.** The cross track error

## 4 FUNDAMENTALS

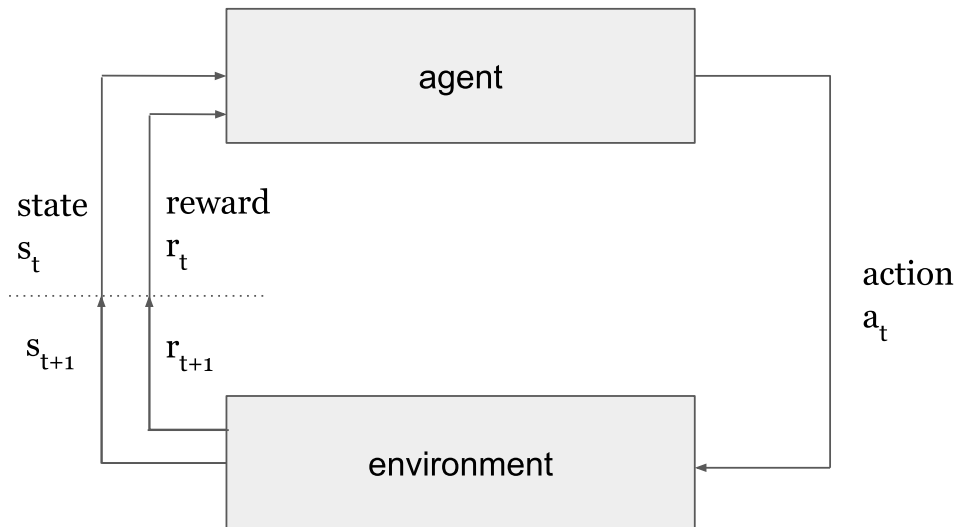
### 210 4.1 Flight Dynamics Model Aircraft

211 Flight Dynamics Model (FDM) is used for representing an aircraft in the simulated environments. The  
 212 (FDM) is the mathematical model that determines the physics of a flying aircraft in simulation. It is therefore  
 213 responsible for providing the aircraft's equations of motion and calculating the forces and moments acting  
 214 upon the aircraft. In this work, the widely used JSBSim library was applied as FDM with the included  
 215 Cessna127P aircraft (J Berndt 2004; Perry 2004; J Berndt and De Marco 2009; J S Berndt et al. 2011).  
 216 The method should work also when training the agent with other types of aircraft, since it is based on  
 217 model-free DRL.

### 218 4.2 Track Errors & Path Following

219 In this section, the principles of the *cross track error* and *vertical track error* are described, since they  
 220 are used to follow a path in the created environments. In a 3D Cartesian coordinate system, consider a  
 221 straight line path in the XY-plane between two points  $P1(x_1, y_1)$  and  $P2(x_2, y_2)$ , as can be seen in Figure  
 222 2. Furthermore, the angle  $\gamma$  of the path and the position of the aircraft at the point  $A(x, y)$  are given. The  
 223 cross track error  $e_{\text{cross}}$  is the normal from the aircraft to the straight line path and is calculated as shown in  
 224 Equation 1, as described in (Martinsen and Lekkas 2018a):

$$e_{\text{cross}} = -(x - x_1) \sin(\gamma) + (y - y_1) \cos(\gamma) \quad (1)$$



**Figure 3.** The interaction of the agent with the environment

225 In this work, the vertical track error is obtained by using a similar calculation as for the cross track error.  
 226 The only real difference is that the path lies in the YZ-plane instead:

$$e_{\text{vertical}} = -(y - y_1) \sin(\gamma) + (z - z_1) \cos(\gamma) \quad (2)$$

227 The angle  $\gamma$  is always known in this work (the runway heading or approach slope), therefore there is no  
 228 need to calculate it from the two points P1 and P2 first.

### 229 4.3 Reinforcement Learning

230 Reinforcement learning allows automating decision-making by maximizing a reward signal, a simple  
 231 scalar (Sutton and Barto 2018, p. 1). This is accomplished by a decision maker, called *agent*. The agent  
 232 learns a specific behavior from experience, which is generated by interacting with an *environment*.

233 A problem in RL is often stated formally as finite horizon Markov Decision Process (MDP), which  
 234 consists of the tuple  $(S, A, R, P)$ . At each time step  $t$ , the agent interacts with the environment by taking  
 235 an *action*  $a_t \in A$ , where  $A$  is a set of actions. After doing so, the environment transitions from the  
 236 current state  $s_t \in S$  into a next state  $s_{t+1} \in S$ , where  $S$  is a set of states. The transition probability  
 237 function  $P(s_{t+1}|s_t, a_t)$  represents the probability for that transition, given the current state  $s_t$  and the  
 238 agent's chosen action  $a_t$ . Subsequently, the agent receives a reward  $r_{t+1}$  produced by the reward function  
 239  $R(s_t, a_t, s_{t+1})$ . Furthermore, the agent receives the new state  $s_{t+1}$ . The state  $s_{t+1}$  becomes the current  
 240 state  $s_t$  and the reward  $r_{t+1}$  the current reward  $r_t$ . Based on those received values, once more, the  
 241 agent selects the next action to interact with the environment in the next time step  $t + 1$ . In the finite  
 242 MDP, the environment terminates, when a terminal state is reached. There, the time horizon from  $t_0$   
 243 to the last time step is called an *episode*. Tasks with episodes are called *episodic tasks*. The sequence  
 244  $(s_0, a_0, r_0), (s_1, a_1, r_1), \dots, (s_n, a_n, r_n)$  of an episode is called *trajectory*. Figure 3 demonstrates the typical  
 245 control loop of RL.

246 The *goal* of the agent is to maximize the received reward. However, the agent tries not to maximize its  
 247 immediate reward, but the cumulative long term reward, called *expected return*  $G_t$ . The agent achieves  
 248 this by taking the actions for a given state that lead to the highest cumulative reward. The function which



249 produces the actions of the agent for a given state is called *policy*  $\pi(a|s)$ . If the agent follows the policy  $\pi$   
 250 at time step  $t$ , then  $\pi(a|s)$  represents the probability for the agent to select action  $a$  given state  $s$ .

$$G_t \doteq r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \quad (3)$$

251 where  $\gamma \in [0, 1]$  is the discount rate. The discount rate influences how the agent values future rewards.

## 5 METHODS

### 252 5.1 Custom Environments

253 To apply Deep Reinforcement Learning to the problem of emergency flight guidance, three custom  
 254 environments were created:

- 255 •**Default**: The `Default` environment contains no wind. Moreover, the action space has one dimension.
- 256 •**Wind**: The `Wind` environment is an extension of the `Default` environment. It is extended by adding  
 257 constant wind during the training of an agent. The intensity of the wind is increased over the training  
 258 time. As in the `Default` environment, the action space has one dimension.
- 259 •**TwoActions**: A third environment was created where the action space consists of two dimensions.  
 260 There is no wind during the training in the `TwoActions` environment.

261 Our software `guidance_flight_env`, which provides the custom environments, started as a git fork  
 262 of Gordon Rennie’s open-source software *Gym-JSBSim* (Rennie 2018), but was adapted for flight guidance.

### 263 5.2 Main Components of the Custom Environments

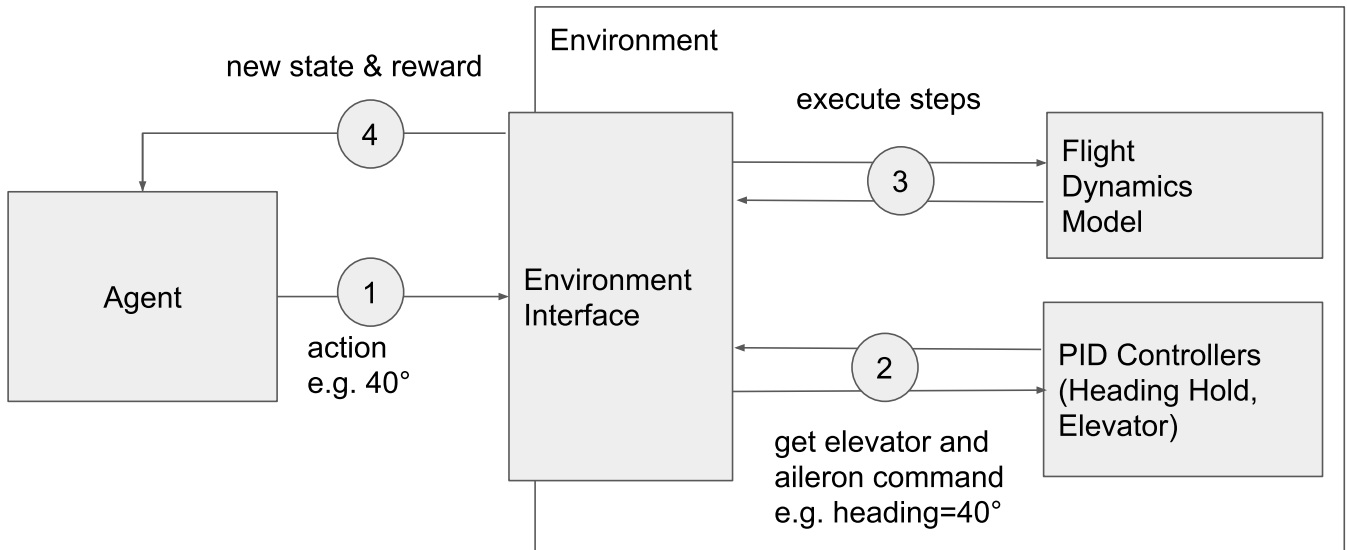
264 In this section, the most important components of the custom environments and the interaction with the  
 265 agent are described. The interaction is illustrated in Figure 4.

266 As reported above in Section 4, the agent interacts with the *environment* by selecting an action  $a_t$  from  
 267 the environment’s action space  $A$  at each time step  $t$ . For the `Default` environment and the `Wind`  
 268 environment the action  $a_t$  is a simple angle value, the desired heading, between  $0^\circ$  and  $359^\circ$ . After the  
 269 agent has selected the heading, it is passed to the environment, where it is then internally processed by a  
 270 heading hold PID controller. The controller, returns the aileron command necessary for pointing the aircraft  
 271 to the desired angle. The aileron command is set in the Flight Dynamics Model (FDM) component, and  
 272 subsequently a number of simulation steps are executed to update the FDM. Next, a new state is obtained,  
 273 and the corresponding reward is calculated. Both are then passed back to the agent. In addition to the  
 274 heading hold PID controller, a second controller, more specifically, a pitch angle hold controller, is used to  
 275 keep also the aircraft’s pitch angle in a fixed position of  $0^\circ$ . The reason for this was shortly described in  
 276 Section 2.

