

**Low-Level control of small scale helicopter using Soft Actor-Critic
method**

by

Majid Kamyab

A thesis submitted in partial fulfillment of the requirements for the degree of

Master of Science

in

Chemical Engineering

Department of Chemical Engineering

University of Alberta

© Majid Kamyab, 2021

Abstract

Unmanned Aerial Vehicles (UAVs), or drones, have been employed in a variety of applications, ranging from surveillance to emergency operations. These systems comprise an "inner loop" that provides stability and control and an "outer loop" in charge of mission-level tasks, such as way-point navigation. Despite their inherent instability, different techniques for controlling these robots have been devised under stable environmental conditions. However, these algorithms must know a robot's dynamics to be effective; furthermore, more complex control is necessary for UAVs to perform in unstable environmental conditions. In this research, a simulated drone has been successfully controlled using model-free reinforcement learning with no prior knowledge of the robot's model. Soft Actor-Critic (SAC) method is trained to perform low-level control of a small-scaled helicopter in a set-point control system. First, a simulation environment is created in which all tests were carried out and then it is shown that SAC can not only develop a strong policy, but it can also deal with unknown circumstances. The result obtained by the SAC agent is also compared to a sliding mode controller to compare the capability of this method to a traditional nonlinear control method. The SMC method proved to be superior by a steady state error of 0, compared to a steady state error of 0.05% for the SAC agent. However, the SAC agent is a model free technique which does not have access to the model of the helicopter, on the other hand the SMC is a model based technique which needs the system identification of the helicopter system.

Acknowledgements

I want to express my gratitude to my distinguished supervisor, Dr. Dubljevic, for his essential supervision, support, and instruction throughout my Master's degree. My thanks go to the Faculty of Engineering for providing me with the money to pursue my studies at the University of Alberta's Department of Chemical Engineering. In addition, I'd like to thank Dr. Koch for his invaluable assistance. I'd want to thank my friends, lab mates, colleagues, and research team – Hamid Khatibi – for a memorable time spent together in academic and social contexts. My thanks also go to my mother, brother, and father for their support during my education.

TABLE OF CONTENTS

1	Introduction	1
1.1	Autonomous UAV	1
1.2	Traditional Control Systems	3
1.3	The Use of Reinforcement Learning as an Optimal Control Method	11
1.4	Simulation Environment for RL	14
1.5	Thesis Objective and outline	15
1.5.1	Chapter 2: Reinforcement Learning Background	15
1.5.2	Chapter 3: Simulation environment	15
1.5.3	Chapter 4: Result and discussion	16
1.5.4	Chapter 5: Conclusion and future work	16
2	Review of Reinforcement Learning	17
2.1	Introduction and terminology	17
2.1.1	Markov Decision Process	17
2.1.2	Dynamic programming	21
2.1.3	Monte Carlo methods	21
2.1.4	Temporal Difference	22
2.2	Policy search	23
2.2.1	Deep Reinforcement Learning	24
2.3	Actor-critic methods	25
2.4	Soft Actor Critic	25
2.4.1	Entropy-Regularized Reinforcement Learning	26
2.4.2	SAC algorithm	26

3	Simulation Environment	29
3.1	Introduction	29
3.2	Governing equations	29
3.2.1	States and control input	29
3.2.2	State-space equations	30
3.2.3	Blade flapping	31
3.2.4	Force derivation	32
3.2.5	Moment derivation	35
3.2.6	Induced velocity	36
3.3	Sliding mode control	39
3.3.1	Force derivation in control affine form	39
3.3.2	Moment derivation in control affine form	42
3.3.3	Control point state space equations	44
3.3.4	Control point position	44
3.3.5	Implementation of sliding mode controller	45
3.4	Environment setup	47
3.4.1	Environment Reset	49
3.4.2	Step	49
3.4.3	Observation	50
3.4.4	Reward	50
3.4.5	Checking for a terminal state	51
3.4.6	Summary	51
4	Result and discussion	52
4.1	SAC agent	52
4.1.1	Architecture	53
4.1.2	Hyper-parameters	53
4.2	Results	56
4.2.1	Training and evaluation	56
4.2.2	Comparison of Controllability and stability to the SMC	56

5	Conclusions and Future Directions	66
5.1	Conclusions	66
5.2	Future work	66
	Bibliography	68

LIST OF TABLES

3.1	Constant parameters in the helicopter modeling	38
4.1	Hyper parameters of SAC agent.	55
4.2	Comparison of SAC and SMC based on the response characteristics of the helicopter dynamic system for the case of initial point set to $[-1,-1,-1]$	58
4.3	Robustness response characteristics of the helicopter dynamic system for the SAC policy by a simulated wind for the case of initial point set to $[-1,-1,-1]$	59

LIST OF FIGURES

2.1	Reinforcement learning schematic and the agent environment interaction.	17
3.1	Function approximation of v_a by v'_a , $R^2 = 0.947$.	38
3.2	Environment Flowchart	48
4.1	SAC controller schematic	52
4.2	Actor diagram for SAC agent, the inputs are the 16 states of the helicopter $X(t)$ and the 4 control inputs $U(t)$, the outputs are the average μ and the standard deviation σ of the agent actions	53
4.3	Critic diagram for SAC agent, inputs are the states and actions while the output is the Q value for the given input.	54
4.4	Actor and critic neural network diagram	54
4.5	Averaged discounted return μ_G and standard deviation σ_G of each iteration using the random actions of the policy.	57
4.6	Helicopter SMC and SAC positional states by initial position $[-1,-1,-1]$.	60
4.7	Helicopter SMC and SAC angular velocities by initial position $[-1,-1,-1]$.	61
4.8	Helicopter SMC and SAC Euler angles by initial position $[-1,-1,-1]$.	62
4.9	Helicopter SMC and SAC velocities by initial position $[-1,-1,-1]$.	63
4.10	Helicopter SMC and SAC control input by initial position $[-1,-1,-1]$.	64
4.11	Helicopter SMC and SAC flapping states by initial position $[-1,-1,-1]$.	65

Nomenclature

Subscripts

cg	Center of gravity	[–]
d	Desired value for the variable	[–]
fus	Fuselage	[–]
ht	Horizontal tail	[–]
mr	Main rotor	[–]
tr	Tail rotor	[–]
vt	Vertical tail	[–]

Symbols

$\mathbb{E}[X]$	Expected value of a random variable X	
α_1	Stabilizer bar rate derivative	53 [–]
α_2	Stabilizer bar input derivative	55 [–]
α_{tail}	Slope of the tail servo angle to the PW of the signal	–1698.5 [rad/s]
α_{tail}	Slope of the tail servo angle to the PW of the signal	[rad/s]
$\delta_{0_{tail}}$	Y-intercept of the tail servo angle to the PW of the signal	[rad]
$\delta_{0_{tail}}$	Y-intercept of the tail servo angle to the PW of the signal	1.4724 [rad]
δ_{coll}	Main rotor collective pitch input	[rad]

δ_{lat}	Lateral cyclic pitch input	[rad]
δ_{lon}	Longitudinal cyclic pitch input	[rad]
δ_{ped}	Tail rotor blade pitch input	[rad]
δ_x, δ_y	Euler rotation angles of the swashplate	[rad]
η	Sliding surface reach time factor	[-]
γ	Discounted rate in RL.	[rad]
λ	surface convergence rate factor	[-]
λ_0	Main rotor inflow ratio	[-]
\mathbb{V}	State value	[-]
\mathcal{D}	Buffer in RL.	[-]
μ	Advance ratio	[-]
μ_x	Non-dimensional airflow components along x axis	[-]
μ_y	Non-dimensional airflow components along y axis	[-]
μ_z	Non-dimensional airflow components along z axis	[-]
Ω	Nominal main rotor speed	115 [rad/s]
$\bar{\mu}$	Normalized Advance ratio	[-]
Φ	Angular velocity transformation matrices from the body to inertial coordinates	[-]
ϕ, θ, ψ	Euler angles	[rad]
Ψ	Blade azimuth angle	[rad]
ρ	Air density	1.107 [kg/m ³]
σ_{mr}	Main rotor solidity factor	[-]

τ_f	Main rotor flapping time-constant	0.04 [s]
τ_{mr}	Main rotor blade element radial distance ratio	[-]
τ_s	Stabilizer bar flapping time-constant	0.2 [s]
τ_{tr}	Tail rotor blade element radial distance ratio	[-]
$\Theta_{0_{tail}}$	Zero pitch angle of the tail blade	[rad]
$\theta_{0_{tail}}$	Zero pitch angle of the tail blade	0.1169 [rad]
Θ_{mr}	Pitch angle of the main rotor	[rad]
Θ_{tr}	Pitch angle of the tail rotor	[rad]
ξ	Ranodm number generated in a normal distribution.	[-]
a_1	Coefficient of the first harmonic approximation in the Fourier series representation of the rotor flapping equations in x direction	[rad/s]
a_κ	Tunable parameter in policy search of RL.	[-]
a_μ	Mean (expectation) of action in a gaussian distribution.	[-]
a_σ	Standard deviation of action in a gaussian distribution.	[-]
A_b	Lateral flapping cross-coupling derivative	-0.1 [-]
a_t	Action of the agent in RL	[-]
a_v	Longitudinal translational velocity contributions to the flapping of the main rotor	[rad/s]
A_{lon}	Longitudinal cyclic to flap gain at nominal rpm	1 [-]
B	Body coordinates	
b_1	Coefficient of the first harmonic approximation in the Fourier series representation of the rotor flapping equations in y direction	[rad/s]
B_A	Longitudinal flapping cross-coupling derivative	0.1 [-]

B_e	bound on b	[−]
b_v	Lateral translational velocity contributions to the flapping of the main rotor	[rad/s]
B_{lat}	Lateral cyclic to flap gain at nominal rpm	0.9875 [−]
C_{D0tr}	Tail rotor zero lift drag coefficient	0.06 [−]
C_{D0}	Main rotor blade zero lift drag coefficient	0.01 [−]
C_{L0}	Main rotor blade zero lift curve slope	0.008 [−]
$C_{L\alpha_{tr}}$	Tail rotor blade lift curve slope	4.95 [rad ^{−1}]
$C_{L\alpha}$	Main rotor blade lift curve slope	5.49 [rad ^{−1}]
c_{mr}	Main rotor chord	0.082 [m]
c_{tr}	Tail rotor chord	0.025 [m]
d	a boolean indicating wheter it is a terminal state or not	[−]
F	Vectors of external forces	[kg.m/s ²]
F_e	bound on f	[−]
F_x	Force along x axis	[kgm/s ²]
F_y	Force along y axis	[kgm/s ²]
F_z	Force along z axis	[kgm/s ²]
G	Expected return of MDP process.	[−]
H	Entropy of a stochastic policy.	[−]
I	Inertia coordinate	
I_s	Equivalent moment of inertia tensor of the TPP rotor disk	
I_{xx}	Rolling moment of inertia	0.3 [kg.m ²]