### 277 5.3 Overview of Core Concept

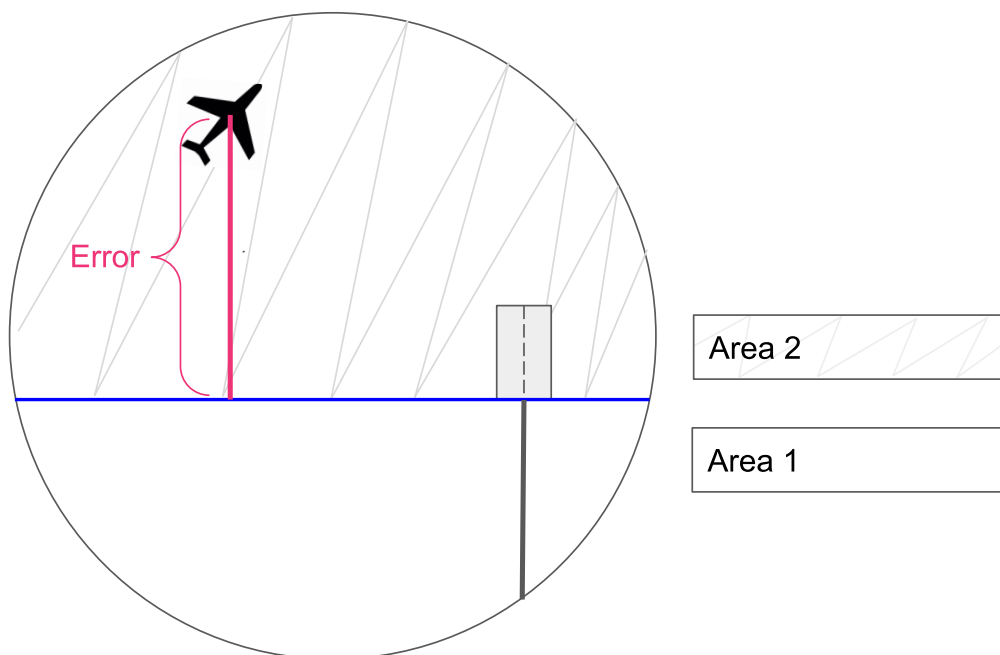
278 In this section, the core concept which is common in all three environments is shortly presented. Details  
 279 are then described in the proceeding sections.

280 A line, perpendicular to the extended center line of the runway, divides the XY-plane of the Cartesian  
 281 coordinate system into two different areas: *Area 1* and *Area 2*. In the case of the aircraft’s location being in



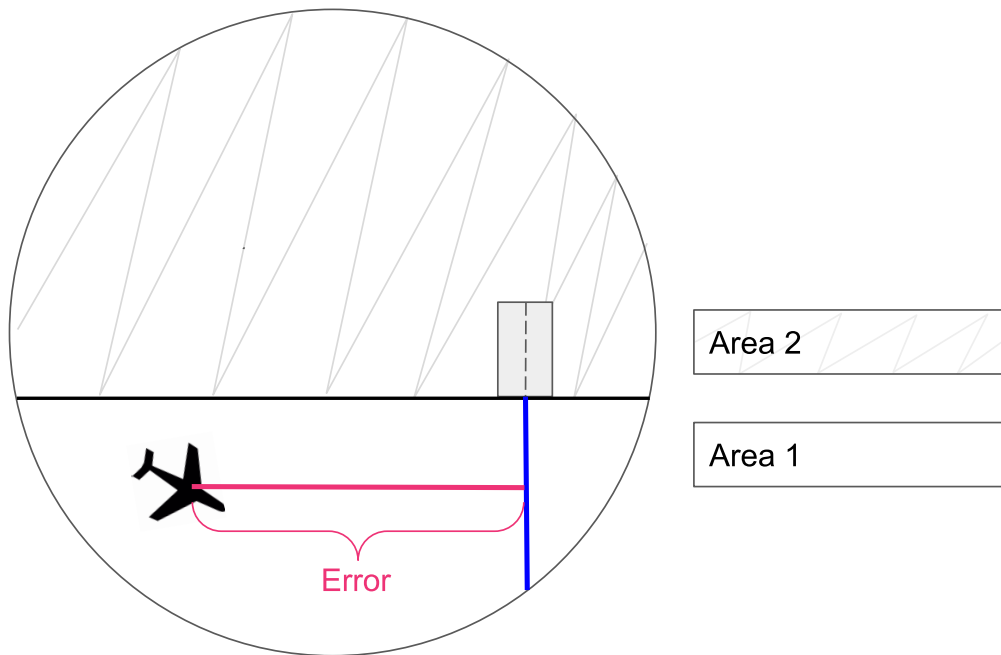
**Figure 4.** The interaction of the agent with the custom environments

282 *Area 2*, the cross track error  $e_{\text{cross}}$  of the aircraft to the perpendicular line is calculated and is included in  
 283 the environment's state at time step  $t$ . The agent then observes this state and, while trying to maximize  
 284 the reward by reducing the error, it guides the aircraft to *Area 1*. The described procedure is illustrated in  
 285 Figure 5.



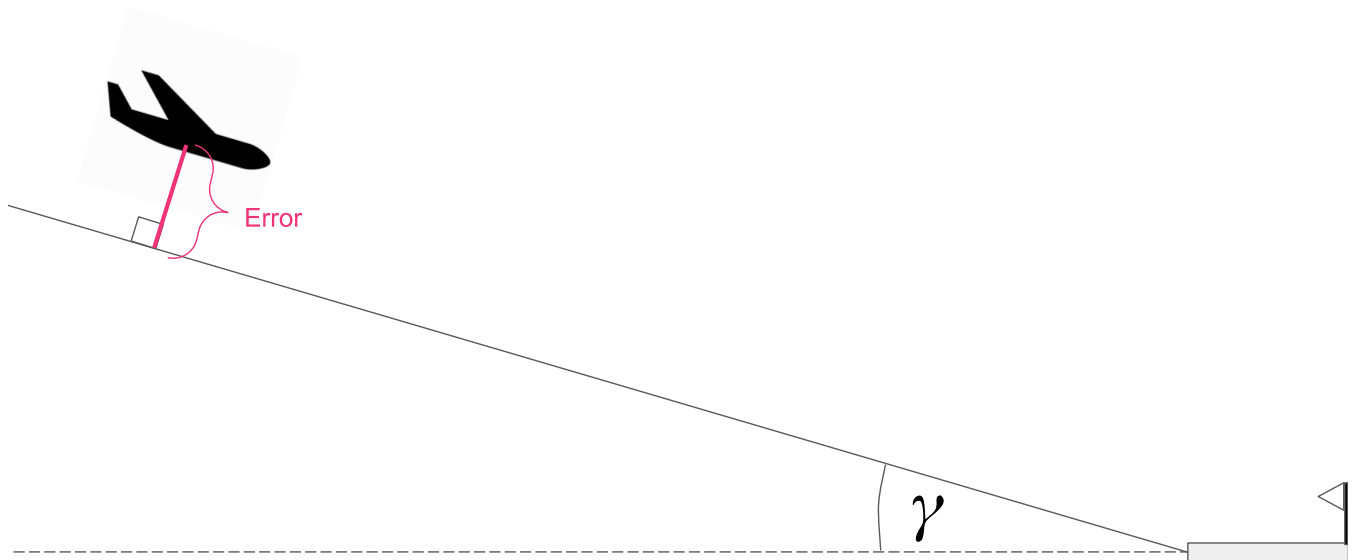
**Figure 5.** The aircraft is in *Area 2*. The cross track error to the perpendicular line (blue) to the extended center line of the runway (black) is calculated.

286 Once the aircraft is in *Area 1*, the cross track error to the extended center line of the runway is added to  
 287 the state instead. Now, the agent is incentivized to guide the aircraft to the target position. See Figure 6



**Figure 6.** The aircraft is in *Area 1*. The cross track error to the extended center line of the runway (blue) is calculated and added to the state.

288 Similar to the cross track error in the XY-plane, the error  $e_{\text{vertical}}$  in the YZ-plane is calculated and  
 289 added to the state. The error  $e_{\text{vertical}}$  is obtained by calculating the vertical distance of the aircraft to an  
 290 imagined line created by a specific approach slope  $\gamma$  (Havenstrøm, Rasheed, and San 2021). An example is  
 291 demonstrated in Figure 7



**Figure 7.** Here, the  $e_{\text{vertical}}$  is calculated from the aircraft's position in the YZ-plane to the line with the slope angle

292 The agent receives positive and negative rewards depending on the magnitude of both errors. The rewards  
 293 should motivate the agent to guide the aircraft closer to the target, while maintaining the correct heading and

294 approach slope angle. When the aircraft reaches the target, the episode terminates. It does also terminate  
 295 when the agent has no remaining altitude left before achieving the target.

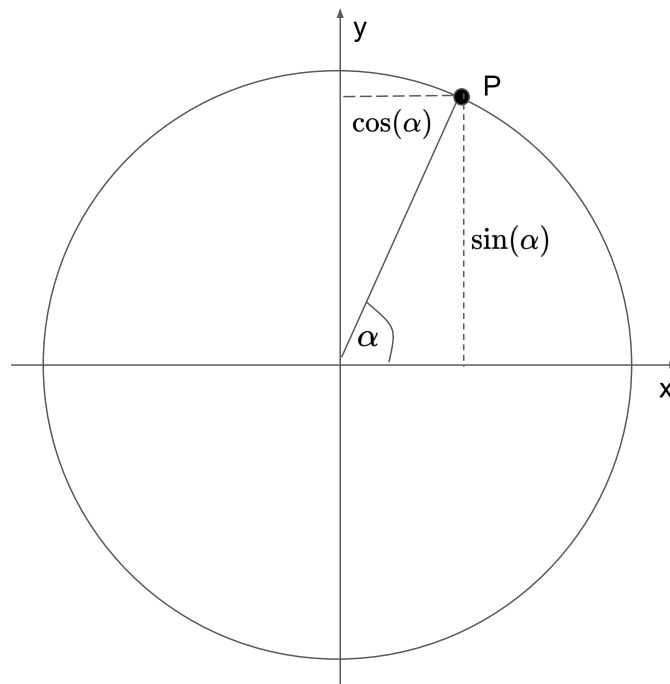
#### 296 5.4 Angles and Discontinuity: TwoActions-Environment

297 For humans, it is mostly clear that an angle of  $359^\circ$  is relatively close to an angle of  $1^\circ$ . However, for a  
 298 neural network this relation seem to be hard to learn, since discontinuity occurs between  $2\pi$  ( $360^\circ$ ) and 0  
 299 (Zhou et al. 2019); numerically the values 360 and 0 are very far apart. Therefore, the following assumption  
 300 is made: Using angle values in the action and observation space could introduce noise during the training  
 301 process, which would lead to worse performance of the RL agent.