I_{yy}	Pitching moment of inertia	1.6 [kg.m ²]
I_{zz}	Yawing moment of inertia	2.0 [kg.m ²]
J	Jacobian matrix	[-]
K_λ	Main rotor downwash factor at fuselage	1 [-]
K_c	Longitudinal flapping due to the stabilizer bar factor	[-]
K_d	Lateral flapping due to the stabilizer bar factor	[-]
K_s	Longitudinal flapping cross-coupling derivative	0.3 [-]
K_u	Flapping due to the forward velocity factor	[-]
K_v	Flapping due to the sideways velocity factor	[-]
K_β	Hub torsional stiffness	255 [N.m]
K_μ	Scaling of flap response to speed variation	[-]
K_{lat}	Lateral cyclic to lateral flap gain	0.98 [-]
K_{lon}	Longitudinal cyclic to longitudinal flap gain	1 [-]
L	Mean Squared Bellman Error MSBE.	[-]
l	learning rate	0.3 [kg.m ²]
M	Vectors of external moments	[kg.m ² /s ²]
m	Helicopter mass	11.5 [kg]
M_x	Moment along x axis	[kgm ² /s ²]
M_y	Moment along y axis	[kgm ² /s ²]
M_z	Moment along y axis	[kgm ² /s ²]
n_{tr}	Gear ratio of tail rotor to main rotor	6 [-]

p	Angular rate component (pitch) along x-axis of the CG in I rotated into B	$[rad/s]$
q	Angular rate component (roll) along y-axis of the CG in I rotated into B	$[rad/s]$
Q_{mr}	Drag torque of main rotor	$kg.m/s^2$
r	Angular rate component (yaw) along z-axis of the CG in I rotated into B	$[rad/s]$
R_b^I	Linear velocity transformation matrices from the body to inertial coordinates	$[-]$
R_t	Reward of the environment in RL	$[-]$
R_{mr}	Main rotor radius	$0.95 [m]$
R_{tr}	Tail rotor radius	$0.15 [m]$
S_h	Horizontal tail area	$[m^2]$
s_t	Observation of the environment in RL	$[-]$
S_{vt}	Vertical tail area	$[m^2]$
S_x^{fus}	Frontal fuselage area	$0.1 [m^2]$
S_y^{fus}	Side fuselage area	$0.83 [m^2]$
S_z^{fus}	Vertical fuselage area	$0.51 [m^2]$
T	Thrust	$kg.m/s^2$
U	Vector of input $[\delta_{col}, \delta_{lat}, \delta_{lon}, \delta_{ped}]^T$	$[-]$
u	Velocity component along x-axis of the CG in I rotated into B	$[m/s]$
u_n	Normalized air relative velocity comp along n-axis. in the main rotor rpn coord.	$[-]$

u_p	Normalized air relative velocity comp along p-axis. in the main rotor rpn coord.	[–]
u_r	Normalized air relative velocity comp along r-axis. in the main rotor rpn coord.	[–]
u_{ntr}	Normalized air relative velocity along n-axis in the tail rotor rpn coord.	[–]
u_{ptr}	Normalized air relative velocity along pr-axis in the tail rotor rpn coord.	[–]
u_{rtr}	Normalized air relative velocity along r-axis in the tail rotor rpn coord.	[–]
u_{wind}	Wind velocity along x-axis of the CG in I rotated into B	[m/s]
v	Velocity component along y-axis of the CG in I rotated into B	[m/s]
V_a	Normal-to-the-disk component of the free stream velocity normalized by V_h	[–]
v_a	Axial inflow ratio	[–]
V_{fus}	Dynamic pressure of the fuselage.	[m/s]
V_h	Main rotor induced velocity in hover	[m/s]
V_{itr}	tail rotor induced velocity	[m/s]
V_i	Main rotor induced velocity	[m/s]
v_{wind}	Wind velocity along y-axis of the CG in I rotated into B	[m/s]
w	Velocity component along z-axis of the CG in I translated into B	[m/s]
w_{wind}	Wind velocity along z-axis of the CG in I rotated into B	[m/s]
x, y, z	Position of CG in I coordinates	[m]
x_{fus}	Tail rotor hub offset from CG along x-axis	–1.22 [m]
x_{ht}	Horizontal tail offset from CG along x-axis	[m]
x_{vt}	Horizontal tail offset from CG along x-axis	[m]

z_{cg}	Main rotor hub height from CG	$-0.32 [m]$
z_c	Vertical displacement of the swashplate	$[m]$
z_{vt}	Vertical tail offset from CG along z-axis	$[m]$

Acronyms

CG	Center of gravity	$[-]$
------	-------------------	-------

Chapter 1

Introduction

1.1 Autonomous UAV

Unmanned aerial vehicles (UAVs) are aircraft with no human on board. They are controlled remotely or automatically. Unmanned Aerial Vehicles (UAV) are gaining popularity, both in terms of academic research and potential applications [1].

Classification of the UAVs has two major sub-classes of fixed-wing and rotary-wing. the rotary-wing UAVs received growing attention in recent years thanks to the improvements in embedded microprocessors and batteries. surveillance [2, 3], disaster management [4, 5], and rescue missions [6] are only a few numbers of examples of the broad implementation field of the rotary-wing UAVs.

The majority of recent years' research is focused on quadcopters which are rotary-wing aircraft with four rotors [7, 8, 9, 10] Thanks to their agility and ease of control. On the other hand, single rotor helicopters have gotten less attention from researchers, mainly because they are intrinsically unstable; they have highly coupled nonlinear dynamics, and wind gusts can easily disturb them.

The helicopter is the principal representation of the rotary wing family. The conventional helicopter layout has two engine-powered rotors: the main rotor and the tail rotor. The main rotor generates the thrust power for the helicopter's elevation. The tail rotor offsets the main rotor torque and maintains the helicopter orientation. The change in body orientations of the helicopter results in the inclination of the main rotor, and therefore generating the propulsive force for the helicopter's longitudinal/lateral movement.

All flying features and physical principles of their full-sized counterpart are re-

tained by small helicopters. Moreover, in comparison to full scale helicopters, they are inherently more manoeuvrable and competent. Due to their satisfactory flying ability, size and low expense, UAV science community has engaged in developing minimal cost and reliable autonomous navigation technologies.

Four control inputs are used for the helicopter. Two cyclic controls which handle the helicopter's longitudinal/lateral movement, a collective control of vertical movement and, lastly, the control of pedal control of the helicopter's heading movement. unrestrained helicopter movement is governed by an underactuated structure, in which the number of control inputs (4) is less than the number of degrees of freedom to be controlled (6 DOF), making it difficult to use the traditional approach for controlling Euler-Lagrange systems (which is usually used in the industrial automation). For these reasons, much research has concentrated on control method for unmanned drones that ensures stability and durability. These factors lead to a complex control problem for single rotor small-scaled helicopters. However, the payload capacity of these helicopters is superior to quadcopters, making them more suitable for transportation in emergency situations [11]. As single rotor small, scaled helicopters received less attention, in this study, we will focus on this type of UAVs.

The exact dynamics of the helicopter are unidentified and represented using mechanical relevant mathematical formulas of lesser order, as in most engineering disciplines. It should be emphasized that the estimated model is simply a "abstract concept" since a comprehensive description of the real dynamics of the helicopter is almost infeasible [12].

As a single-rotor helicopter is unstable by nature, it requires a flight control system that operates the vehicle, which is like a human pilot in a large, scaled helicopter. As a result, the flight control can either accept remote control input from an operator or operate autonomously. Remote control of single rotor helicopters is not economically viable, so autonomous control is preferred for most commercial applications. Therefore, the autonomous control of unmanned aerial vehicles (UAVs) is the goal of this research.

1.2 Traditional Control Systems

Control of single rotor helicopters is studied through classic (continuous) or modern (digital) control approaches. Most helicopter systems are inherently non-linear, with non-linear differential equations specified for their dynamics. Researchers, however, generally construct linearized helicopter systems models for analytical purposes. In particular, if this system runs around an operational point and the signals involved are minor, a linear model that estimates a certain non-linear helicopter system may be produced. A large number of approaches have been suggested by researchers for the design and study of control systems for linear systems.

Traditional flight control systems are primarily classified as linear or nonlinear. This categorization is often based on the rotor-craft model expression provided by the controller. Linearization designs are more application-focused and have been used on the majority of helicopter models. Their appeal derives from the ease of control, which reduces both computation cost and duration of the project.

In general, most control systems are based on the broadly established idea of stabilization derivatives, utilizing a linear system of helicopter dynamics. However, a substantial study has been carried out in recent years on non-linear dynamic formulations in the context of helicopter control flight. The concepts of nonlinear controllers are mostly assessed for their conceptual framework to the problem of helicopter navigation. Their application remains a major issue, mostly because of the control system's increasing order and complex nature. Its contribution, however, is crucially important to understand the constraints and possibilities of helicopter navigation.

A linear Multiple-Input Multiple-Output (MIMO) coupled helicopter model serves as the foundation for the linear controller architecture. The internal model method and integral control design are two common design strategies for dealing with the trajectory tracking of linear systems. The proposed control method has the drawback of being complicated to build, whereas integral control is limited to instances where the reference output is a continuous signal. The key principle underlying the linear controller design is to identify the desired state vector for each of these two subsystems, such that when the helicopter status variables converge with their intended state values, the tracking error asymptotically converges to zero. For each

subsystem, the desired state vectors are components and higher derivatives of the reference output vectors.

The linear H_∞ control theory is used for a linear helicopter model such as the one done. However, control laws based on linear helicopter dynamics is not globalized since it shows desired behavior just around a region of operation. This has led to a large number of studies using non-linear control approaches to implement dynamic helicopter models. The feedback linearity control for trajectory tracking was implemented based on a lower order component of the Lagrangian helicopter model [13].

Because of its highly cross-coupling nature of single rotor small scale helicopters (SRSSH), usually, a MIMO approach is implemented [14, 15]. H_∞ method is also used in [16, 17] using a 30-state nonlinear model by an inner loop and outer loop technique. Sliding mode controller is also used for control of SRSSH [18].

Controller design approaches ignore the multivariate character of rotor-craft dynamics as well as the strong link between rotorcraft variables and control inputs. In this sort of framework, each control input is in charge of regulating a single rotorcraft outlet. interconnections between rotorcraft outputs are ignored, and each control input is linked to a SISO feedback loop. The SISO feedback mechanisms associated with the control inputs are totally independent of one another. The SISO feedback mechanisms are built using standard looping platforms [19]. The amplitude and gain tolerances of a feedback loop determine the other's stability. These tolerances define the amount of amplitude and timing that the controller may inject to keep the feedback cycle dynamics constant. However, in the case of multivariable systems, these tolerances can readily lead to erroneous findings.

An 11 state linear model was developed to examine the feedback controller features of the PID technique [20]. Based on the prediction error technique, a time-domain identification procedure was used to identify the set of parameters. The PID design proved unable to reduce the mutual coupling among helicopter's lateral and longitudinal movements, and the aircraft control system was confined to standstill flying. The obtained findings revealed that SISO strategies have mediocre reliability and that multidimensional procedures are essential to minimize the helicopter dynamics' intrinsic strongly coupled impact.

Because of the lag time between the helicopter’s translational and attitude subsystems, most linear control schemes employ a multi-loop control method [21, 22, 23]. Each input controls one helicopter output via a single-input single-output (SISO) feedback system, and the helicopter’s attitude equations are separated from translational motion using two primary control loops. The slower outer-loop regulates the helicopter’s heaving, longitudinal, and lateral movements by computing the needed collective input and attitude angles to guide the aircraft along its intended route. The basis inputs to the inner feedback loop are then these desirable attitude angles. The inner-loop is used to regulate the helicopter’s attitude, which moves at a considerably quicker rate than the translational motion.

A linearized model of the helicopter dynamics is used in the multi-loop approach and the cross-couplings between different DOFs are neglected. Since the cross-coupling dynamics are important, this often results in poor performance of the controller. To account for the cross-couplings that exist between different DOFs of the helicopter, a multi-input multi-output (MIMO) control approach has been used in recent years [14, 24].

Koo et al. use the input-output feedback linearization technique to provide a MIMO solution for the control of small-scale helicopters. The helicopter dynamics are not linearized by the accurate input-output linear system, resulting in instability zero dynamics. The zero dynamics are then stabilized in the simulated world by ignoring the connections between moments and forces and utilizing approximate input-output linearization to obtain limited tracking. Instead of controllable inputs like the collective, cyclic, and pedal inputs, unrealistic control inputs like the gradients of the main and tail rotor thrust and the flapping angles are employed to describe the system [14].

The influence of thrust force components associated with the primary rotor disc displacement is ignored by most nonlinear dynamic systems. These parasite forces have a minimum impact on movement dynamics. This is standard procedure. This approximation leads to several mathematical models with a response form appropriate for backstep control designs laid forth in [25] and numerous researchers used this procedure [26, 27, 28].

Mahony et al. described a MIMO strategy for controlling small-scale aircraft

in hover using a backstepping mechanism [15]. To do this, the flapping behaviors and friction forces are ignored, and the control design is based on a mathematical model of the helicopter dynamics around hover. In a study done by Raptis et al., a time-dependent backstepping approach is used to create a MIMO control scheme for a small-scale helicopter [24]. Simplifying hypotheses are used to generate the helicopter’s dynamic model in a cascading design appropriate for the backstepping control scheme. For instance, in all aviation phases, induced velocity is considered to be constant and the impacts on the thrust computations of the vehicle velocity are disregarded so that main and tail rotors are respectively proportionate in proportion to the input of collectives and pedals. The main and tail rotors’ drag torque is also disregarded.

Another non-linear control scheme is given in a work by Godbolt et al. [29] employing a cascade method. In order to unite attitude and movement dynamics, the internal loop control mechanism is utilized. The control design uses simplification principles. For example, due to the rigidity of the main rotor shaft, the contributions of the rolling and pitching moments to the fuselage dynamic attitude are ignored. Also, because of the rotor blowing in the translational dynamics, it neglects the influence of smaller body forces. A nonlinear control technique is then taken into account to offset the tail rotor’s impacts to small friction forces.

An H_∞ controller’s usual construction consists of two components. The first element consists of Proportional Integral compensators and low pass filters in a manner similar to the traditional approaches of single input single output systems. The Proportional Integral compensators enhance the system’s low-frequency gain, reducing disturbances, and attenuating steady-state error. The low pass filters are generally employed for noise reduction. The second element of the control is the H1 synthesis component, which is determined by a constant signal gain for stabilizing multi-functional dynamic response, as well as being appropriate for a performance criterion [21, 30].

A single value loop forming process based on two degrees of H_∞ freedom was created in the research done by Walker et al. [19] which is an observation basis multivariate controller. The controls were to build a complete autopilot system for a helicopter. The flying system is incorporated with piloted aviation operations, as

opposed to automated flight technologies. The aim of the remote control is for the helicopter to monitor the pilot's control input and speed control. The control scheme is designed to eliminate the connection between axes of helicopter dynamics, therefore lowering the burden of the pilot. The pilot is alone responsible for generating the benchmark and high-speed controls that are required to move the aircraft.

An innovative architecture of static H_∞ output controls was given to stabilize an autonomous helicopter in a hovering citegadewadikar2009h. The optimum control technology makes it possible to devise multivariate feedback systems that enhance the rank of the control unit utilizing fewer states. The structure of the controller feedback loops coincided with the actual flight experience of the helicopter such that the controller's design was acceptable. The H_∞ control system form decreases the influence on high-frequency Helicopters of un-modeled dynamics.

In a research by Kendoul et al. [31] the control design for a Yamaha R-50 helicopter using H_∞ loop forming technology is provided. The control design is composed of non 30-state model of helicopter dynamics in an internal loop approach that is linearized by various operational positions in the desired trajectory envelope. Then an H_∞ loop-fitting controller is built to cover this required flying area based on the acquired linear models.

The UAV control scheme is studied in [32] for a non-linear trajectory tracking control. The non-linear model of helicopter dynamics is discretized and the tracking control issue is then formulated to reduce costs using a quickly converged steepest descent approach. The primary problems of application are the coordination of the cost weight matrix and the constants in the probability density.

In the majority of situations, three nonlinear matrix expressions are required to solve the final loop control issue. In [33], the H_∞ synthesis portion of the controller was resolved by solving just two paired matrix formulas that do not need the information of the initial stabilization gain. There are two principal loops in the control system framework. The first loop is capable of stabilizing the dynamic behavior of the arrangement, and the second loop is for position monitoring. A 13-state linear model of the coupling fuselage and rotor dynamic is the architecture of the control unit. The sequence and structure of the model were adopted in [34].

In another study, Riccati Equation concept is provided [35]. The complicated

dynamics of the helicopter are modified to a pseudo linear, state-dependent (SDC) coefficient and a feedback-optimum matrix is produced at all times by solving the LQR equation. Because there are many non-parametric terms in SDC form and the fact that the helicopter model is not aligned in terms of the control system, it is ignored to achieve a control-affine SDC helicopter dynamics framework necessary for SDRE control designs in certain non-linearity models. A non-linear compensation is then built to increase the control signal to roughly cancel ignored nonlinear effects.

Owing to its resilience with boundary parameter uncertainties, the sliding mode controller can be another non-linear, small-scale, unmanned helicopter management MIMO method. A robust, nonlinear, sliding mode controller flight control is given in [36] for a compact, standalone hover helicopter. The dynamics of the nonlinear helicopter are initially oversimplified by disregarding the drag torque of the rotors and the rear and the connections of the aerodynamic forces and momentum. Then the linearized model is transformed for a squared model into a linear system. For input refined systems, untrue control inputs such as the rolling, pitch, and yaw moments are considered instead of the actual control inputs, and the gradient is considered to be the primary rotor thrust.

The Translational Rate Control (a technique for a UAV is detailed in a study by Pieper et al. [37]) for another sliding mode controller approach in hovers. A fundamental, linearized model of the hovering helicopter and a Sliding mode controller is built to comply with the operating quality requirements for the Translational Rate Control hover control system.

A reference model sliding mode controller design is detailed in a study by Wang et al. [38] and a multi-loop control method is employed to regulate the hover of a UAV. The non-linear helicopter model is modeled linearly around the hover and the coupling movement of the helicopter is ignored, to treat every DOF as a self-contained SISO system. The PID technique is then developed for each of the longitudinal and the lateral controller designs and heavily loaded loops.

Another sliding mode controller technique is presented in a research for controlling a UAV [39]. In this method, the DOFs of helicopter movement are decoupled in these three principal feedback loops: position, speed, and orientation loop. To get an appropriate form for the sliding mode controller method, the Equations of each loop

are simplified. For instance, for the Euler angles in the speed cycle, the small-angle presumption is utilized to linearize the equations and get an input-affined shape. A sliding mode controller for each loop is then designed.

A small-scale autonomous helicopter group control is presented adopting a sliding mode controller approach [40]. In order to produce arbitrary tri-dimensional formations, a sliding mode controller is established for each technique, and the training will be maintained by two leaders/follower controllers. The rotor's flapping complexities are ignored and unrealistic control inputs including the main and tail rotor thrust and pitch and roller moments are applied to describe the system in an input-affine manner instead of actual controlled inputs. The square shape is then exploited to get the control design using a reference points technique.

The aerodynamics of the helicopter is separated into three components with slow, medium, and rapid modes with a multiple time control based on the technique of the slider mode controller [41]. In all flying regimes, nonlinearities of the main and tail rotor intake are removed and the induced speed is presumed to be fixed. A nonlinear controller is built with a sliding mode controller for each mode and results for simulation are provided. Nevertheless, for controls that may result in a non-unique solution, the slow mode controller requires an iterative process.

It is vital that the control architecture is strong enough in the case of the helicopter which has considerable uncertainty. In the presence of parametric and model uncertainty, there is a design that ensures limited traceability [42]. The suggested control scheme includes stabilizing strategies for input saturation feedback systems as well as adaptive nonlinear output control techniques.

In another study, the helicopter model includes the dynamic behavior of the helicopter movement equations that are augmented by a modified aerodynamic force and torque generating model. The Helicopter Dynamics nonlinear model is presented in [43]. In most studies into the design of a non-linear helicopter controller, this particular model was used. The precise linearization input-output fails to linearise the model of the helicopter which leads to instability of zero dynamics. The usage of the approximation model, which does not consider the thrust forces created by the main rotor flap movement, has also been demonstrated to be fully linearized.

In [14], an approximation linearization in input-output was used to achieve a

helicopter system that is dynamically linear without zero dynamics and that has the required characteristic of relative smoothness. The difficulty of an oscillatory ship deck helicopter landing [44] has been appropriately controlled using a conceptual representation. In [45], the design of a floating flight controller for the unmanned APID-MK3 helicopter is described with a unique approach.

In the literature, the majority of control schemes, including Multi-loop and MIMO, are implementing the linear model of the helicopter under various trimming requirements, instead of using the non-linear model directly. This confines the correctness of the linear model to the neighborhood of its linearization of the trimmed requirements. Several linear models are therefore necessary to cover a variety of flying regimes and several gain programmed controllers are required in all such regimes to control the helicopter [46].

Aerodynamic forces and moments fluctuate substantially across different flying circumstances due to the complicated aerodynamic performance of helicopter thrust output. These approximations are not desired for managing an autonomous helicopter over a broad variety of flying phases, through linear system and/or rejection of non-linear components [47].

The issue of optimal control methods is that they all necessitate knowledge of the robot's dynamics, requiring system identification and model derivation for each UAV. Depending on the task, this can become tedious, if not impossible. Notably, the final control system will be a one-of-a-kind solution to a specialized study. These strategies may be insufficient to deal with changing conditions, unanticipated events, and stochastic environments [48].

Previous approaches to nonlinear control using neural networks and nonlinear inversion were published in [22]. Nonlinear control approaches have also been presented. In all situations, the requirements for nonlinear inversion and the increase of a NN raise the controller's order substantially. In this way, it becomes impossible to derive the controller from the helicopter's non-linear governing equations. Consequently, these cases have used developed controls based on the helicopter's linearization dynamics. In the research of Hovakimyan et al. [49] the reduced model uses just the heavy and longitudinal mobility of the helicopter, which further restricts it.

In order to obtain adequate efficiency, the control strategies presented in the re-

search stated above require accurate knowledge of the dynamic models involved. The issue is how to manage unforeseen disturbances to the nominal model in helicopter operations. Unexpected disturbances of this nature usually involve parameters and analytical uncertainty, unmodelled dynamics, and environmental disturbance. The existence of uncertainties and external disturbances can disrupt the feedback controller's operation and lead to significant deterioration. Approximation approaches utilizing artificial neural networks (NN) were suggested to address the presence of model uncertainty. In [50], approximated NN-enhancing dynamic reversal was presented, while in [51] neuronal dynamic programming was demonstrated to be beneficial in the monitoring and trim control of the helicopter.

On this basis, the following question is posed: What if the vehicle teaches itself how to perform a task optimally without using a model? This leads to the next section on reinforcement learning.

1.3 The Use of Reinforcement Learning as an Optimal Control Method

Artificial intelligence (AI) has lately caused a breakthrough in various industries worldwide, ranging from engineering to medical services. Recent advancements in computer technology and data storage, along with AI's learning capacities, have propelled AI to the forefront of numerous applications, such as object recognition and natural language processing. AI is expected to contribute more than 15 trillion USD to the global economy while increasing GDP by 26% by 2030. Overall, artificial intelligence (AI) is a powerful tool that covers many aspects of nowadays scientific achievements [52].

Machine learning (ML) is arguably the most significant branch of AI. It is described as an ability in computer systems that allows them to learn without the need for continuous control over it [53]. The area of machine learning may be divided further into supervised learning, unsupervised learning, semi-supervised learning, and reinforcement learning.

The term "supervised learning" refers to a situation in which the "experience," or training example, provides essential information that is absent in the unknown "test

examples” whereby the learned knowledge is to be implemented. An expert provides the additional information in experience. It tries to generalize across experiences and then applies this knowledge to predict labels for test examples [54]. Since the agent tries to mimic the expert, it will not wholly provide the same response as the expert. This error is called the Bayes error rate [55].

In unsupervised learning, there is no distinction between training data and test data. A typical example of such a job is grouping data collection into subgroups of related objects. Semi-supervised learning is a combination of supervised learning and unsupervised learning. During training, semi-supervised learning mixes a small quantity of labeled data with a lot of unlabeled data, which will improve learning accuracy.

Ideally, supervised learning or semi-supervised learning can completely replicate the supervisor. However, it cannot outperform the supervisor in terms of outcomes. Reinforcement learning (RL) attempts to solve this dilemma by substantial changes to the learning process. Ultimately, the objective of RL is to enable machines to outperform all existing approaches. The RL agent tries to achieve a better result than the currently feasible ones by learning the best mapping of states to actions using a reward signal as a criterion. RL methods allow a vehicle to discover an optimal behavior on its own through trial-and-error interactions with its surroundings. This is based on the commonsense idea that if an action results in a satisfactory or better situation, the tendency to perform that action in the initial situation is *reinforced*.

RL is like classical optimal control theory [56] in engineering platform. Both theorems deal with the problem of determining an input (i.e., optimal controller in control theory or optimal policy in RL) for solving the optimization problem. Furthermore, both rely on a system’s notation being described by an underlying set of states, actions, and a model that captures transitions between one state and the other. So RL can tackle the same problem that optimal control does [57, 58]. However, because the agent does not have access to the state vector dynamics, the agent must learn the repercussions of its actions via trial and error while interacting with the environment.

Although there are some recent achievements on model-based RL [59], most of the RL algorithms are model-free. They attempt to control without the knowledge

of a dynamic model; in other words, it only receives the current states* and a reward from the environment (helicopter in this case) in each step.

This framework has received much attention in recent years, with promising outcomes in a range of domains, including outperforming human specialists on Atari games [60], Go [61], and replicating complex helicopter maneuvers. [62, 63, 64] . A remarkable range of robotics challenges may be conveniently formulated as reinforcement learning problems dating back to 1992 when the OBELIX robot is trained to push objects [65]. A model-free policy gradient technique was used to teach a Zebra Zero robot arm how to perform a peg-in-hole insertion task [66].