302 To avoid this noise, some researchers, (Berner et al. 2019, Xi and Liu 2020), encode an orientation angle  
 303 value  $\alpha$  as shown in Equation 4.

$$\alpha \rightarrow (\sin(\alpha), \cos(\alpha)), \text{ where } \alpha \in [0, 2\pi] \quad (4)$$

304 They then used the tuple of  $(\sin(\alpha), \cos(\alpha))$  in the observation- or action space to represent one angle.  
 305 See Figure 8. This implies that the neural network only received the tuple, instead of simply a single angle  
 306 value.



**Figure 8.** An angle  $\alpha$  can be represented by the two values:  $\sin(\alpha)$  and  $\cos(\alpha)$

307 To test the above assumptions, the `TwoActions` environment was created and there all angles in the  
 308 state and action space were encoded as in Equation 5. To obtain the actual angle again, both values were  
 309 decoded in the environment using the `arctan2` function as shown in Equation 5

$$(\sin(\alpha), \cos(\alpha)) \rightarrow \arctan2(\sin(\alpha), \cos(\alpha)), \text{ where } \alpha \in [0, 2\pi] \quad (5)$$

---

310 The `Default` environment, on the other hand, uses the single angle values directly. Later, in Section 6  
311 the two environments are compared by demonstrating the performance of the trained agent in both.

## 312 5.5 Initial Setup

313 At the beginning of each episode, the following steps are executed:

314 **1. Cartesian Coordinates & Aircraft's Position:** The altitude and geographical coordinates (latitude,  
315 longitude) of the aircraft are converted into coordinates of the 3D Cartesian coordinate system ( $x, y, z$ ).  
316 This conversion simplifies calculations, since now simple vector algebra can be used. The conversion  
317 is justified, since the resulting conversion errors are kept to a minimum when distances between objects  
318 are relatively low, which is mostly true in the case of an engine's failure. After the conversion, the origin  
319 of the XY-plane ( $x=0, y=0$ ) represents the aircraft's initial  $x, y$  coordinates. Next, the aircraft's altitude ( $z$   
320 coordinate) is chosen arbitrary between a minimum and maximum value. Now, the target position and all  
321 future positions of the aircraft are calculated relative to the origin. Using relative positions, instead of  
322 absolute positions, improves the generalization capacities of the agent's neural network (ibid.).

323 **2. Aircraft Heading:** The heading of the aircraft is chosen arbitrary in the range of  $[0, 2\pi]$ .

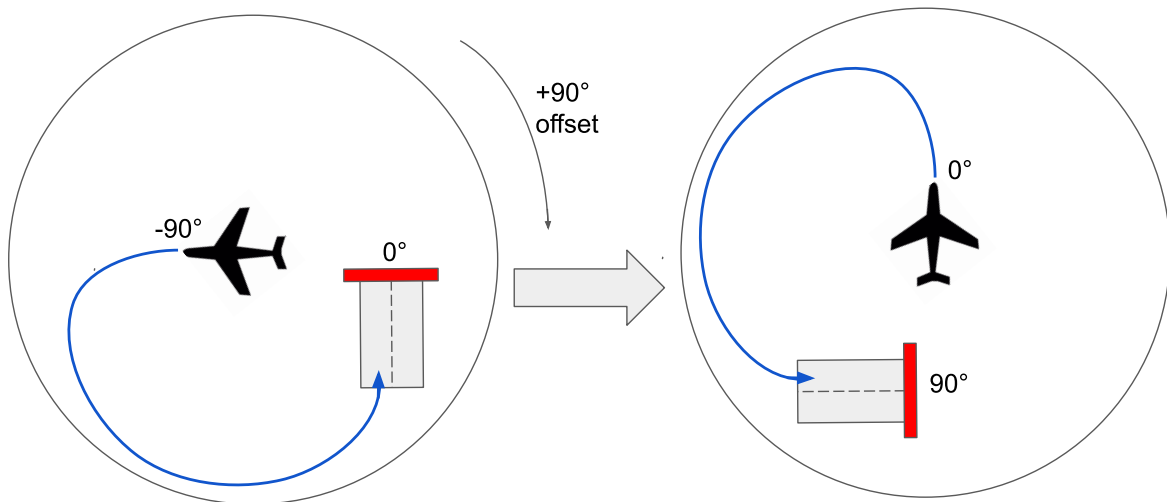
324 **3. target position:** After setting the origin to the aircraft's initial position, the position of target position is  
325 chosen arbitrary in a specified radius from the origin.

326 **4. Landing Field's Heading:** The heading of the landing field is kept at  $0^\circ$  for all episodes during the  
327 training. This simplification reduces training time, since the agent does not need to learn to adapt to  
328 different targets. However, this simplification does not restrict the capabilities of the agent. A simple  
329 offset to all angles can be added when using the agent after the training, as can be seen in the example  
330 in Figure 9. The value of the offset is the real heading of the landing field. Since all distances are not  
331 absolute but relative, after adding the offset everything works as it would for the case of the  $0^\circ$  headings  
332 during the training.

333 **5. Approach Slope:** The approach slope of the target landing field is set to a fixed value of  $4^\circ$  for all  
334 episodes, since approach slope angles of most runways are around  $3^\circ$ - $5.5^\circ$  Allerton 2009, p.193.

335 **6. Engines Off:** The simulated aircraft's engines are turned off to simulate an emergency setting.

336 **7. Curriculum Learning:** The technique of curriculum learning (Florensa et al. 2017) is used to train the  
337 agent. Curriculum learning was incorporated into the environment by dividing the training into five phases.  
338 When the agent's mean reward surpasses a certain threshold, then the phase number is incremented,  
339 which in turn increases the environment's difficulty. The difficulty is specifically increased by extending  
340 the radius of the circle, where the target positions are arbitrarily generated. In addition, in the `Wind`  
341 environment, the wind speed is changed in each phase.



**Figure 9.** An example of a landing field with a real heading of  $90^\circ$ : After the training, adding the offset (real heading of the landing field) of  $90^\circ$  to all angles in the environment, allows the agent to guide the aircraft to the landing field.

## 342 5.6 Action Space

343 To control the flight path of the aircraft, the agent needs to be able to select the desired heading  $\psi$  of the  
 344 aircraft. The aircraft's pitch angle is fixed at  $0^\circ$  and is therefore not part of an action.

345 In the case of the default and `Wind` environment, the action space has one dimension. In every time step  
 346  $t$ , the agent selects a single action  $a$ , which represents the angle value for the heading. The action could be  
 347 a value ranging from  $0 - 2\pi$ . However, it is good practice to rescale actions to be in the range of  $[-1, 1]$   
 348 (Raffin et al. 2019). The agent's selected value is then rescaled again back into the range of  $[0, 2\pi]$  inside  
 349 the environment.

350 In Section 5.4, we've made the assumption that to avoid problems with discontinuity, the desired angle  $\psi$   
 351 should be encoded using the function in Equation 4. Therefore, to the test the action  $a$  in the `TwoActions`  
 352 environment, is formulated as a tuple of two values  $a = [x, y]$ , where  $x, y \in [-1, 1]$ . Inside the environment,  
 353 the actual angle  $\psi$  is then obtained by  $\psi = \arctan2(y, x)$

## 354 5.7 State Space

### 355 5.8 State Space of the Default environment

356 In the case of the `Default` environment, the state consists of the following 13 elements, shown in Table  
 357 1.

358 The state elements are described in detail in the following subsections.

State Variable	Symbol	Unit
cross track error	$e_{\text{cross}}$	km
vertical track error	$e_{\text{vertical}}$	km
in area	$a$	-
difference x	$x_{\text{diff}}$	km
difference y	$y_{\text{diff}}$	km
difference z	$z_{\text{diff}}$	km
distance to target	$d_{3D}$	km
remaining altitude	$h$	ft
descent rate	$u$	ft / s
turn rate	$tr$	rad / s
true airspeed	$v_{tas}$	ft / s
true heading	$\psi$	rad
heading error	$e_{\psi}$	rad

**Table 1.** The state space of the `Default` environment

## 359 5.9 State Space of the Wind Environment

360 The state space of the `Wind` environment includes all elements from Table 1. Furthermore, the following  
361 elements from Table 2 are added.

State Variable	Symbol	Unit
total wind north	$w_n$	ft / s
total wind east	$w_e$	ft / s
drift	$\lambda$	rad

**Table 2.** The additional state elements of the `Wind` environment

362 A total constant wind originating in the East or West is added to the state of the `Wind` environment.  
363 Furthermore, the drift angle  $\lambda$  is added. The drift angle is the difference between the aircraft's true heading  
364 and the track angle. The track angle is the angle of the lateral track the aircraft actually flies over the ground.  
365 Under the occurrence of wind, the track angle and aircraft's true heading differ by a relatively high drift  
366 value  $\lambda$ , whereas in the case of no wind they should be equal. Figure 10 demonstrates this relationship.

367 The total wind values and drift, allow the agent to adjust the heading to react to the wind accordingly.

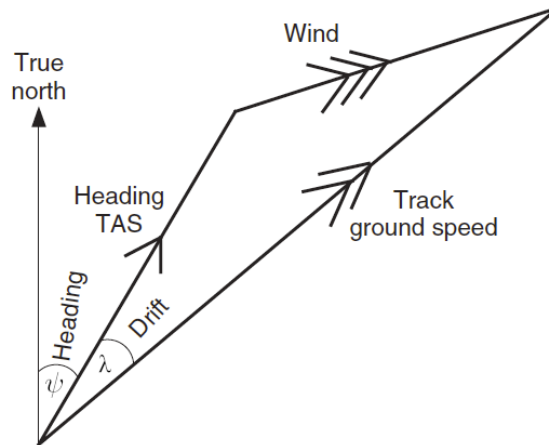
## 368 5.10 State Space of the TwoActions environment

369 The only difference in the `TwoActions` environment compared to the `Default` environment is that all  
370 the angles in the state are represented with two values ( $\sin(\alpha)$ ,  $\cos(\alpha)$ ) as described in Section 5.4. For this  
371 reason, the true heading and the heading error are replaced by two values each. This leads to four additional  
372 elements in the state space of the `TwoActions` environment, raising the number of state elements to 17.

## 373 5.11 Terminal States

374 A terminal state is reached when one of the following is true:

- 375 •The aircraft's remaining altitude is lower than zero. This happens when the aircraft's z-coordinate value  
376 is lower or equal to the target's z-coordinate value.



**Figure 10.** The aircraft's real track angle over the ground differs to its heading by the drift angle  $\lambda$ . (Figure adapted from Allerton 2009, p.250.)

- 377 •The aircraft reaches the target position. This is true, when the Euclidean distance from the aircraft's  
 378 position to the target position is lower than the threshold  $\theta_{\text{target}} = 0.1$  km. The value 0.1 was selected  
 379 during preliminary experiments. A lower value for  $\theta_{\text{target}}$  could have been selected, but it would have  
 380 increased training times. However, except for longer training times, it is expected that the performance  
 381 of the trained agent would not decline. Nonetheless, using lower values for  $\theta_{\text{target}}$  should be studied in  
 382 future research.
- 383 •The remaining time  $t$  is equal to zero. It should be noted that, in all experiments of this work, the altitude is  
 384 chosen such that the engines-off aircraft reaches the ground before  $t$  reaches the value of zero. Therefore,  
 385 in this work, the remaining time  $t$  was mostly ignored. However, it could be considered in future research,  
 386 if an aircraft with temporary working engines is investigated or if a higher remaining altitude is chosen.