Recently, RL-based UAV control has received a lot of interest. The initial research generated an engineered reward function. They developed a model of robot dynamics through demonstration but then employed the model in simulation, leading to the simulation of robot state while using RL to optimize a NN controller for autonomous helicopter flying [67] or inverted helicopter maneuver [63]. However, defining the reward function could be an arduous task. One solution would be to utilize an expert and award the helicopter for emulating the expert’s behavior. Abbeel et al. used this approach to perform aerobatic helicopter flight [62].

In recent years, deep learning has been shown to improve the RL field [68]. Deep learning relies on neural networks’ powerful function approximation properties, which can automatically find compact low-dimensional representations of high-dimensional data (e.g., images). This enabled reinforcement learning methods to scale up to previously unreachable problems.

Deep reinforcement learning has also gained attention recently in UAV control, William Koch et al. [69] compared Deep Deterministic Policy Gradient (DDPG) [70], Trust Region Policy Optimization (TRPO) [71] and Proximal Policy Optimization (PPO) [72] algorithms on the Iris quadcopter and then comparing the result to a PID controller. Although TRPO and DDPG failed to reach stability, they have shown that PPO results are powerful enough to be comparable to a PID controller. Barros and Colombini [73] also proved that the Soft Actor-Critic (SAC) [74] method can perform a low-level control on a commercial quad-rotor Parrot AR Drone 2.0. However, there

*in the fully observable Markov decision process (FOMDP). In the partially observable Markov decision process, a history of states is required in each step.

is still a lack of research on a small-scaled single-rotor helicopter.

1.4 Simulation Environment for RL

In RL, the amount of try and error required to learn beneficial actions is usually high. As a result, sampling the environment is the primary challenge with reinforcement learning. One way to approach this is by having parallel similar real-world environments doing the same thing [75]. However, in the case of the UAV, failure means the loss of a UAV, and hence it is costly. This problem is exacerbated by several real-world factors that make UAVs a problematic domain for RL [76]. UAVs are frequently dangerous and costly to run during the initial training such that the aircraft will fail several times until it reaches a satisfactory performance. This will need high maintenance costs in addition to the original hardware expenditures. Moreover, Robotic have continuous high-dimensional state and action spaces, and finally, it requires a fast online response. As a result, the use of a simulation environment seems necessary for the initial learning procedure of an RL algorithm.

To compensate for the expense of real-world interactions, the UAV must first learn the behavior in simulation and then transfer it to the real vehicle. Usage of a simulator provides an affordable approach in order to create samples. In a simulation, it is possible to crash the UAV as many times as needed; In addition, no safety measures must be taken for, and there would also be no lag in the process due to maintenance or any other real-world issues. Simulations are also more reproducible; For example, wind gusts are not easy to reproduce in the real world, while in simulation, the wind gust random model can be saved and reused elsewhere.

The issue with using a simulation environment is that none of them can completely capture real-world complexity. When a policy is trained in simulation, it usually is not optimal to use in the real world [77]. One possible solution would be to initially train the policy in simulation and then perform tuning in the real world [78, 79].

1.5 Thesis Objective and outline

In this research, we wish to expand on recent research in RL, especially Deep Reinforcement Learning (DRL) to control a SRSSH. More precisely, low-level control rules are learnt directly from the UAV simulation. Notably, the purpose of this thesis is only to train the DRL technique in a simulated setting and providing proof that the aforementioned method is capable of stabilizing the unmanned small scaled helicopter in an acceptable way, leaving future work to examine the transfer to the actual world or produce more complicated maneuvers. In the following, the outline of this thesis is included.

1.5.1 Chapter 2: Reinforcement Learning Background

A wrong choice of RL method or its hyper-parameters can be time-consuming or even impossible to reach good stability of the UAV. This is because it mainly necessitates an extensive exploration of the state-space in order to extract acceptable policies. So, in the second chapter, a review of reinforcement learning methods is discussed. By providing a mathematical framework and describing essential components, this chapter includes a formal introduction to RL. Following that, the chapter provides an overview of Value-based and policy-based methods. Finally, the chapter introduces the DRL algorithm, SAC, which will subsequently be used for UAV control.

1.5.2 Chapter 3: Simulation environment

This chapter introduces the Simulated environment used for interaction with the RL method. First, the helicopter dynamics are discussed, including the forces applied to the UAV, such as fuselage and main rotor forces. Secondly, its effect on the 6 degrees of freedom (DOF) UAV is discussed. In addition a traditional control approach, more specifically sliding mode controller is introduced to compare the result of RL policy with the optimal control theory method. Finally, the environment setup is discussed, including the actions and rewards in the RL platform.

1.5.3 Chapter 4: Result and discussion

This chapter contains the results of applying the RL algorithm on a simulated environment, as well as a discussion and comparing the result of the obtained policy with the one generated by the sliding mode controller. In addition the effect of disturbance on the controller is discussed.

1.5.4 Chapter 5: Conclusion and future work

The conclusion and recommendations for future work are given in the final section of this chapter.

Chapter 2

Review of Reinforcement Learning

2.1 Introduction and terminology

In section 1.3 the RL framework was briefly discussed. In this chapter, the details of this methodology are explained.

2.1.1 Markov Decision Process

MDP is consecutive decision-making in which actions impact immediate rewards and later states, and hence future rewards. In other words, MDP is a stochastic control process using a discrete-time framework. An MDP system consist of 4 components (figure 2.1):

- *states* ($S_t \in \mathcal{S}$): A state (s) is a collection of all essential information about the current situation that can be used to forecast future states. For example, in the case of a robot arm trying to grab a box, the current position of the robot arm could be the state. States can be a multidimensional discrete or continuous set.

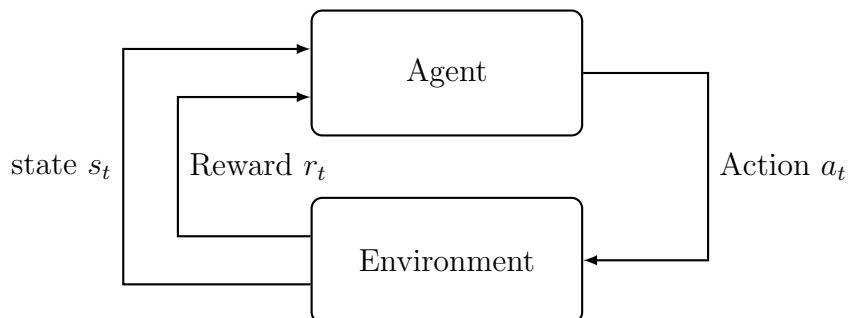


Figure 2.1: Reinforcement learning schematic and the agent environment interaction.

Sometimes, an observation of the states is available instead of the states themselves. For example, instead of a current position of the arm, a snapshot picture is available.

- *action* ($A_t \in \mathcal{A}$): Actions are utilized to control the states by the agents *policy*, which is a mapping from states to actions. It can be either stochastic $a = \pi(\cdot|s)$ or deterministic $a = \mu(s)$. Actions somehow can be compared to the control input in the feedback of a control system. As an example, in a navigation problem, the actions are the torque applied to the wheels. Actions might belong to a discrete or continuous set, and they can also be multidimensional.
- *Reward* ($R_t \in \mathcal{R} \subset \mathbb{R}$): It is the measure of how well the agent is choosing the actions, to put it another way, how well is its *policy*. For example, in the robot arm problem, it could be how close it is to grab the box.
- *environment* (p): The environment is fully described by its *dynamics* (distribution) which can be stochastic $s_{t+1} = p(\cdot|s_t, a_t)$ or deterministic $s_{t+1} = p(s_t, a_t)$. Environment could be any sort of system in which a reward could be defined for a set of given actions applied to the environment.

The MDP framework is conceptual and adaptable, and it may be widely used in a variety of situations in several ways, including the stock price prediction [80] to low-level control of UAVs [81]. Therefore, the definitions are different compared to a control platform. In an MDP, the interaction between the agent (controller) and the environment (the plant, controlled unit) happens in a discrete-time steps platform. The agent performs actions (control signal), receives the reward, and ends up being in a new state. Each interaction between the environment and the agent is called a *step*. in each step the agent receives states S_t and reward $R(s_t)$ from environment and generates a set of action(s) A_t based on its policy which would transform the environment states to a new one S_{t+1} based on transition probability $P(s_{t+1}|s_t, a_t)$ and consecutively provide with a R_{t+1} . So MDP can be defined as a tuple [56]:

$$D \equiv (S, A, P, R) \tag{2.1}$$

Expected Reward can be based on the current state and action $r = r(s, a)$:

$$r(s, a) = \mathbb{E}[R_t | S_{t-1} = s, A_{t-1}] \quad (2.2)$$

Or be based on the state-action-next state:

$$r(s, a) = \mathbb{E}[R_t | S_{t-1} = s, A_{t-1}, S_t = s'] \quad (2.3)$$

Expected return

Broadly speaking, the goal of a policy is to maximize the average reward or discounted *return* (a weighted average in which distant rewards have a less impact) in an episode*. In other words, the goal is to maximize the expected return G_t . There are different ways of defining the expected return [82], here we discuss the one with discounted rate $\gamma \in (0, 1)$ in an episode with T as final time step:

$$G = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots + R_T \quad 0 \leq \gamma \leq 1 \quad (2.4)$$

γ is usually a number close to one since a low γ can result in an instability [76]. If γ is chosen to be 1, then the approach is called average-reward criterion [83]. In this case, it usually cannot distinguish between short-term transient reward, and it is mostly dominated by the steady-state region. If the policy achieves both acceptable short-term and long-term optimal behavior, then it is known as bias optimal [84].

Value function

The *Value function* specifies how good a state is in an episode while a specific policy π is followed. It can be based only on the state $V^\pi(s)$:

$$V^\pi(s) = \mathbb{E}^\pi[G | s_t = s] \quad (2.5)$$

where \mathbb{E}_π denotes the expected return given that the agent follows policy. Note that the value of expected return should be calculated until terminal state is reached. In

* *episode* consists of steps, starting from initial to terminal state, when the terminal state is reached, the process starts over from the initial state.

a similar way the *State-Action Value function* of acting a in state (s) is defined as:

$$Q^\pi(s, a) = \mathbb{E}^\pi[G|s_t = s, a_t = a] \quad (2.6)$$

The value functions are policy dependent, meaning that the value of a state could be low in a policy while it would be high in another one; having this in mind, it is obvious that the optimal value functions are the ones obeying the optimal policy π^* :

$$V^{\pi^*}(s) = \mathbb{E}^{\pi^*}[G|s_t = s] \quad (2.7)$$

$$Q^{\pi^*}(s, a) = \mathbb{E}^{\pi^*}[G|s_t = s, a_t = a] \quad (2.8)$$

$Q^{\pi^*}(s, a)$ is the same as having the optimal policy because given the state (s) we can obtain the optimal policy from the below equation:

$$a^*(s) = \arg \max_a Q^*(s, a) \quad (2.9)$$

Bellman equation

Bellman equation expresses the value of a state, based on the value of its successor states. Bellman equation is obeyed in all the above equations for example in the Value function we have:

$$V^\pi(s_t) = R(s_t, \pi(s_t)) + \gamma \sum_{s_{t+1}} P(s_{t+1}|s_t, \pi(s_t)) V^\pi(s_{t+1}) \quad (2.10)$$

$$V^{\pi^*}(s_t) = R(s_t, \pi^*(s_t)) + \gamma \sum_{s_{t+1}} P(s_{t+1}|s_t, \pi^*(s_t)) V^{\pi^*}(s_{t+1}) \quad (2.11)$$

In a situation with discrete actions, determining the optimal policy is simple, since an exhaustive search is possible if the optimal value function and the transition probabilities for the following states are known, however, in case of continuous spaces, function approximation methods are utilized.

There are numerous value function-based methods which has 3 major classes of:

1. Dynamic programming-based methods.

2. Monte Carlo methods
3. Temporal difference methods.

2.1.2 Dynamic programming

Dynamic Programming (DP) is well suited in a discrete scheme [85]; however, it is possible to use it in a continuous framework. DP uses value functions to arrange and guide the search for optimal policies. The transition probability of the environment should be available, or it could be determined from experience.

In a DP algorithm *policy iteration* is used, which is a process that alternates between *policy evaluation* and *policy improvement*. Initially, a random policy is used to start the approach, then the value function for the current policy is determined by policy evaluation. Each value of state in the current iteration is updated based on the values of the state in the previous iteration (bootstrapping), the policy π , and transition probability p . Finally, the policy is improved based on the most recent value function.

2.1.3 Monte Carlo methods

Unlike the DP, Monte Carlo methods learn directly from *experience*[†] with no prior knowledge of MDP transitions. They carry out rollouts by executing the existing policy on the system, which is referred to as operating on-policy. The value function is updated after an episode is ended. This process is done using the average returns using the current experiences. The frequency of transitions and rewards is recorded and utilized to calculate value function estimates. As more episodes are produced, the average value will converge. The policy is improved by making it greedy regarding the value of the states. Although the method is quite simple, it is pretty powerful; for example, in the game of Tetris, this method outperforms most of the other ones [86].

[†]sampled episodes from environment

2.1.4 Temporal Difference

Temporal Difference (TD) is a generalization of the Monte Carlo method. It also utilizes the bootstrapping of the DP so that TD(1) is the same as the Monte Carlo method, updating the values only when the episode is ended. TD(0) only considers the sampled successor states rather than the full distribution over the successor states in DP. In TD(λ) ($0 \leq \lambda \leq 1$), values are updated before the end of the episode, and more than 1 step ahead is used.

Two popular TD approaches exist, with slightly different update procedures, state-action-reward-state-action (SARSA) [87] and Q-learning [88]. SARSA uses the below equation for updating the Q value.

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[R_{t+1} + \gamma Q(x_{t+1}, u_{t+1}) - Q(s_t, a_t)]$$

While Q-learning uses:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[R_{t+1} + \gamma \max_{a_{t+1}} Q(x_{t+1}, a_{t+1}) - Q(s_t, a_t)] \quad (2.12)$$

SARSA is an on-policy algorithm, which means that its behavior and target policy are the same. Target policy is the output policy of the agent, which is used for evaluating the algorithm. The behavior policy π_b is how the agent acts in exploration. Since exploratory policies are not optimal, *on-policy* methods such as SARSA may quickly converge to a local optimum.

Off-policy agents, such as Q-learning, employ different target and behavior policies; hence, they may use equal probability for taking actions in each state, so $\pi_b(a^*|s) > 0$; as a result, they would find the optimal policy given enough time [89].

Value function methods struggle with the challenges of RL in robotic because they demand data to be filled into the entire state-action space, and they are intrinsically unstable [90]. In addition, the bootstrapping will result in a bias if we want to

use function approximation techniques which is inevitable in continuous spaces of robotics. As a result, value-based methods are not suitable for robotic applications, so we introduce a new family of RL methods called the policy search in the next section.

2.2 Policy search

Policy search approaches do not require using a value function model and instead search for the optimal policy. The concept behind this method is that it is feasible to enhance an episode’s return without knowing the value of each state. The disadvantage of this technique is that it requires evaluating the policy and calculating the return in order to determine if the chosen policy is superior or not. Usually, a parameterized policy is chosen, and the parameters are tuned to maximize the expected return. This is usually done by methods such as gradient ascent [91] or hill climbing [92].

Policy searches provide many advantages. For example, it is feasible to take advantage of an expert for parameter initialization [93], or it is possible to choose the suitable policy parameter structure, ensuring robustness and stability [83]. Therefore, making policy search, a well-suited method for robotic which is proven by real system applications [94, 95].

In the case where gradient ascent is used for the optimization of the policy, we have:

$$J(\pi_\kappa) = \mathbb{E}^{\pi_\kappa}[G_t] \quad (2.13)$$

$$\kappa_{k+1} = \kappa_k + \alpha(\nabla_\kappa J(\pi_\kappa)|_{\kappa=\kappa_k}) \quad (2.14)$$

In which, $\nabla_\kappa J(\pi_\kappa)$ is called *policy gradient* [96]. Which can be expressed as:

$$\nabla_\kappa J(\pi_\kappa) = \gamma^t G_t \frac{\nabla_\kappa \pi_\kappa(a_t|s_t)}{\pi_\kappa(a_t|s_t)} \quad (2.15)$$

$$\nabla_\kappa J(\pi_\kappa) = \gamma^t G_t \nabla_\kappa \log \pi_\kappa(a_t|s_t) \quad (2.16)$$

The term $\nabla_{\kappa} \log \pi_{\kappa}(a_t|s_t)$ is referred to as *eligibility vector*. Equation 2.16 is first introduced by [97] known as REINFORCE algorithm. This algorithm needs the episode to be terminated to calculate G_t , which is why this algorithm is considered a Monte Carlo algorithm. Methods such as Trusted Region Policy Optimization (TRPO) [71] or Proximal Policy Optimization (PPO) [72] are examples of using such methodology.

For continuous actions, instead of learning the probability of the infinite number of actions, usually a Gaussian distribution is used:

$$\pi(a|s, \kappa) = \frac{1}{a_{\sigma}(s, \kappa)\sqrt{2\pi}} \exp\left(-\frac{(a - a_{\mu}(s, \kappa))^2}{2a_{\sigma}(s, \kappa)^2}\right) \quad (2.17)$$

One of the methods to parameterize the policy is using a neural network named Deep Reinforcement Learning.

2.2.1 Deep Reinforcement Learning

In RL, neural networks (NN) are function approximation tools when the state or action space is continuous or too large. In some instances, it is simpler to approximate the value function, whereas, in others, it is easier to approximate policy. In latter cases, policy-based methods are more favorable as they yield a better asymptotic policy [98]. In both cases, a neural network can be employed for value approximation or policy approximation.

Neural networks can learn to map states to values or state-action pairs to Q values. Instead of using a lookup table to store, index, and update all possible states and their values - which is impossible with huge problems- We can train a neural network on samples from the state and action space to predict the value of states or which actions to take given a state.

Now that policy search is introduced, it is possible to discuss the next generation of RL, actor-critic methods.

2.3 Actor-critic methods

Actor-critic methods are policy search methods in which a bias is introduced through bootstrapping in order to improve learning speed and reduce variance. The actor-critic method to the REINFORCE is like the TD algorithm to the Monte Carlo methods.

If only one step of the return is considered, (like TD(0)) the general formula for an actor-critic method can be given as:

$$\kappa_{k+1} = \kappa_k + \alpha(R_{t+1} + \gamma\hat{v}_\omega(S_{t+1}) - \hat{v}_\omega(S_t))\nabla_{\kappa} \log \pi_{\kappa}(a_t|s_t) \quad (2.18)$$

2.4 Soft Actor Critic

Soft Actor-Critic (SAC) is an actor-critic off-policy algorithm with a stochastic policy [99, 74]. It is inspired by stochastic policy optimization and Deep Deterministic Policy Gradient (DDPG) approaches [70]. It has similarities to Twin Delayed DDPG (TD3) method [100] such that both use two clipped Q approximators. Since it is a stochastic method, it also benefits from something similar to target policy smoothing. Which makes it a potent tool in the robotic control field [101].

The main feature of the SAC algorithm is that it tries to balance a trade-off between expected return and entropy [102]. The more the entropy, the higher the exploration, and the less the entropy, the higher the expected return in the short term. This is related to the exploration-exploitation trade-off: increasing entropy leads to more exploration, speeding up learning later. It can also prevent converging to futile local optimums.

Before we can discuss the further details of the algorithm, it is necessary to discuss the details of the usage of entropy in RL.

2.4.1 Entropy-Regularized Reinforcement Learning

The entropy-regularized reinforcement learning changes the goal of RL by including an entropy term, so that the optimal policy not only aims to increase the reward but also tries to increase its entropy at each visited state [103]. The temperature parameter α balance between exploration and exploitation in such way that by increasing α the policy would try to explore more by adding a stochastic term the reward importance. The formula for this method is:

$$\pi^* = \arg \max_{\pi} \mathbb{E}_{\tau \sim \pi} \left[\sum_{t=0}^T \gamma^t (r(s_t, a_t, s_{t+1}) + \alpha H(\pi(\cdot | s_t))) \right] \quad (2.19)$$

in which $H(\pi(\cdot | s_t))$ is the entropy of a stochastic policy, given by:

$$H(\pi(\cdot | s_t)) = \mathbb{E}[-\log \pi(\cdot | s_t)] \quad (2.20)$$

So, comparing a deterministic policy to an entropy regularized policy, when multiple actions are almost equally valuable, the policy commits equal probability mass to the actions instead of choosing the most valuable action. In this framework, the state value and the state-action value should be modified:

$$V_{\pi}(s) = \mathbb{E}_{\tau \sim \pi} \left[\sum_{t=0}^T \gamma^t \left(r(s_t, a_t, s_{t+1}) + \alpha H(\pi(\cdot | s_t)) \right) \middle| s_0 = s \right] \quad (2.21)$$

$$Q_{\pi}(s) = \mathbb{E}_{\tau \sim \pi} \left[\sum_{t=0}^T \gamma^t \left(r(s_t, a_t, s_{t+1}) + \alpha H(\pi(\cdot | s_t)) \right) \middle| s_0 = s, a_0 = a \right] \quad (2.22)$$

2.4.2 SAC algorithm

The SAC algorithm is given in Algorithm 2. The Q functions are updated using the Mean Squared Bellman Error (MSBE)

$$L(\delta_i, \mathcal{D}) = \mathbb{E}_{(s, a, r, s_{t+1}, d) \sim \mathcal{D}} \left[\left(Q_{\delta_i}(s, a) - \underbrace{\left(r + \gamma \min_{j} Q_{\delta_{tar, j}}(s_{t+1}, a_{t+1}) - \alpha \log \pi_{\kappa}(a_{t+1} | s_{t+1}) \right)}_{y(r, s_{t+1}, d)} \right)^2 \right], \quad (2.23)$$

In which D is the buffer of the algorithm in which the transitions are stored. Hence the Q functions are updated by the following gradient:

$$\delta_{i,new} = \delta_{i,old} + l_\delta \nabla_{\delta_{i,old}} \frac{1}{|B|} \sum_{(s,a,r,s_{t+1},d) \in B} (Q_{\delta_{i,old}}(s,a) - y(r, s_{t+1}, d))^2 \quad \text{for } i = 1, 2 \quad (2.24)$$

The policy is updated given:

$$\max_{\kappa} \mathbb{E}_{s \sim \mathcal{D}, \xi \sim \mathcal{N}} \min Q_{\delta_i}(s, a_\kappa(s, \xi)) - \alpha \log \pi_\kappa(a_\kappa(s, \xi) | s), \quad (2.25)$$

The policy is updated by:

$$\kappa_{new} = \kappa_{old} + l_\kappa \nabla_{\kappa} \frac{1}{|B|} \sum_{s \in B} \left(\min Q_{\delta_i}(s, a_\kappa(s)) - \alpha \log \pi_\kappa(a_\kappa(s) | s) \right) \quad (2.26)$$

sampling $a_\kappa(s)$ from Gaussian distribution of policy $\pi_\kappa(\cdot | s)$ is done by the squashed Gaussian function:

$$a_t = f_\kappa(s_t, \xi_t) \quad (2.27)$$

$$a_t = \tanh(\mu_\kappa(s_t) + \sigma_\kappa(s_t) \cdot \xi_t), \quad \xi \in \mathcal{N}(0, I) \quad (2.28)$$

However, after convergence is reached in order to evaluate the policy, the randomness term of the action is omitted to improve performance:

$$\bar{a}_t = \tanh(\mu_\kappa(s_t)) \quad (2.29)$$

Algorithm 2 sac algorithm

```
1: Initialization: initialize policy parameters  $\theta$ 
2: initialize Q-function parameters  $\delta_1, \delta_2$ 
3: initialize target network parameters  $\delta_{\text{targ},1} \leftarrow \delta_1, \delta_{\text{targ},2} \leftarrow \delta_2$ 
4: initializing the replay pool  $\mathcal{D}$ 
5: repeat
6:   repeat
7:     sample action  $a \sim \pi_\theta(\cdot|s)$ 
8:     observe next state  $s_{t+1}$  reward  $r$  and done signal  $d \in [\text{TRUE}, \text{FALSE}]$ 
9:     save  $(s_t, a_t, r(s_t, a_t), s_{t+1}, d)$  in replay pool  $\mathcal{D}$ 
10:    if  $d$  is TRUE then
11:      reset environment state.
12:    end if
13:    until  $\mathcal{D} > \mathcal{D}_{\text{min}}$ 
14:    sample action  $a \sim \pi_\theta(\cdot|s)$ 
15:    observe next state  $s_{t+1}$  reward  $r$  and done signal  $d \in [\text{TRUE}, \text{FALSE}]$ 
16:    save  $(s_t, a_t, r(s_t, a_t), s_{t+1}, d)$  in replay pool  $\mathcal{D}$ 
17:    if  $d$  is TRUE then
18:      reset environment state.
19:    end if
20:    if it's time to update the parameters then
21:      for  $j$  in range(however many updates) do
22:        sample a batch of transitions,  $\mathcal{B} = \{(s, a, r, s_{t+1}, d)\}$  from  $\mathcal{D}$ 
23:        Update Q-functions.
24:        Update policy.
25:        Update target networks by linearization
                
$$\delta_{\text{targ},i} \leftarrow \eta\delta_{\text{targ},i} + (1 - \eta)\delta_i \quad \text{for } i = 1, 2$$

                
$$\kappa_{\text{targ}} \leftarrow \eta\kappa_{\text{targ}} + (1 - \eta)\kappa$$

26:      end for
27:    end if
28:    evaluate the policy to check the convergence using  $\bar{a}_t$ 
29:  until convergence
30: Return  $\kappa, \delta_1$  and  $\delta_2$ 
```

Chapter 3

Simulation Environment

3.1 Introduction

Small-scale helicopters are highly nonlinear systems with complex coupling. Analyzing velocity fields around the rotor requires complicated experiments and numerical methods, which differ in each flight regime such as hover, stall, etc. [104, 105, 106]. There are numerous studies on mathematical models of the helicopter dynamics and governing equations of the forces and moments applied to it [107, 108, 109]. For this study, we have used the model already developed for the Evolution-EX helicopter in [110]. Here we briefly discuss the model development of this helicopter.