387 When one of the terminal states is reached, the episode ends and the agent is rewarded with sparse  
 388 rewards as described in Section 5.12

## 389 5.12 Reward

390 A reward function was developed to motivate the agent to behave in the following way: Leading the  
 391 aircraft to the target position in the 3D space. The target position should be reached with the expected  
 392 heading. The reward function was constructed by combining sparse and dense rewards:  $r = r_{\text{sparse}} + r_{\text{dense}}$ .

393 In (Henderson et al. 2018), Henderson et al. showed that reward clipping or reward scaling effected  
 394 learning of the agent significantly, when neural networks and gradient-based methods are used (which is  
 395 the case in most state of the art DRL-algorithms). Therefore, all rewards in the three custom environments  
 396 are kept in the range of  $[-10, 10]$ .

### 397 5.12.1 Sparse Rewards

398 In this section, the different sparse rewards, which are returned on the terminal states, are presented in  
 399 detail.

400 **Terminal State at Target:** When the aircraft reaches the target from the wrong side, coming from *Area*  
 401 2, then the reward is -10. However, when the aircraft reaches the target from the correct side, coming from



402 *Area 1*, then a positive reward of 9 is given. Furthermore, the agent is rewarded for the correct heading.  
 403 This is done in the following way: The heading error  $e_h$  in the range of  $[x_{min}, x_{max}]$  is rescaled into the  
 404 range of  $[0, 1]$ , where  $x_{max}$  represent the threshold (the tolerance value) for the correct heading. The value  
 405  $x_{min}$  is set to 0 and  $x_{max}$  to 10 (degrees). The scaling results in  $e_{h_{scaled}}$

406 Then, the bonus reward is calculated in the following way:

$$r_b = 1 - e_{h_{scaled}}, \text{ where } r_b \in [0, 1] \quad (6)$$

407 The final reward at the target terminal state is:

$$r_{target} = \begin{cases} 9 + r_b & \text{if } area = 1 \\ -10 & \text{else} \end{cases} \quad (7)$$

408 **Remaining Altitude is Zero:** When the remaining altitude is lower than zero, then the current episode  
 409 is terminated. The Euclidean distance of the aircraft's position  $(x_1, y_1, z_1)$  to the target  $(x_2, y_2, z_2)$  is  
 410 calculated and multiplied by a constant. The negative reward  $r_{altitude}$  is calculated as in Equation 8:

$$r_{altitude} = -6 * \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2 + (z_2 - z_1)^2} \quad (8)$$

411 **Time is Zero:** When the time step  $t$  is equal to zero, the current episode is terminated. The reward  $r_{time}$   
 412 is calculated, as in the case of remaining altitude equal to zero. (See Equation 8).

### 413 5.12.2 Dense rewards

414 In addition to sparse rewards, dense rewards were incorporated into all the environments to additionally  
 415 encourage the agent to behave as expected. The agent receives a dense reward on each time step  $t$ . By  
 416 adding dense rewards, the training time is reduced significantly.

417 **Reach Track:** As described in Section 5.3, depending on the aircraft being in *Area 1* or *Area 2*, different  
 418 track errors are calculated:

419 •**Aircraft is in Area 2:** The cross track error  $e_{cross}$  is calculated as described in Equation 1 (Section 4) to  
 420 the line perpendicular to the centerline of the landing field. The vertical track error  $e_{vertical}$  in *Area 2* is set  
 421 to zero, which is done to motivate the agent to set focus on leading the aircraft to *Area 1* first.

422 •**Aircraft is in Area 1:** Only now, in *Area 1*, the aircraft needs to position the aircraft for the final approach,  
 423 where it is important to reduce the vertical track error. Therefore, the vertical track error  $e_{vertical}$  is  
 424 calculated as in Equation 2 (Section 4). Moreover, the cross track error  $e_{cross}$  is calculated relative to the  
 425 extended center line of the landing field.

426 Next,  $e_{to\_track}$  is calculated as in:

$$e_{to\_track} = \begin{cases} e_{cross} + e_{vertical} & \text{if } area = 1 \\ e_{cross} & \text{else} \end{cases} \quad (9)$$

427 The negative reward  $r_{to\_track}$  is now derived from  $e_{to\_track}$  by applying the exponential function to the  
 428 track as in:

$$r_{\text{to\_track}} = -\exp(e_{\text{to\_track}}) \quad (10)$$

429 By doing so, the negative reward increases exponentially for a larger  $r_{\text{to\_track}}$ . This should additionally  
 430 motivate the agent to reach the track without many deviations. Finally, the resulting reward value  $r_{\text{to\_track}}$  is  
 431 then scaled to the range  $[0, 1]$ .

432 In addition, to  $r_{\text{to\_track}}$ , in each time step, the agent receives a constant negative reward  $r_{\text{area2}} = -2$  when  
 433 the aircraft is in *Area 2*, else  $r_{\text{area2}} = 0$ . Therefore, the agent receives more negative rewards when the  
 434 aircraft is in *Area 2* compared to when it is in *Area 1*. Whereas,  $r_{\text{to\_track}}$  acts as a helping leading signal for  
 435 the agent to guide the aircraft to the track, the constant  $r_{\text{area2}}$  adds an urgency to reach the *Area 1* as soon  
 436 as possible.

437 **Keep aircraft on-track** Once the agent guided the aircraft to reach the desired track, the agent then  
 438 receives additional rewards to continue keeping the aircraft on that track. In this work, the aircraft is  
 439 *on-track* when the following three conditions hold:

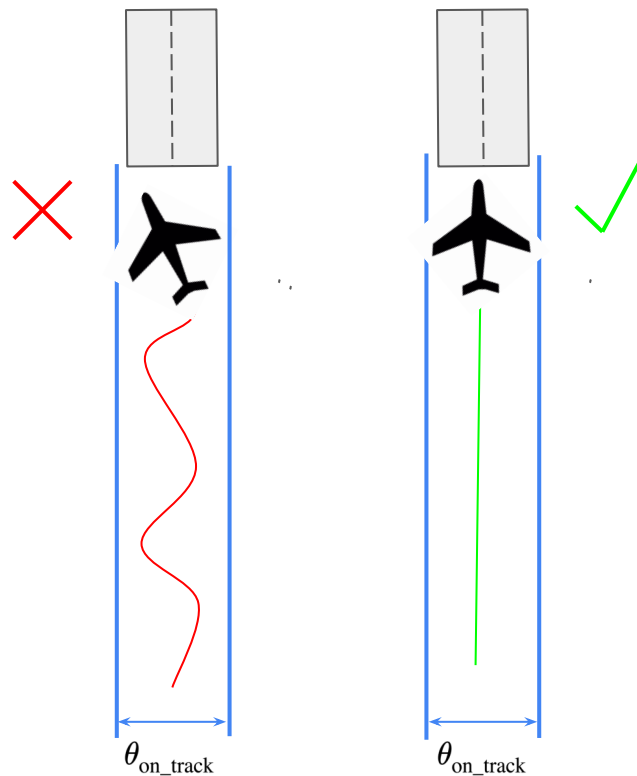
- 440 1. The medium track error  $e_{t\_m}$  is lower than the threshold  $\theta_{\text{on\_track}} = 0.1$ . The medium track error  $e_{t\_m}$  is  
 441 calculated as in  $e_{t\_m} = (|e_{t\_last}| + |e_t|)/2$ , where  $e_{t\_last}$  is the track error of the previous time step  $t - 1$   
 442 and  $e_t$  the track error of the current time step  $t$ .
- 443 2. The Euclidean distance in time step  $t$  is lower than the distance in the previous time step  $t - 1$ . In both  
 444 time steps, the Euclidean distance is calculated from the aircraft to the target position. This condition is  
 445 important, since it ensures that the distance of the aircraft to the target is reducing in every time step.
- 446 3. The heading error  $e_h$  is lower than  $90^\circ$ . The reason is to only allow the agent to receive a positive reward  
 447 if the aircraft reaches the target position from the correct side (from *Area 1*).

448 Therefore, when the aircraft is *on-track*, the agent receives a positive reward. However, to avoid positive  
 449 reward cycles (or reward hacking) (Randløv and Alstrøm 1998; Ng, Harada, and Russell 1999), it is  
 450 necessary to add additional negative rewards when the above *on-track* conditions do not hold anymore (i.e.,  
 451 when the aircraft leaves the track). For this reason, the agent receives the reward of  $r_{\text{on\_track}} = 1$  if the  
 452 aircraft is on-track, else the agent receives a negative reward in the form of  $r_{\text{on\_track}} = -e_{\text{track\_diff}}$ , where  
 453 the value of  $e_{\text{track\_diff}}$  is the difference of the track error  $e_{\text{track}}$  at time step  $t - 1$  and the current at time step  
 454  $t$ . Hence, the agent is punished for the aircraft moving away from the track and rewarded for staying on it.  
 455 This is presented in the following Equation 11.

$$r_{\text{on\_track}} = \begin{cases} 1 & \text{if aircraft is on-track} \\ -e_{\text{track\_diff}} & \text{else} \end{cases} \quad (11)$$

### 456 5.12.3 Heading

457 When the aircraft is *on-track*, it could still occasionally oscillate horizontally in the XY-plane, for a short  
 458 period of time. The oscillating is normally not a problem if the aircraft is still relatively distant. However,  
 459 when close to the landing field it's important that the aircraft does not oscillate, as is illustrated in Figure 11



**Figure 11.** Given the threshold  $\theta_{\text{on\_track}}$ , the aircraft should fly in a straight line when on-track (right side), to avoid reaching the landing field with a wrong heading (left side, intentionally shown exaggerated). To encourage the agent to achieve this, an additional reward is given.

460 Avoiding oscillation is necessary, to guarantee that when reaching the target position, the heading of the  
 461 aircraft is as expected. Therefore, too additionally minimize the heading error, when *on-track*, the agent  
 462 should be encouraged to keep the heading fixed on the desired direction.

463 A first idea, to avoid oscillation, could be the following: Reducing the threshold  $\theta_{\text{on\_track}} = 0.1$  described  
 464 in the first condition in 5.12.2 even further. However, this stricter condition would increase the difficulty for  
 465 the agent to receive any positive reward ( $r_{\text{on\_track}} = 1$ ) at all. The increased difficulty could lead to longer  
 466 training times, or even to the issue that the agent does never learn to stay on-track.