3.2 Governing equations

A combination of four subsystems describes the Evolution-EX helicopter's (EEH) dynamics: the rigid-body dynamics of the fuselage, the main rotor, the tail rotor, and the empennage. Two frameworks are defined like other dynamic problems: the body (B) and the Inertia (I) framework.

3.2.1 States and control input

The states regarding the UAV dynamics include the velocity vector $[3 \times 1]$:

$$V = [u, v, w]^T \tag{3.1}$$

In this equation, u denotes the velocity in the x direction, v represents the velocity in the y direction, and w is the velocity in the z direction. and the angular velocity

vector $[3 \times 1]$:

$$\omega = [p, q, r]^T \quad (3.2)$$

In this equation, p represents represent the angular velocity in the p direction, q denotes the angular velocity in the y direction, and r is the angular velocity in the z direction. with respect to body coordinates B and the position vector $[3 \times 1]$:

$$p = [x, y, z]^T \quad (3.3)$$

and the Euler angles vector which is roll pitch and yaw $[3 \times 1]$:

$$\Theta = [\phi, \theta, \psi]^T \quad (3.4)$$

with respect to inertia framework I and the input vector $[4 \times 1]$:

$$U = [\delta_{col}, \delta_{lat}, \delta_{lon}, \delta_{ped}]^T \quad (3.5)$$

the δ_{col} is the collective input which is responsible for increasing the angle of attack in all the angles of the blade plane, the lateral input δ_{lat} is responsible to increase the angle of attack for a lateral movement of the helicopter and δ_{lon} on the other hand is doing the same thing for a longitudinal movement. The δ_{ped} increases the angle of attack of the blades for the tail rotor in the same way that δ_{col} does for the main rotor blades. In the next section the governing equations regarding the states are discussed

3.2.2 State-space equations

The Newton-Euler equations of motion of the helicopter fuselage are defined as:

$$\dot{V} = \frac{1}{m}F - \omega \times V \quad (3.6)$$

$$\dot{\omega} = I^{-1}M - I^{-1}(\omega \times I\omega) \quad (3.7)$$

$$\dot{\Theta} = \Phi(\Theta)\omega \quad (3.8)$$

$$\dot{p} = R_b^I(\Theta)V \quad (3.9)$$

F and M are defined as vector of external forces and moments respectively. Derivation of F and M are elaborated in 3.2.4 and 3.2.5 respectively. R_b^I and Φ are linear and angular velocity transformation matrices given as follows:

$$R_b^I = \begin{bmatrix} s(\theta)c(\psi) & -c(\phi)\sin(\psi) + s(\phi)s(\theta)c(\psi) & s(\phi)s(\psi) + c(\phi)s(\theta)c(\psi) \\ c(\theta)s(\psi) & c(\phi)c(\psi) + s(\phi)s(\theta)s(\psi) & -s(\phi)c(\psi) + c(\phi)s(\theta)s(\psi) \\ -s(\theta) & s(\phi)c(\theta) & c(\phi)c(\theta) \end{bmatrix} \quad (3.10)$$

$$\Phi = \begin{bmatrix} 1 & s(\phi)t(\theta) & c(\phi)t(\theta) \\ 0 & c(\phi) & -s(\phi) \\ 0 & \frac{s(\phi)}{c(\theta)} & \frac{c(\phi)}{c(\theta)} \end{bmatrix} \quad (3.11)$$

in which s, c and t stands for "sin", "cos" and "tan" respectively. The I is the moment of inertia in which the off-diagonal terms are neglected:

$$I = \begin{bmatrix} I_{xx} & 0 & 0 \\ 0 & I_{yy} & 0 \\ 0 & 0 & I_{zz} \end{bmatrix} \quad (3.12)$$

The I_{xx} , I_{yy} and I_{zz} are the rolling, pitching and yawning moment of inertia.

In the following sections, equations regarding the derivation of forces and moments in 3.6 and 3.7 are introduced.

3.2.3 Blade flapping

The dynamics of main rotor and stabilizer bar of the EEH is modeled by *hybrid model approach* [111]. In this approach \dot{a} and \dot{b} are tip-path-plane (TPP) longitudinal and lateral flapping angles respectively and the coefficients of first harmonic approximation in the Fourier series form. The rotor flapping state equations are:

$$\dot{a} = -q - \frac{a}{\tau_f} + \frac{1}{\tau_f}(K_u\mu_x + K_w\mu_z) + \frac{A_{lon}}{\tau_f}(\delta_{lon} + K_c c) + A_b \frac{b}{\tau_f} \quad (3.13)$$

$$\dot{b} = -p - \frac{b}{\tau_f} + \frac{1}{\tau_f}(K_v\mu_y) + \frac{B_{lat}}{\tau_f}(\delta_{lat} + K_d d) + B_a \frac{a}{\tau_f} \quad (3.14)$$

K_u is the flapping due to the forward velocity factor, the K_v on the other hands denotes the flapping due to the side-way velocity factor. In the same way, K_w is considered to be the flapping due to downward velocity factor. K_c reflects the longitudinal flapping due to the stabilizer bar factor and K_d depicts the lateral flapping due to the stabilizer bar factor. A_{lon} denotes the longitudinal cyclic to flap gain at nominal rpm and B_{lat} is the lateral cyclic to flap gain at nominal rpm. Last but not least, the τ_f is the main rotor flapping time-constant.

μ_x, μ_y and μ_z are the non-dimensional airflow components defined as:

$$\begin{aligned}\mu_x &= \frac{u - u_{wind}}{\Omega R_{mr}} \\ \mu_y &= \frac{v - v_{wind}}{\Omega R_{mr}} \\ \mu_z &= \frac{w - w_{wind}}{\Omega R_{mr}}\end{aligned}\tag{3.15}$$

And the K_u, K_v and K_w are given by:

$$K_u = 2K_\mu \left(\frac{4}{3} \delta_{col} - \frac{Vi}{\omega R_{mr}} \right)\tag{3.16}$$

$$K_v = -K_u\tag{3.17}$$

$$K_w = 16K_\mu \mu_{mr}^2 \frac{\text{sign}(\mu_{mr})}{(1 - \mu_{mr}^2/2) * (8\text{sign}(\mu_{mr}) + CL_\alpha \sigma)}\tag{3.18}$$

The stabilizer bar state equations c and d are TPP longitudinal and lateral flapping angles of the stabilizer bar given by:

$$\dot{c} = -q - \frac{c}{\tau_s} + \frac{C_{lon}}{\tau_s} \delta_{lon}\tag{3.19}$$

$$\dot{d} = -p - \frac{d}{\tau_s} + \frac{D_{lat}}{\tau_s} \delta_{lat}\tag{3.20}$$

3.2.4 Force derivation

The force is derived as follows:

$$F = \begin{bmatrix} F_x \\ F_y \\ F_z \end{bmatrix} + R_b^I \begin{bmatrix} 0 \\ 0 \\ mg \end{bmatrix} \quad (3.21)$$

in which F_x , F_y and F_z are defined as:

$$F_x = -a T_{mr} + F_{x,fus} \quad (3.22)$$

$$F_y = b T_{mr} + T_{tr} + F_{y,fus} + F_{y,vt} \quad (3.23)$$

$$F_z = -T_{mr} + F_{z,fus} + F_{z,ht} \quad (3.24)$$

T_{mr} is the main rotor thrust and T_{tr} denotes the tail rotor thrust. T_{mr} is given by:

$$T_{mr} = \mathbf{f}_{T_{mr}} + \mathbf{b}_{T_{mr}} U; \quad (3.25)$$

in which $\mathbf{f}_{T_{mr}}$ is the autonomous term of the main rotor thrust:

$$\mathbf{f}_{T_{mr}} = \frac{1}{4} \rho \pi R_{mr}^4 \Omega^2 \sigma_{mr} (C_{L0} (\frac{2}{3} + \mu_x^2 + \mu_y^2) + C_{L\alpha} (\mu_z - \lambda_0)) \quad (3.26)$$

In which R_{mr} is the main rotor radius and Ω represents the nominal main rotor speed. σ_{mr} depicts the main rotor solidity factor and C_{L0} and $C_{L\alpha}$ are the main rotor blade zero lift curve and blade lift curve slope respectively.

λ_0 is the inflow ratio expressed as:

$$\lambda_0 = \frac{V_i}{\Omega R_{mr}} \quad (3.27)$$

The solidity factor derived by:

$$\sigma_{mr} = \frac{N c_{mr}}{\pi R_{mr}} \quad (3.28)$$

and $b_{T_{mr}}$ in 3.25 is control input coefficient term given by:

$$\mathbf{b}_{T_{mr}} = \frac{1}{4} \rho \pi R_{mr}^4 \Omega^2 \sigma_{mr} C_{L\alpha} \begin{bmatrix} \mu_x^2 + \mu_y^2 + \frac{2}{3} & -\mu_y & \mu_x & 0 \end{bmatrix} \quad (3.29)$$

Similarly, it is possible to derive the tail rotor thrust T_{tr} :

$$T_{tr} = \mathbf{f}_{T_{tr}} + \mathbf{b}_{T_{tr}} U; \quad (3.30)$$

In which $\mathbf{f}_{T_{tr}}$ is the autonomous term of the tail rotor thrust:

$$\mathbf{f}_{T_{tr}} = -\frac{1}{4}\rho\pi R_{tr}^4 n_{tr}^2 \Omega^2 \sigma_{tr} C_{L\alpha_{tr}} v_{tail} \quad (3.31)$$

And input coefficients $\mathbf{b}_{T_{tr}}$ is given by:

$$\mathbf{b}_{T_{tr}} = -\frac{1}{4}\rho\pi R_{tr}^4 n_{tr}^2 \Omega^2 \sigma_{tr} C_{L\alpha_{tr}} \left[0 \quad 0 \quad 0 \quad u_{tail}^2 + w_{tail}^2 + \frac{2}{3} \right] \quad (3.32)$$

The normalized velocities at tail rotor can be given as:

$$\begin{aligned} u_{tail} &= \frac{u - u_{wind}}{\Omega_{tr} R_{tr}} \\ v_{tail} &= \frac{v - v_{wind} - V_{i_{tr}} + x_{fus} r}{\Omega_{tr} R_{tr}} \\ w_{tail} &= \frac{w - w_{wind} - K_\lambda V_i + x_{fus} q}{\Omega_{tr} R_{tr}} \end{aligned} \quad (3.33)$$

The tail rotor nominal speed is given by:

$$\Omega_{tr} = n_{tr} \Omega \quad (3.34)$$

In which the n_{tr} is gear ratio of tail rotor to main rotor. $F_{y,vt}$ is the vertical tail force derived by:

$$F_{y,vt} = \frac{1}{2}\rho S_{vt} \left(C_{L\alpha}^{vt} V_{vt} (v - v_{wind}) + v_{tail}^2 \right) \quad (3.35)$$

In the above equation, the ρ is the air density. S_{vt} denotes the vertical tail area, And horizontal tail force $F_{z,ht}$ is:

$$F_{z,ht} = \frac{1}{2} S_{ht} \left(C_{L\alpha}^{ht} |u - u_{wind}| w_{ht} + w_{ht}^2 \right) \quad (3.36)$$

In equation 3.35 V_{vt} and v_{tail} are axial and normal velocities in vertical tale defined as:

$$\begin{aligned} V_{vt} &= \sqrt{(u - u_{wind})^2 + (w - w_{wind} + x_{vt} q - K_\lambda V_i)^2} \\ v_{tail} &= v - v_{wind} + x_{vt} r - V_{i_{tr}} \end{aligned} \quad (3.37)$$

Similarly in equation w_{ht} is the horizontal tail velocity in the z direction 3.36:

$$w_{ht} = w - w_{wind} - x_{ht}q - K_{\lambda}V_i \quad (3.38)$$

The K_{λ} is the main rotor downwash factor at fuselage.

$F_{z,fus}$, $F_{y,fus}$ and $F_{x,fus}$ are drag forces derived by:

$$F_{x,fus} = -\frac{1}{2}\rho S_x^{fus} V_{fus}(u - u_{wind}) \quad (3.39)$$

$$F_{y,fus} = -\frac{1}{2}\rho S_y^{fus} V_{fus}(v - v_{wind}) \quad (3.40)$$

$$F_{z,fus} = -\frac{1}{2}\rho S_z^{fus} V_{fus}(w - w_{wind} + V_i) \quad (3.41)$$

the dynamic pressure of the fuselage V_{fus} in expression 3.67 is defined as:

$$V_{fus} = \sqrt{(u - u_{wind})^2 + (v - v_{wind})^2 + (w - w_{wind} + V_i)^2} \quad (3.42)$$

3.2.5 Moment derivation

Moment includes 3 terms roll, pitch and yaw:

$$M = \begin{bmatrix} M_{roll} \\ M_{pitch} \\ M_{yaw} \end{bmatrix} \quad (3.43)$$

The three terms in the above equation are given as:

$$M_{roll} = (K_{\beta} - T_{mr}Z_{cg})b; \quad (3.44)$$

$$M_{pitch} = (K_{\beta} - T_{mr}Z_{cg})a; \quad (3.45)$$

$$M_{yaw} = Q_{mr} + T_{tr}x_{fus}; \quad (3.46)$$

Z_{cg} is the Main rotor hub height from center of gravity and x_{fus} is the tail rotor hub offset from center of gravity along x-axis. Main rotor drag torque Q_{mr} is derived by:

$$Q_{mr} = \mathbf{f}_{Q_{mr}} + \mathbf{b}_{Q_{mr}} U; \quad (3.47)$$

In which the main rotor drag torque terms are given as:

$$\mathbf{f}_{Q_{mr}} = \frac{1}{8} \rho \pi R_{mr}^5 \Omega^2 \sigma_{mr} C_{L\alpha} \left(\frac{C_{D0}}{C_{L\alpha}} (\mu_x^2 + \mu_y^2 + 1) - 2(\mu_z - \lambda_0)^2 \right) \quad (3.48)$$

$$\mathbf{b}_{Q_{mr}} = \frac{1}{8} \rho \pi R_{mr}^5 \Omega^2 \sigma_{mr} C_{L\alpha} (\lambda_0 - \mu_z) \begin{bmatrix} \frac{4}{3} - \mu_y & \mu_x & 0 \end{bmatrix} \quad (3.49)$$

3.2.6 Induced velocity

As indicated in [112] and [113] the blade element analysis considers each blade element as a two-dimensional airfoil. The aerodynamic behavior of neighboring blade elements is independent of each other. An induced inflow velocity on each blade element should be accounted for, which is a product of the rotor wake. Analytical ways of calculating the induced velocity may be found using momentum theory, vortex theory, or nonuniform inflow calculations [112].

In general, the inflow velocity calculation is a challenging task due to its non-uniformity across the blade span; mathematical simplifications should be applied to minimize the complexity of the analysis. Finally, after determining the velocity components of the blade element, the aerodynamic forces acting on this element are calculated. The complete dynamic behavior of the blade is obtained by integrating the applied forces of the individual elements throughout the blade span. Here we use an experimental approach for induced velocity.

V_i and V_{itr} are induced velocity in main rotor and tail rotor respectively. V_i is given by:

$$V_i = \frac{v_a}{\sqrt{1 + \bar{\mu}^2}} \quad (3.50)$$

In which:

$$V_h = \sqrt{mg/(2\rho\pi R_{mr}^2)} \quad (3.51)$$

$$\mu = \sqrt{\mu_x^2 + \mu_y^2} \quad (3.52)$$

$$\bar{\mu} = \frac{\mu}{V_h/(\Omega R_{mr})} \quad (3.53)$$

$$V_a = -\frac{w - w_{wind}}{V_h} \quad (3.54)$$

$$v_a = \begin{cases} -\frac{1}{2}V_a - \sqrt{\frac{V_a^2}{4} - 1} & \text{if } V_a \leq -2 \\ 1 - \frac{1}{2}V_a + \frac{25}{12}V_a^2 + \frac{7}{6}V_a^3 & \text{if } -2 < V_a < 0 \\ -\frac{1}{2}V_a + \sqrt{\frac{V_a^2}{4} + 1} & \text{if } V_a \geq 0 \end{cases} \quad (3.55)$$

Similarly, V_{itr} is:

$$V_{itr} = \frac{v_{atr}}{\sqrt{1 + \mu_{tr}^2}} \quad (3.56)$$

In which:

$$V_{h_{tr}} = \sqrt{f_{E_y, mr}/(2\rho\pi R_{tr}^2 x_{fus})} \quad (3.57)$$

$$\mu_{tr} = \sqrt{u_{tail}^2 + w_{tail}^2} \quad (3.58)$$

The $\bar{\mu}_{tr}$ is the advance ratio normalized by $V_{h_{tr}}/(\Omega R_{tr})$:

$$\bar{\mu}_{tr} = \frac{\mu_{tr}}{V_{h_{tr}}/(\Omega R_{tr})} \quad (3.59)$$

$$V_{atr} = -\frac{v - v_{wind} + x_{fus}r}{V_h} \quad (3.60)$$

$$v_{atr} = \begin{cases} -\frac{1}{2}V_{atr} - \sqrt{\frac{V_{atr}^2}{4} - 1} & \text{if } V_{atr} \leq -2 \\ 1 - \frac{1}{2}V_{atr} + \frac{25}{12}V_{atr}^2 + \frac{7}{6}V_{atr}^3 & \text{if } -2 < V_{atr} < 0 \\ -\frac{1}{2}V_{atr} + \sqrt{\frac{V_{atr}^2}{4} + 1} & \text{if } V_{atr} \geq 0 \end{cases} \quad (3.61)$$

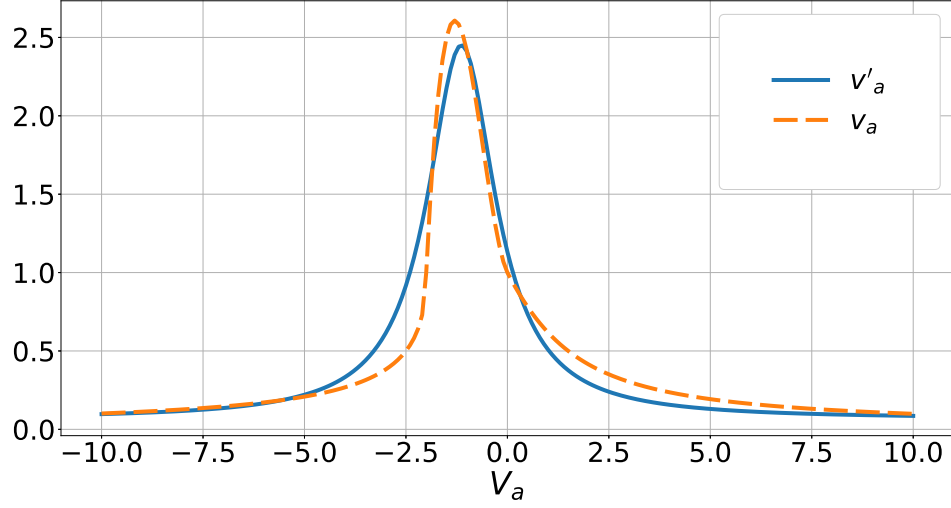


Figure 3.1: Function approximation of v_a by v'_a , $R^2 = 0.947$.

Instead of using equation 3.55 and 3.61, we use the following approximate functions by regression as they provide a faster calculation time in simulations:

$$v'_a = \frac{4.055}{(1.28V_a + 1.45)^2 + 1.7} + 0.066 \quad (3.62)$$

$$v'_{atr} = \frac{4.055}{(1.28V_{atr} + 1.45)^2 + 1.7} + 0.066 \quad (3.63)$$

Figure 3.1 depicts the two functions v_a and v'_a based on V_a .

The constant parameters in the above modeling is given in the table 3.1.

Table 3.1: Constant parameters in the helicopter modeling

Parameter	Value	Dimension	Parameter	Value	Dimension
α_1	53	[-]	α_2	55	[-]
Ω	115	[rad/s]	ρ	1.107	[kg/m ³]
τ_f	0.04	[s]	τ_s	0.2	[s]
τ_s	0.2	[s]	A_b	-0.1	[-]
A_{lon}	1	[-]	B_a	0.1	[-]
B_{lat}	0.9875	[-]	C_{D0tr}	0.06	[-]
C_{D0}	0.01	[-]	C_{L0}	0.008	[-]
$C_{L\alpha_{tr}}$	4.95	[rad ⁻¹]	$C_{L\alpha}$	5.49	[rad ⁻¹]

Table 3.1 Continued:

c_{mr}	0.082	[m]	c_{tr}	0.025	[m]
I_{xx}	0.3	[kg.m ²]	I_{yy}	1.6	[kg.m ²]
I_{zz}	2.0	[kg.m ²]	K_λ	1	[-]
K_s	0.3	[-]	K_β	255	[N.m]
K_{lat}	0.98	[-]	K_{lon}	1	[-]
m	11.5	[kg]	n_{tr}	6	[-]
R_{mr}	0.95	[m]	R_{tr}	0.15	[m]
S_x^{fus}	0.1	[m ²]	S_y^{fus}	0.83	[m ²]
S_z^{fus}	0.51	[m ²]	x_{fus}	-1.22	[m]
z_{CG}	-0.32	[m]			

3.3 Sliding mode control

In order to compare the results obtained by the RL method to a traditional control method, we here provide the details of a sliding mode controller as a nonlinear method. As it is an under actuated system which means that there are only have 4 inputs and 6 states to control. Based on [114, 110] in order to use sliding mode controller on a helicopter, first we have to change it to a square and affine in control form so we discuss how this is achieved.

3.3.1 Force derivation in control affine form

In order to provide a control input by sliding mode controller, each of the force and moment should be linearized based on the control input. Control-affine form of each term is given as a input coefficient "b" and a "f" term which is related to autonomous response of the system. Force is given as:

$$F = \mathbf{f}_F + \mathbf{b}_F U \quad (3.64)$$

f_F is the autonomous term of the force expressed as:

$$\mathbf{f}_F = \begin{bmatrix} \mathbf{f}_{F_{x,mr}} + F_{x,fus} \\ \mathbf{f}_{F_{y,mr}} + \mathbf{f}_{F_{y,tr}} + F_{y,vt} + F_{y,fus} \\ \mathbf{f}_{F_{z,mr}} + F_{z,ht} + F_{z,fus} \end{bmatrix} + R_b^{I^{-1}} \begin{bmatrix} 0 \\ 0 \\ mg \end{bmatrix} \quad (3.65)$$

in which:

$$\mathbf{f}_{F_{x,mr}} = \mathbf{f}_T(\tau_f q - a_v) \quad (3.66)$$

$$F_{x,fus} = -\frac{1}{2}\rho S_x^{fus} V_{fus}(u - u_{wind}) \quad (3.67)$$

$$\mathbf{f}_{F_{y,mr}} = (\tau_f p + b_v)\mathbf{f}_T \quad (3.68)$$

$$\mathbf{f}_{F_{y,tr}} = \mathbf{f}_{T_{tr}} \quad (3.69)$$

$$\mathbf{f}_{F_{z,mr}} = -\mathbf{f}_T \quad (3.70)$$

a_v and b_v in 3.66 and 3.68 are respectively longitudinal and lateral translational velocity contributions to the flapping of the main rotor and defined by the following terms:

$$\begin{aligned} a_v &= \frac{\partial a_1}{\partial \mu_x} \mu_x + \frac{\partial a_1}{\partial \mu_z} \mu_z \\ b_v &= \frac{\partial a_1}{\partial \mu_y} \mu_y \end{aligned} \quad (3.71)$$

the dynamic pressure of the fuselage in expression 3.67 is defined as:

$$V_{fus} = \sqrt{(u - u_{wind})^2 + (v - v_{wind})^2 + (w - w_{wind} + V_i)^2} \quad (3.72)$$

$\mathbf{f}_{T_{mr}}$ is the main rotor thrust autonomous term in the control affine form given by 3.26. and λ_0 is the inflow ratio expressed as:

$$\lambda_0 = \frac{V_i}{\Omega R_{mr}} \quad (3.73)$$

σ_{mr} is the solidity factor derived by:

$$\sigma_{mr} = \frac{N c_{mr}}{\pi R_{mr}} \quad (3.74)$$

The velocities at tail rotor can be normalized given as:

$$\begin{aligned}
u_{tail} &= \frac{u - u_{wind}}{\Omega_{tr} R_{tr}} \\
v_{tail} &= \frac{v - v_{wind} - V_{itr} + x_{fus} r}{\Omega_{tr} R_{tr}} \\
w_{tail} &= \frac{w - w_{wind} - K_\lambda V_i + x_{fus} q}{\Omega_{tr} R_{tr}} \\
\Omega_{tr} &= n_{tr} \Omega
\end{aligned} \tag{3.75}$$

similarly velocities at vertical tail or horizontal tail can be defined.