467 For this reason, the alternative solution is to add an additional reward for the correct (or wrong) heading  
 468 instead. The agent obtains the additional positive or negative reward  $r_{\text{heading}} \in [-0.5, 0.5]$  only when the  
 469 aircraft is on-track. Since the  $r_{\text{heading}}$  is in  $[-0.5, 0.5]$  and  $r_{\text{on\_track}} = 1$ , the agent would, in the worst case  
 470 (heading error around  $90^\circ$ ), still receive a reward of 0.5. Therefore, the additional dense heading reward can  
 471 be seen as a kind of bonus reward for the agent. Nonetheless, the agent trying to maximize the reward in  
 472 general will try to maximize both  $r_{\text{on\_track}}$  and  $r_{\text{heading}}$ , and therefore reaching the target with the expected  
 473 heading, by avoiding oscillating.

474 The reward  $r_{\text{heading}}$  is obtained as follows. First, the difference  $e_{\text{heading\_diff}} \in [0, \frac{\pi}{2}]$  of the heading error  
 475 in time step  $t$  and the heading error in time step  $t - 1$  is calculated. Next, if the difference in time step  $t$  is  
 476 smaller than the difference in time step  $t - 1$ , then the reward  $r_{\text{heading}}$  with a positive value is returned, else  
 477 a negative a shown in the following Equation 12:

$$r_{\text{heading}} = \begin{cases} \frac{e_{\text{heading\_diff}}}{\pi} & \text{if aircraft is on track} \\ -\frac{e_{\text{heading\_diff}}}{\pi} & \text{else} \end{cases} \quad (12)$$

478 Having an additional negative reward avoids positive reward cycles (Randløv and Alstrøm 1998; Ng,  
479 Harada, and Russell 1999) as described previously in Section 5.12.2.

### 480 5.13 Final Combined Reward Function

481 To obtain the final combined reward function, as shown in the following Equation 13:

$$r = r_{\text{sparse}} + r_{\text{dense}} \quad (13)$$

482 it is therefore needed to add all sparse rewards of the preceding sections, as in Equation 14:

$$r_{\text{sparse}} = r_{\text{target}} + r_{\text{altitude}} + r_{\text{time}} \quad (14)$$

483 Whereas the dense rewards of the preceding sections are combined, as in Equation 15:

$$r_{\text{dense}} = r_{\text{to\_track}} + r_{\text{on\_track}} + r_{\text{area2}} + r_{\text{heading}} \quad (15)$$

## 6 EXPERIMENTS

484 As described in Section 5.1 three custom environments were developed and implemented. Experiments  
485 were conducted in all three environments. In the following sections, first the setup for the experiments are  
486 described. Then the results of those experiments are presented and interpreted.

487 The three experiments have the following setup in common:

488 •**Training in the cloud:** For the training phase, two 2,25 GHz CPUs (AMD EPYC 7B12) and a GPU  
489 (NVIDIA Tesla K80 24 GB) were used.

490 •**Algorithm:** RLlib’s Twin Delayed DDPG (TD3) was used as the agent’s algorithm, since it is a current  
491 state-of-the-art algorithm and a successor of the well-studied Deep Deterministic Policy Gradient (DDPG)  
492 algorithm Lillicrap et al. 2015. Furthermore, TD3 requires less hyperparameter tuning in comparison to  
493 other algorithms like DDPG. The same hyperparameters of the TD3 were used for all experiments. The  
494 hyperparameters of the algorithm are presented in the Appendix 1.

495 •**Environments:** The custom environments described in Chapter 5.1 were used for training and evaluating  
496 the agent.

497 •**Evaluation:** After the agent was trained in one of the environments, the agent was evaluated in the  
498 same environment. For each experiment or sub-experiment, the evaluation was run thrice, each for 100  
499 episodes. For each run, a different seed was used. The results of the three runs were then averaged.

500 •**Aircraft:** Since most fatal accidents in aviation occur when small aircraft are used (*Air safety statistics in*  
501 *the EU - Statistics Explained* 2020; Dorr 2018), the light aircraft single-engine Cessna172P was selected  
502 in the *JSBSim* FDM. Before starting the experiment, the aircraft’s engine is turned off in the *JSBSim*  
503 FDM. Furthermore, throttle and mixture were set to zero. Those steps are needed to simulate broken  
504 engines and the emergency situation.

505 •**Initial Altitude:** The initial altitude at the beginning of each episode is selected arbitrarily in the range of  
506 2500-3500 feet. The reasons for this altitude range are the following: Firstly, at this altitude range, only  
507 a short amount of time is left to reach the landing field. Therefore, it represents a very dangerous part  
508 of the remaining flight and simulates the real-world scenario accordingly. Secondly, a practical reason,  
509 given the relatively low initial altitude, the episodes during simulation are shorter. This led to decreased  
510 training times. Nonetheless, if needed in future research, higher numbers for the altitude range could be  
511 chosen. Except for a longer training phase, it is expected that there should be almost no difference in the  
512 agent's performance after training.

513 •**Curriculum Learning:** Depending on the average reward received by the agent, one of the five phases,  
514 described in Section 5.5, is selected. When entering the next phase, the target spawning range widens, as  
515 can be seen in Table 3.

Phase	Reward threshold	Target Spawn Range (km)
0	$-\infty$	0.5
1	-900	1
2	-750	1.5
3	-600	2
4	-300	2

**Table 3.** The radius of the circle, where the target positions are arbitrarily generated, is adapted for each phase

## 516 6.1 Setup Experiment 1: Default environment

517 The agent was trained for 10,000 episodes (around 10 Million steps). And in addition, for 15,000 episodes  
518 (around 15 Million steps).

## 519 6.2 Setup Experiment 2: TwoActions

520 The agent was trained for 10,000 episodes (around 10 Million steps) in the `TwoActions` environment.

## 521 6.3 Setup Experiment 3: Wind Environment

522 In the `Wind` environment, during the training, the wind speed increases in each phase. See Table 4. For  
523 reducing complexity, the constant wind is set to originate in the East or West only.

Phase	Reward threshold	Speed Range [West, East] (ft/min)
0	$-\infty$	0
1	-900	[-600, 600]
2	-750	[-1200, 1200]
3	-600	[-2100, 2100]
4	-300	[-3000, 3000]

**Table 4.** In the `Wind` environment, for each phase additionally the wind speed changes

524 The agent was trained for 10,000 episodes (around 10 Million steps).

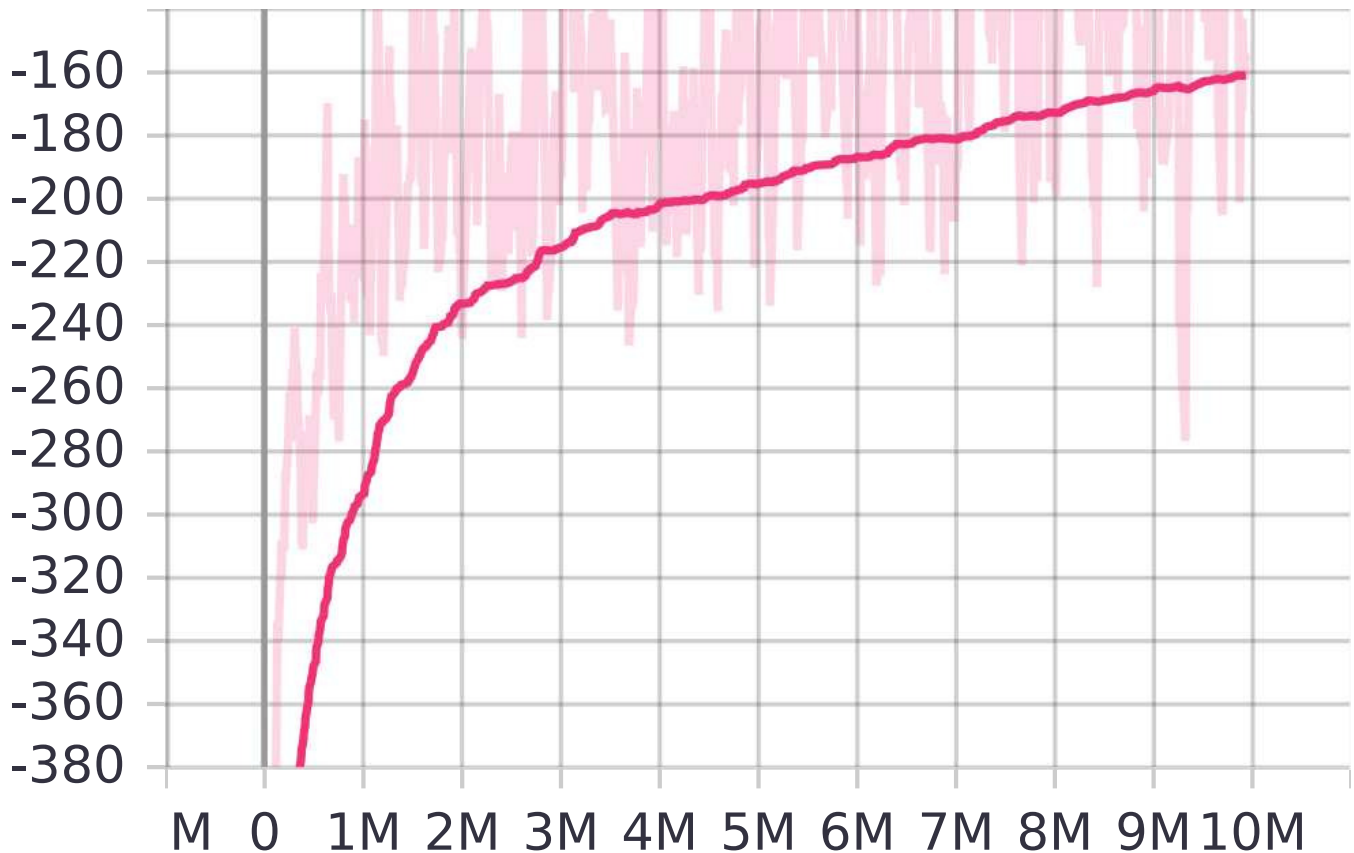
525 Later, the agent trained in the `Wind` environment was evaluated under constant wind originating in the  
526 East or West. Four sub-experiments were conducted, to test the performance of the agent under the four  
527 different wind speed ranges: (1) 600 ft/min (2) 1200 ft/min (3) 2100 ft/min (4) 3000 ft/min.

# 7 RESULTS

## 528 7.1 Experiment 1: Default environment

529 Figure 12 shows the average reward for 10 Million steps. The reward curve there might indicate that the  
530 agent did not conclude its learning process, since a plateau was not yet reached. Therefore, in a second  
531 sub-experiment, the agent was trained for an additional 5 Million steps, in total for 15 Million steps.  
532 However, even at 15 Million steps, the reward curve still seemed not to have reached a plateau, which  
533 implied that further improvements might have been possible. Therefore, probably even longer training  
534 periods would have been necessary. Yet, both agents were evaluated first, and the results of the evaluations

535 were already quite promising. Hence, the agent was not trained for longer than 15 Million steps, since the  
 536 hardware resources for training were relatively limited. For future research, training periods should be  
 537 prolonged.



**Figure 12.** The interaction of the agent with the custom environments

538 First, for the agent trained for 10 Million steps, the aircraft reached the target 62.3% of the time on  
 539 average. 69.3% of the time the agent landed on-track (reaching target included), where the agent landed  
 540 with only an average distance of 0.095 km away from the target position. Furthermore, the average heading  
 541 error was quite low at 4.65° (See Table 5).

<b>At-Target</b>	62.3 %
<b>On-Track (includes at-target)</b>	69.3%
<b>Distance to target (when on-track)</b>	0.095 km $\pm$ 0.086 km
<b>Heading error</b>	4.65° $\pm$ 7.63°
<b>Reward</b>	-122.37 $\pm$ 393.59

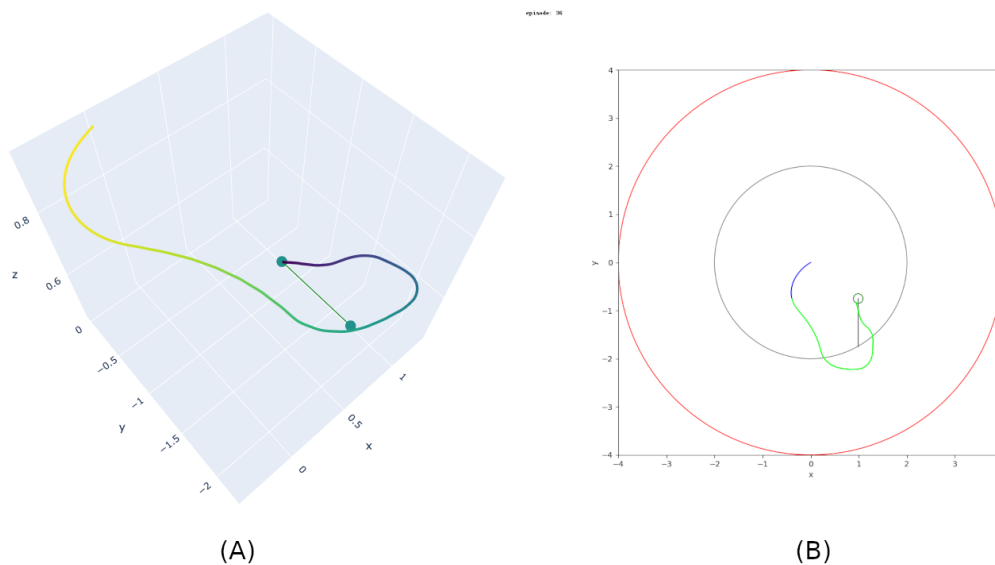
**Table 5.** Experiment 1: Evaluation of the agent trained in the `Default` environment for 10 Million steps. (All values are averages)

542 As expected, the agent trained for 15 Million steps, performed mostly better. This can be seen in Table  
 543 6. The agent reaches the target successfully with a high value of 73 % of the cases on average and ends  
 544 on-track (at-target included) even in 81.3% of the cases.

<b>At-Target</b>	73%
<b>On-Track (includes at-target)</b>	81.3%
<b>Distance to target (when on-track)</b>	0.1 km $\pm$ 0.1 km
<b>Heading error</b>	12.5° $\pm$ 9.4°
<b>Reward</b>	-133.93 $\pm$ 369.357

**Table 6.** Experiment 1 (sub-experiment): Evaluation of the agent trained in the `Default` environment for 15 Million steps. (All values are averages)

545 A typical successful trajectory of the agents, which were trained in the `Default` environment, can be  
 546 seen in Figure 13. The heading error is only at around the value of 1.7°.



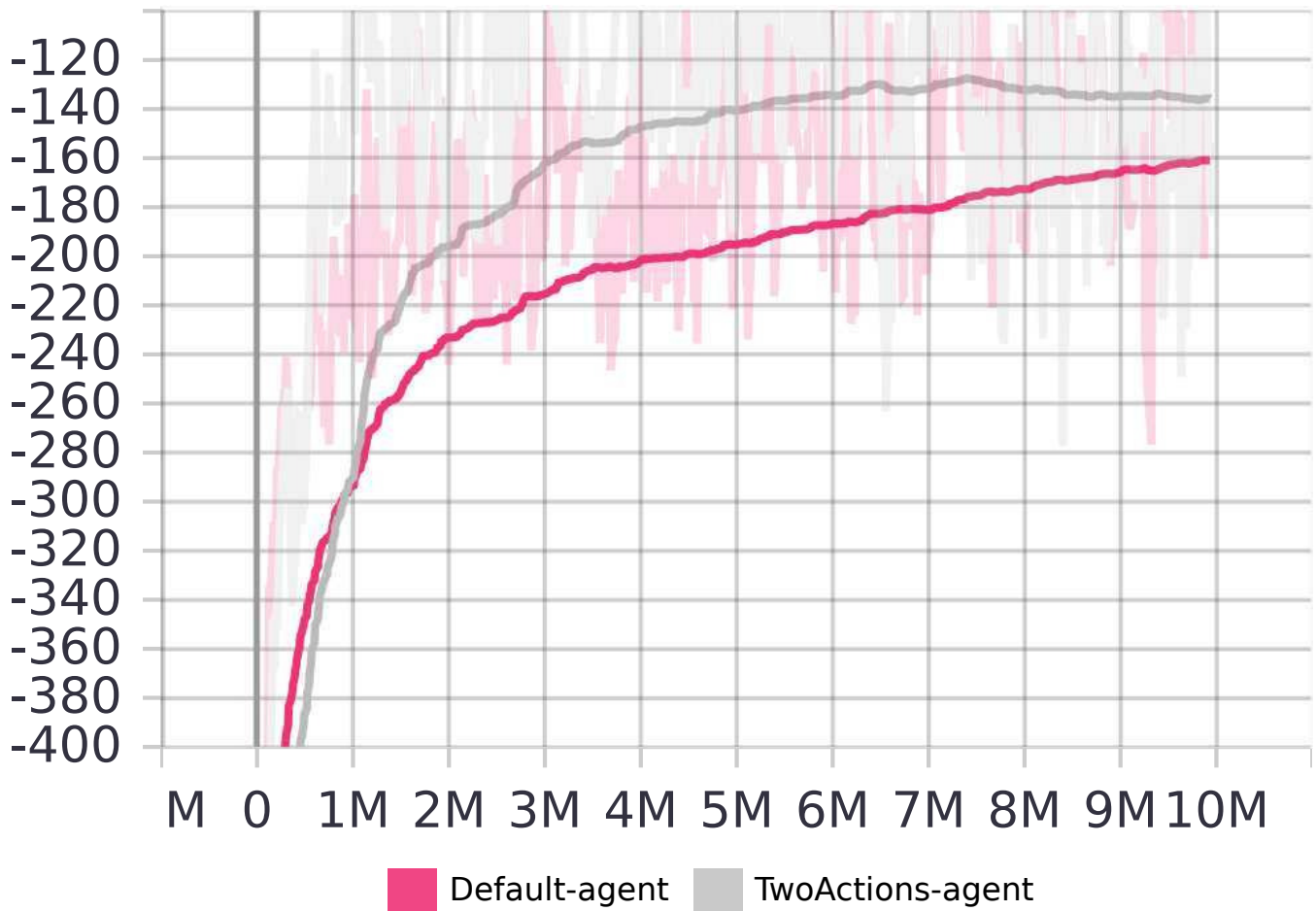
**Figure 13.** Experiment 1 - Typical trajectory example: The aircraft reaches the target successfully. The average heading error is small. (A) 3D view, (B) Top view.

## 547 7.2 Experiment 2: TwoActions environment

548 The average rewards during training in the `TwoActions` environment and the `Default` environment  
 549 were compared. Both were trained for 10 Million steps. Figure 14 shows both reward curves together.  
 550 As expected, the agent in the `TwoActions` environment seemed to receive much higher rewards earlier  
 551 during training. Still, surprisingly, the learning in the `TwoActions` environment seemed to have stopped  
 552 and reached a plateau after 8.5 Million steps. On the other hand, the agent trained in the `Default`  
 553 environment seemed not to have reached a plateau. (Not even at 15 Million steps, as described in Section  
 554 6.1.)

555 The trained `TwoActions` agent was then evaluated and compared to the agent trained in the `Default`  
 556 environment. As can be seen in Table 7 the agent trained in the `Default` environment mostly outper-  
 557 forms the `TwoActions` agent. The `TwoActions` agent's average at-target rate reaches only 20.7%,  
 558 furthermore, its average heading error is around twice as high when compared to the heading error of the  
 559 `Default` agent, and lastly the average distance to the target, when stopping on-track, is also much higher.





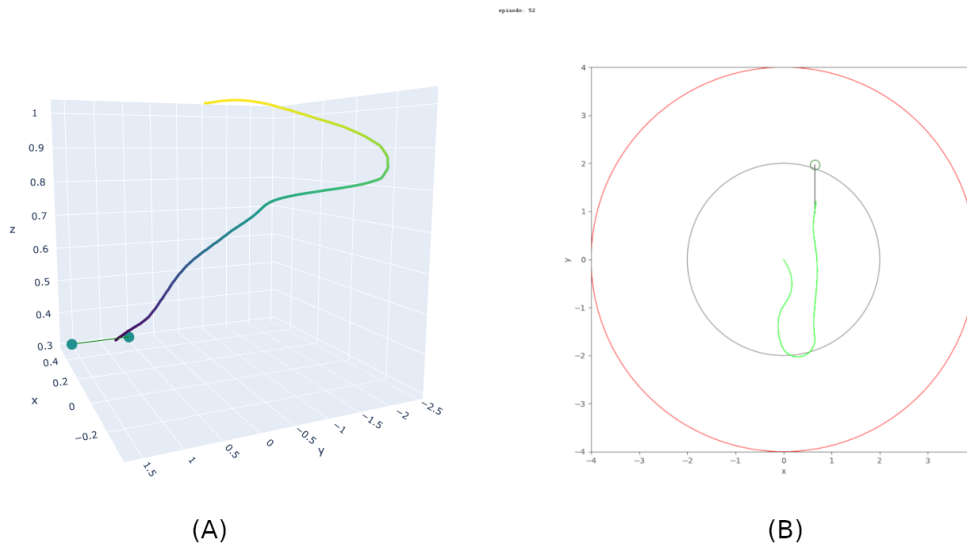
**Figure 14.** Comparison of the agent trained in the `Default` environment (pink curve) with the agent trained in the `TwoActions` environment (gray curve). The learning of the `TwoActions` agent seemed to have reached a plateau after 8.5 Million steps.

560 Nonetheless, the `TwoActions` agent achieved a high reward and on-track rate on average, although the  
 561 on-track rate decreased slightly at around 8.5 Million steps.

	<b>TwoActions</b>	<b>Default</b>
<b>At-Target</b>	20.7%	<b>62.3%</b>
<b>On-Track (includes at-target)</b>	77%	69.3%
<b>Distance to target (when on-track)</b>	0.57 ± 0.45 km	<b>0.095 km ± 0.086 km</b>
<b>Heading error</b>	9.43° ± 5.39°	<b>4.65° ± 7.63°</b>
<b>Reward</b>	<b>-61.57 ± 394.56</b>	-122.37 ± 393.59

**Table 7.** Evaluation comparison of the `Default` agent and the `TwoActions`-agent (All values are averages).

562 In most on-track cases, the agent has led the aircraft first relatively far away from the target and then let  
 563 the aircraft stay on-track for the rest of the episode, as can be seen in Figure 15.



**Figure 15.** Experiment 2 - On-Track: The aircraft lands on-track. The agent seems to have led the aircraft far away first and then kept the aircraft on-track for the rest of the episode.

### 564 7.3 Experiment 3: Wind Environment

565 The evaluation results for the different wind speed ranges are presented in Table 8. As can be seen in the  
 566 table, with increasing wind speed, the performance of the agent decreased moderately, since the on-track  
 567 rate and the reward decreased. Moreover, the heading error and the distance to the target (when the aircraft  
 568 terminates on-track) increased.

569 Nevertheless, for *all* wind ranges together, on average, the agent reached the target 26.25% of the time.  
 570 Furthermore, on average, the on-track rate was quite acceptable at 45.25%. This is also true for the heading  
 571 error at  $8.71^\circ$  and the average value of 0.28 km for the distance to the target.

	0-600 ft/min	0-1200 ft/min	0-2100 ft/min	0-3000 ft/min
<b>At-Target</b>	37.0%	30.3%	22.7%	15%
<b>On-Track</b> (includes at-target)	53.7%	48%	44%	35.3%
<b>Distance to target (if on-track)</b>	0.17 km $\pm$ 0.22 km	0.23 km $\pm$ 0.28 km	0.30 km $\pm$ 0.32 km	0.40 km $\pm$ 0.41 km
<b>Heading error</b>	$6.98^\circ \pm 9.57^\circ$	$7.75^\circ \pm 8.79^\circ$	$8.0^\circ \pm 8.9^\circ$	$12.21^\circ \pm 11.40^\circ$
<b>Reward</b>	$-49 \pm 286$	$-99.70 \pm 330.52$	$-159.89 \pm 374.84$	$-183.89 \pm 370.46$

**Table 8.** Comparison of the wind-agent under different values for the wind speed. The wind is originating from West or East. (All values are averages)

---

## 8 DISCUSSION

572 Three environments were developed to train agents via deep reinforcement learning to achieve the following  
573 goal: Guiding an aircraft to an arbitrary target in the 3D space, after a failure of the aircraft's engines. After  
574 the training was concluded, the trained agents were evaluated in different experiments.

575 Overall, in the windless case, it was found that the agent successfully generated trajectories that guided  
576 the aircraft to reach the target; or at least terminating very close to it. This was especially true when  
577 trained in the `Default` environment. Even when the agent was trained in the `TwoActions` environment,  
578 the aircraft stopped mostly on-track and very close to the target. For the agents, trained in the windless  
579 environments, the aircraft's heading error at the end of the episode was very low. This was also true for the  
580 average distance to the target when ending on-track. These results show that reinforcement learning might  
581 have the potential to be a useful tool for flight guidance in an emergency situation.

582 Experiment 1 showed a relatively high performance after training the agent for around 10 Million steps,  
583 and even a superior performance when trained for 15 Million steps. Even at 15 Million steps, no plateau  
584 for the average reward was reached. This suggests that there is probably still more room for improvement.  
585 For future research, it should be tried to increase the number of training steps.

586 The findings of experiment 2 were quite surprising. As described in Section 5.4 in detail, the  
587 `TwoActions` environment was created to test the following assumption: It was assumed that the agent  
588 trained in the `TwoActions` environment should be superior to the agent trained in the `Default` en-  
589 vironment. However, the results of experiment 2, unexpectedly, demonstrated the contrary. The agent  
590 trained in the `Default` environment (for 10 Million steps) performed better than the `TwoActions` agent.  
591 Nevertheless, the `TwoActions` agent reached a higher average reward and high on-track rate after 10  
592 Million steps. When stopping on-track, the distance to the target was relatively low (Table 7), which means  
593 that the `TwoActions` agent was capable of successfully leading the aircraft closely to the target position.

594 To further investigate the reason for this, the generated on-track trajectories were carefully examined,  
595 and the following was observed: The `TwoActions` agent first guided the aircraft away from the target  
596 and then led it to stay as long as possible on-track, as can be seen in Figure 15 above. This could suggest  
597 that the `TwoActions` agent did indeed learn how to receive higher rewards sooner than the agent of the  
598 `Default` environment. Interestingly, the `TwoActions` agent probably did this by exploiting the fact  
599 that additional rewards were obtained when keeping the aircraft on-track (See Section 5.12). Hence, the  
600 `TwoActions` agent could have learned to collect more rewards by staying on-track as long as possible  
601 instead of ever reaching the target.

602 The agent's learned behavior could be a typical example of reward hacking (Ibarz et al. 2018), where the  
603 intended behavior, specified by the designer of the reward function, differs from the behavior of the agent  
604 (Amodei et al. 2016). Some further research would be needed to confirm this assumption. However, if it is  
605 true, then the reward function needs further adjustments. Solving the problem of reward hacking is by itself  
606 still an open research question (Clark and Amodei 2016; Amodei et al. 2016). Therefore, for now, there is  
607 no simple answer to tackle the issues of reward hacking. Nevertheless, in future research, potential-based  
608 reward shaping (Ng, Harada, and Russell 1999) could be applied to adjust the reward function and to  
609 probably prevent reward hacking by the agent in this work.

610 When wind was involved, disappointingly, the agent's performance dropped. An explanation for the drop  
611 could be that more training time was necessary, since the `Wind` environment has an elevated difficulty.

612 Moreover, this indicates, that probably wind-specific dense rewards should be added to the reward function.  
613 The added dense rewards could improve the learning process of the agent and reduce training time.

## 9 CONCLUSION

614 This work has studied the potential of applying Deep Reinforcement Learning (DRL) to emergency flight  
615 guidance after a complete engines-failure of a fixed-wing aircraft. A DRL agent was trained to act as the  
616 real-time guidance system in a previously unknown environment. The agent learned to generate a 3D path  
617 that led the aircraft from the point of failure to a near, arbitrary-selected landing field.

618 Three different environments were developed and implemented to train the DRL agent. The integration of  
619 curriculum learning and a carefully constructed reward system, which combined sparse and dense rewards,  
620 exceedingly improved the agent's learning performance and reduced training times.

621 Subsequently, the performance of the trained agents were evaluated in the same environments. The  
622 results of the evaluation showed the following: In the windless environments, the agents were capable of  
623 successfully guiding the aircraft to reach the target in a high number of cases. On the other hand, when  
624 wind was involved, the aircraft reached the target successfully in a lower but still significant number of  
625 cases. Furthermore, in around half of all cases the aircraft still landed on-track, meaning, mostly in line  
626 with the runway center line and with a relatively low heading error.

627 The results imply that a guidance system based on Deep Reinforcement Learning is potentially a  
628 valuable instrument, for guiding an aircraft to a landing field after a complete engines-failure. It could  
629 be a significant alternative to systems based on conventional methods. Still, further research is needed,  
630 especially in situations where wind is involved. Nonetheless, the results of this work give a first hint of the  
631 correct direction.

632 To our knowledge, no previous work exists that uses DRL to specifically address the presented problem  
633 above: Generating a 3D path to a landing field, for a fixed-wing aircraft in case of total loss of thrust.  
634 Therefore, the work's presented method and results should contribute to this particular field of study.  
635 Furthermore, the developed environments are available as open-source package and can be used for further  
636 research <sup>2</sup>.

637 Future research should focus on improving the performance of the agent in occurrence of constant wind.  
638 This could be achieved by adapting the reward system and the curriculum learning phases. Moreover, wind  
639 coming from different directions and with changing intensities (e.g., wind gradients) should be investigated,  
640 to imitate the real-world more precisely.

641 As described above, the agent was trained using a lightweight aircraft in simulation. Training another  
642 agent using the here presented method on other aircraft is relatively easy. Nonetheless, future work could  
643 focus on extending the training method to improve the generalization capabilities of one single agent model.  
644 By doing so, different aircraft types could be used during deployment, without the need of training and  
645 using different agent models for every aircraft.

646 In this work, the default hyperparameters of the RLlib's TD3 algorithm were mostly left unchanged.  
647 Although the used TD3 algorithm requires less parameter hyperparameter tuning than it's predecessor  
648 DDPG, future research should study different hyperparameter settings. Automatic hyperparameter search  
649 frameworks like Ray Tune (Liaw et al. 2018) or Optuna (Akiba et al. 2019) could be used.

---

<sup>2</sup> [https://github.com/lauritowal/guidance\\_flight\\_env](https://github.com/lauritowal/guidance_flight_env)

650 Furthermore, a different algorithm, for example Proximal Policy Optimization (PPO) algorithm (Schul-  
651 man et al. 2017), could be investigated to improve the agent's performance. PPO is another popular  
652 state-of-the-art algorithm, that performed exceptional in many other domains.

653 Lastly, in a real-world flight, static and dynamic obstacles (mountains, other aircraft, etc.) might be  
654 present. For this reason, future work should additionally focus on extending the environments to include  
655 obstacles.

## 1 APPENDIX 1

656 The hyperparameters of the Rllib's TD3 algorithm used for the agent.

```
"twin_q": False,
"policy_delay": 1,
"smooth_target_policy": False,
"target_noise": 0.2,
"target_noise_clip": 0.5,
"evaluation_interval": None,
"evaluation_num_episodes": 10,
"use_state_preprocessor": False,
"actor_hiddens": [400, 300],
"actor_hidden_activation": "relu",
"critic_hiddens": [400, 300],
"critic_hidden_activation": "relu",
"n_step": 1,
"exploration_config": {
    "type": "OrnsteinUhlenbeckNoise",
    "random_timesteps": 1000,
    "ou_base_scale": 0.1,
    "ou_theta": 0.15,
    "ou_sigma": 0.2,
    "initial_scale": 1.0,
    "final_scale": 1.0,
    "scale_timesteps": 10000,
},
"timesteps_per_iteration": 1000,
"evaluation_config": {
    "explore": False
},
"buffer_size": 50000,
"prioritized_replay": True,
"prioritized_replay_alpha": 0.6,
"prioritized_replay_beta": 0.4,
"prioritized_replay_beta_annealing_timesteps": 20000,
"final_prioritized_replay_beta": 0.4,
"prioritized_replay_eps": 1e-6,
```

```

"compress_observations": False,
"training_intensity": None,
"critic_lr": 1e-3,
"actor_lr": 1e-3,
"target_network_update_freq": 0,
"tau": 0.002,
"use_huber": False,
"huber_threshold": 1.0,
"l2_reg": 1e-6,
"grad_clip": None,
"learning_starts": 1500,
"rollout_fragment_length": 1,
"train_batch_size": 256,
"num_workers": 0,
"worker_side_prioritization": False,
"min_iter_time_s": 1,

```

## REFERENCES

- 657 (ASA) Federal Aviation Administration (FAA)/Aviation Supplies & Academics (2017) *Pilot's Handbook*  
658 *of Aeronautical Knowledge: FAA-H-8083-25B*. Aviation Supplies & Academics, Incorporated.
- 659 Administration, Federal Aviation (2004) "FAA-H-8083-3A" in: *US Department of Transportation Federal*  
660 *Aviation Administration*, pp. 1–33.
- 661 *Air safety statistics in the EU - Statistics Explained* (Aug. 2020) [https://ec.europa.eu/eurostat/statistics-explained/index.php?title=Air\\_safety\\_statistics\\_in\\_the\\_EU](https://ec.europa.eu/eurostat/statistics-explained/index.php?title=Air_safety_statistics_in_the_EU). (Accessed on 07/16/2021).
- 664 Akiba, Takuya et al. (2019) "Optuna: A Next-generation Hyperparameter Optimization Framework". In:  
665 *Proceedings of the 25rd ACM SIGKDD International Conference on Knowledge Discovery and Data*  
666 *Mining*.
- 667 Akkaya, Ilge et al. (2019) "Solving rubik's cube with a robot hand". In: *arXiv preprint arXiv:1910.07113*.
- 668 Allerton, David (2009) *Principles of flight simulation*. John Wiley & Sons.
- 669 Amodei, Dario et al. (2016) "Concrete problems in AI safety". In: *arXiv preprint arXiv:1606.06565*.
- 670 Berndt, Jon (2004) "JSBSim: An open source flight dynamics model in C++". In: *AIAA Modeling and*  
671 *Simulation Technologies Conference and Exhibit*, p. 4923.
- 672 Berndt, Jon and Agostino De Marco (2009) "Progress on and usage of the open source flight dynamics  
673 model software library, JSBSim". In: *AIAA modeling and simulation technologies conference*, p. 5699.
- 674 Berndt, Jon S et al. (2011) "JSBSim reference manual". In: *Work in progress, available at*  
675 <http://jsbsim.sourceforge.net/JSBSimReferenceManual.pdf>.
- 676 Berner, Christopher et al. (2019) "Dota 2 with large scale deep reinforcement learning". In: *arXiv preprint*  
677 *arXiv:1912.06680*.
- 678 *Best Glide Speed and Distance* (2018) [https://www.faa.gov/news/safety\\_briefing/2018/media/SE\\_Topic\\_18-05.pdf](https://www.faa.gov/news/safety_briefing/2018/media/SE_Topic_18-05.pdf). (Accessed on 07/20/2021).
- 680 Bouhamed, O. et al. (2020) "Autonomous UAV Navigation: A DDPG-Based Deep Reinforcement Learning  
681 Approach". In: *2020 IEEE International Symposium on Circuits and Systems (ISCAS)* pp. 1–5. DOI:  
682 10.1109/ISCAS45731.2020.9181245.
- 683 Brockman, Greg et al. (2016) *OpenAI Gym*. eprint: arXiv:1606.01540.

- 684 Clark, Jack and Dario Amodei (Dec. 2016) *Faulty Reward Functions in the Wild*. <https://openai.com/blog/faulty-reward-functions/>. (Accessed on 07/19/2021).
- 685
- 686 Crocker, David (2007) *Dictionary of aviation: Over 5500 terms clearly defined*. A & C Black.
- 687 Dorr, Les (July 2018) *Fact Sheet – General Aviation Safety*. [https://www.faa.gov/news/fact\\_sheets/news\\_story.cfm?newsId=21274](https://www.faa.gov/news/fact_sheets/news_story.cfm?newsId=21274). (Accessed on 07/16/2021).
- 688
- 689 Florensa, Carlos et al. (13–15 Nov 2017) “Reverse Curriculum Generation for Reinforcement Learning”. In: *Proceedings of the 1st Annual Conference on Robot Learning*. Ed. by Sergey Levine, Vincent Vanhoucke, and Ken Goldberg. Vol. 78. Proceedings of Machine Learning Research. PMLR pp. 482–495. URL: <http://proceedings.mlr.press/v78/florensa17a.html>.
- 690
- 691
- 692
- 693 Goodfellow, Ian, Yoshua Bengio, and Aaron Courville (2016) *Deep Learning*. <http://www.deeplearningbook.org>. MIT Press.
- 694
- 695 Gu, Shixiang et al. (2016) “Continuous deep q-learning with model-based acceleration”. In: *International conference on machine learning*. PMLR pp. 2829–2838.
- 696
- 697 Havenstrøm, Simen Theie, Adil Rasheed, and Omer San (2021) “Deep Reinforcement Learning Controller for 3D Path Following and Collision Avoidance by Autonomous Underwater Vehicles”. In: *Frontiers in Robotics and AI* 7, p. 211.
- 698
- 699
- 700 Henderson, Peter et al. (2018) “Deep reinforcement learning that matters”. In: *Proceedings of the AAAI conference on artificial intelligence*. Vol. 32.
- 701
- 702 Ibarz, Borja et al. (2018) “Reward learning from human preferences and demonstrations in Atari”. In: *arXiv preprint arXiv:1811.06521*.
- 703
- 704 Klein, Marius, Andreas Klos, Jörg Lenhardt, et al. (2018) “Moving target approach for wind-aware flight path generation”. In: *International Journal of Networking and Computing* 8.2, pp. 351–366.
- 705
- 706 Klein, Marius, Andreas Klos, and Wolfram Schiffmann (2020) “A Smart Flight Director for Emergency Landings with Dynamical Recalculation of Stable Glide Paths”. In: *AIAA AVIATION 2020 FORUM* p. 3098.
- 707
- 708
- 709 Liaw, Richard et al. (2018) “Tune: A Research Platform for Distributed Model Selection and Training”. In: *arXiv preprint arXiv:1807.05118*.
- 710
- 711 Lillicrap, Timothy P et al. (2015) “Continuous control with deep reinforcement learning”. In: *arXiv preprint arXiv:1509.02971*.
- 712
- 713 Martinsen, Andreas B and Anastasios M Lekkas (2018a) “Curved path following with deep reinforcement learning: Results from three vessel models”. In: *OCEANS 2018 MTS/IEEE Charleston*. IEEE pp. 1–8.
- 714
- 715 — (2018b) “Straight-path following for underactuated marine vessels using deep reinforcement learning”. In: *IFAC-PapersOnLine* 51.29, pp. 329–334.
- 716
- 717 Mnih, Volodymyr et al. (2013) “Playing atari with deep reinforcement learning”. In: *arXiv preprint arXiv:1312.5602*.
- 718
- 719 Ng, Andrew Y Daishi Harada, and Stuart Russell (1999) “Policy invariance under reward transformations: Theory and application to reward shaping”. In: *Icml*. Vol. 99, pp. 278–287.
- 720
- 721 Perry, Alexander R (2004) “The flightgear flight simulator”. In: *Proceedings of the USENIX Annual Technical Conference*. Vol. 686.
- 722
- 723 Raffin, Antonin et al. (2019) *Stable Baselines3*. <https://github.com/DLR-RM/stable-baselines3>.
- 724
- 725 Randløv, Jette and Preben Alstrøm (1998) “Learning to Drive a Bicycle Using Reinforcement Learning and Shaping.” In: *ICML* vol. 98. Citeseer, pp. 463–471.
- 726

- 
- 727 Rennie, Gordon (2018) “Autonomous Control of Simulated Fixed Wing Aircraft using Deep Reinforce-  
728 ment Learning”. In: [https://purehost.bath.ac.uk/ws/portalfiles/portal/](https://purehost.bath.ac.uk/ws/portalfiles/portal/216919613/Rennie_Gordon.pdf)  
729 [216919613/Rennie\\_Gordon.pdf](https://purehost.bath.ac.uk/ws/portalfiles/portal/216919613/Rennie_Gordon.pdf).
- 730 Schrittwieser, Julian et al. (2020) “Mastering atari, go, chess and shogi by planning with a learned model”.  
731 In: *Nature* 588.7839, pp. 604–609.
- 732 Schulman, John et al. (2017) “Proximal policy optimization algorithms”. In: *arXiv preprint*  
733 *arXiv:1707.06347*.
- 734 Silver, David et al. (2017) “Mastering the game of go without human knowledge”. In: *nature* 550.7676,  
735 pp. 354–359.
- 736 Stephan, Johannes and Walter Fichter (2016) “Fast generation of landing paths for fixed-wing aircraft with  
737 thrust failure”. In: *AIAA Guidance, Navigation, and Control Conference*, p. 1874.
- 738 Sutton, Richard S and Andrew G Barto (2018) *Reinforcement learning: An introduction*. MIT press.
- 739 United States Department of Transportation Federal Aviation Administration, Airman Testing Standards  
740 Branch (2016) *Airplane Flying Handbook FAA-H-8083-3B*.
- 741 Váňa, Petr et al. (2018) “Any-time trajectory planning for safe emergency landing”. In: *2018 IEEE/RSJ*  
742 *International Conference on Intelligent Robots and Systems (IROS)* IEEE pp. 5691–5696.
- 743 Vinyals, Oriol et al. (2019) “Grandmaster level in StarCraft II using multi-agent reinforcement learning”.  
744 In: *Nature* 575.7782, pp. 350–354.
- 745 Wang, Zhuang et al. (2018) “Design of Agent Training Environment for Aircraft Landing Guidance Based  
746 on Deep Reinforcement Learning”. In: *2018 11th International Symposium on Computational Intelligence*  
747 *and Design (ISCID)* vol. 2. IEEE pp. 76–79.
- 748 Xi, Chenyang and Xinfu Liu (2020) “Unmanned Aerial Vehicle Trajectory Planning via Staged Rein-  
749 forcement Learning”. In: *2020 International Conference on Unmanned Aircraft Systems (ICUAS)* IEEE  
750 pp. 246–255.
- 751 Yan, Chao, Xiaojia Xiang, and Chang Wang (2019) “Towards real-time path planning through deep  
752 reinforcement learning for a uav in dynamic environments”. In: *Journal of Intelligent & Robotic Systems*,  
753 pp. 1–13.
- 754 Zhou, Yi et al. (2019) “On the continuity of rotation representations in neural networks”. In: *Proceedings*  
755 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5745–5753.