In equation 3.35 $V_v t$ and v_{tail} are axial and normal velocities in vertical tale defined as:

$$\begin{aligned}
V_{vt} &= \sqrt{(u - u_{wind})^2 + (w - w_{wind} + x_{vt} q - K_\lambda V_i)^2} \\
v_{tail} &= v - v_{wind} + x_{vt} r - V_{itr}
\end{aligned} \tag{3.76}$$

Similarly in equation 3.36:

$$w_{ht} = w - w_{wind} - x_{ht} q - K_\lambda V_i \tag{3.77}$$

b_F in equation 3.64 is the input coefficient of the force which is defined as:

$$\mathbf{b}_F = \begin{bmatrix} \mathbf{b}_{F_{x,mr}} \\ \mathbf{b}_{F_{y,mr}} + \mathbf{b}_{F_{y,tr}} \\ \mathbf{b}_{F_{z,mr}} \end{bmatrix} \tag{3.78}$$

In which the following terms are used:

$$\mathbf{b}_{F_{x,mr}} = (\tau_f q - a_v) \mathbf{b}_T - [0 \ 0 \ K_{lon} \mathbf{f}_T \ 0] \tag{3.79}$$

$$\mathbf{b}_{F_{y,mr}} = (-\tau_f p + b_v) \mathbf{b}_T - [0 \ K_{lat} \mathbf{f}_T \ 0 \ 0] \tag{3.80}$$

$$\mathbf{b}_{F_{y,tr}} = \mathbf{b}_{T_{tr}} \tag{3.81}$$

$$\mathbf{b}_{F_{z,mr}} = -\mathbf{b}_T \tag{3.82}$$

b_T and $b_{T_{tr}}$ is the main rotor thrust input coefficient in the control affine form given as 3.29 and 3.32.

3.3.2 Moment derivation in control affine form

Control affine form of moment is given as:

$$M = \mathbf{f}_M + \mathbf{b}_M U \quad (3.83)$$

\mathbf{f}_M is calculated by:

$$\mathbf{f}_M = \begin{bmatrix} \mathbf{f}_{M_x, mr} + \mathbf{f}_{M_x, tr} + M_{x, vt} \\ \mathbf{f}_{M_y, mr} + \mathbf{f}_{M_y, tr} + M_{y, ht} \\ \mathbf{f}_{M_z, mr} + \mathbf{f}_{M_z, tr} + M_{z, vt} \end{bmatrix} \quad (3.84)$$

in which:

$$\mathbf{f}_{M_x, mr} = (-\tau_f p + b_v)(K_\beta - \mathbf{f}_T * z_{cg}) \quad (3.85)$$

$$\mathbf{f}_{M_x, tr} = -z_{fus} \mathbf{f}_{T_{tr}} \quad (3.86)$$

$$M_{x, vt} = -F_{y, vt} z_{vt} \quad (3.87)$$

$$\mathbf{f}_{M_y, mr} = (-\tau_f q + a_v)(K_\beta - \mathbf{f}_T z_{cg}) \quad (3.88)$$

$$\mathbf{f}_{M_y, tr} = \mathbf{f}_{Q_{tr}} \quad (3.89)$$

$$M_{y, ht} = -F_{z, ht} x_{ht} \quad (3.90)$$

$$\mathbf{f}_{M_z, mr} = \mathbf{f}_Q \quad (3.91)$$

$$\mathbf{f}_{M_z, tr} = x_{fus} \mathbf{f}_{T_{tr}} \quad (3.92)$$

$$M_{z, vt} = F_{y, vt} x_{vt} \quad (3.93)$$

$\mathbf{f}_{Q_{tr}}$ in equation 3.89 is the autonomus term of the tail rotor drag torque:

$$\mathbf{f}_{Q_{tr}} = \frac{1}{8} \rho \pi R_{tr}^5 n_{tr}^2 \Omega^2 \sigma_{tr} C_{L\alpha_{tr}} \left(\frac{C_{D0_{tr}}}{C_{L\alpha_{tr}}} (u_{tail}^2 + w_{tail}^2 + 1) - 2v_{tail}^2 \right) \quad (3.94)$$

similarly \mathbf{f}_Q in 3.91 is derived from 3.48.

b_M in equation 3.83 is:

$$\mathbf{b}_M = \begin{bmatrix} \mathbf{b}_{M_{x,mr}} + \mathbf{b}_{M_{x,tr}} \\ \mathbf{b}_{M_{y,mr}} + \mathbf{b}_{M_{y,tr}} \\ \mathbf{b}_{M_{z,mr}} + \mathbf{b}_{M_{z,tr}} \end{bmatrix} \quad (3.95)$$

in which:

$$\mathbf{b}_{M_{x,mr}} = z_{cg}(\tau_{fp} - b_v)b_T + [0 \quad K_{lat}(K_\beta - \mathbf{f}_{T_{mr}}z_{cg}) \quad 0 \quad 0] \quad (3.96)$$

$$\mathbf{b}_{M_{x,tr}} = -z_{fus}b_{T_{tr}} \quad (3.97)$$

$$\mathbf{b}_{M_{y,mr}} = z_{cg}(\tau_{fq} - a_v)b_T + [0 \quad 0 \quad K_{lon}(K_\beta - \mathbf{f}_{T_{mr}}z_{cg}) \quad 0] \quad (3.98)$$

$$\mathbf{b}_{M_{y,tr}} = \mathbf{b}_{Q_{tr}} \quad (3.99)$$

$$\mathbf{b}_{M_{z,mr}} = \mathbf{b}_Q \quad (3.100)$$

$$\mathbf{b}_{M_{z,tr}} = x_{fus}\mathbf{b}_{T_{tr}} \quad (3.101)$$

in equation 3.99 the $b_{Q_{tr}}$ is defined as:

$$\mathbf{b}_{Q_{tr}} = \frac{1}{8} \rho \pi R_{tr}^5 n_{tr}^2 \Omega^2 \sigma_{tr} C_{L\alpha_{tr}} [0 \quad 0 \quad 0 \quad v_{tail}] \quad (3.102)$$

Similarly for expression 3.100 the b_Q is derived from 3.49. if we combine all the state space equations we have:

$$\underbrace{\begin{bmatrix} \dot{V} \\ \dot{\omega} \end{bmatrix}}_{\dot{x}_{6 \times 1}} = \underbrace{\begin{bmatrix} \frac{\mathbf{f}_E - \omega \times V}{I^{-1}(\mathbf{f}_M - \omega \times I\omega)} \end{bmatrix}}_{f_{6 \times 1}} + \underbrace{\begin{bmatrix} \frac{\mathbf{f}_U}{I^{-1}\mathbf{b}_M} \end{bmatrix}}_{f_{6 \times 4}} U \quad (3.103)$$

3.3.3 Control point state space equations

the control point is set to be a point other than CG. This point is a point in the negative direction of z axis in the body coordinates of the helicopter. So by controlling this new control point position and yaw of center of gravity it is possible to control the UAV by using the following set of equations:

$$\underbrace{\begin{bmatrix} \ddot{X}_{CP} \\ \ddot{\psi} \end{bmatrix}}_{\dot{Y}_{4 \times 1}} = \underbrace{\begin{bmatrix} \mathbf{f}_{1_{3 \times 1}} \\ \mathbf{f}_{2_{1 \times 1}} \end{bmatrix}}_{f_{6 \times 1}} + \underbrace{\begin{bmatrix} \mathbf{b}_{1_{3 \times 4}} \\ \mathbf{b}_{2_{1 \times 4}} \end{bmatrix}}_{b_{4 \times 4}} U \quad (3.104)$$

in which:

$$\mathbf{f}_1 = R_b^I \left(\frac{\mathbf{f}_F}{m} + DI^{-1}(\mathbf{f}_M - \omega \times (I\omega)) \right) \quad (3.105)$$

$$\mathbf{f}_2 = \mathbf{f}_s(I^{-1}f_M - I^{-1}\omega \times (I\omega)) + \mathbf{f}_q q + \mathbf{f}_r r; \quad (3.106)$$

In which:

$$\mathbf{f}_q = \dot{\phi} \cos \phi \sec \theta + \dot{\theta} \sin \phi \tan \theta \sec \theta \quad (3.107)$$

$$\mathbf{f}_r = -\dot{\phi} \sin \phi \sec \theta + \dot{\theta} \cos \phi \tan \theta \sec \theta \quad (3.108)$$

$$\mathbf{f}_s = [0 \quad \sec \theta \sin \phi \quad \sec \theta \cos \phi] \quad (3.109)$$

The \mathbf{b}_1 and \mathbf{b}_2 are the control input coefficients in the control point state space equations given as:

$$\mathbf{b}_1 = R_b^I \left(\frac{F_U}{m} + DI^{-1}\mathbf{b}_M \right) \quad (3.110)$$

$$\mathbf{b}_2 = f_s(I^{-1}\mathbf{b}_M); \quad (3.111)$$

3.3.4 Control point position

In this part of the simulation the current position X_{CP} and velocities \dot{X}_{CP} of the control point system Y is calculated using the current states of the UAV center of gravity X_{CG} :

$$X_{CP} = X_{CG} + R_b^I d_B \quad (3.112)$$

In which:

$$d_B = [0, 0, -d] \quad (3.113)$$

$$X_{CG} = [x, y, z] \quad (3.114)$$

d is the distance from the center of gravity to the control point which is 1 meter in this study. First order derivative of the control point position can be derived by:

$$\dot{X}_{CP} = R_b^I (V + (\omega \times d_B)) \quad (3.115)$$

$$\dot{\Theta} = \Phi(\Theta)\omega \quad (3.116)$$

the yaw of the UAV which is the forth row of the Y is the third row of the angular velocity in inertia coordinates:

$$\dot{\Theta} = \Phi(\Theta)\omega \quad (3.117)$$

so we have:

$$\dot{Y} = \left[\frac{\dot{X}_{CP}}{\dot{\psi}} \right]_{4 \times 1} \quad (3.118)$$

$$Y = \left[\frac{X_{CP}}{\psi} \right]_{4 \times 1} \quad (3.119)$$

3.3.5 Implementation of sliding mode controller

By having 3.118 it is now possible to use sliding mode controller on the helicopter:

$$s_i = \dot{y}_i - \dot{y}_{d,i} + \lambda_i y_i - \lambda_i y_{d,i} \quad (3.120)$$

The λ_i are the convergence rates which is supposed to be strictly positive.

$$\lambda = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 2 \end{bmatrix} \quad (3.121)$$

The objective is to control the 3.118 instead of the state space equations of 3.103:

$$\dot{y}_r = \dot{y}_{d,i} - \lambda_i y_i + \lambda_i y_d \quad (3.122)$$

The parameters for sliding mode controller are given next. The first one is the surface reach time given by the following equation:

$$\eta = [1 \quad 1 \quad 1 \quad 1]^T \quad (3.123)$$

The F_e is a vector used as a bound on F:

$$F_e = [10 \quad 10 \quad 5 \quad 5]^T \quad (3.124)$$

The B_e is the bound on b matrix defined by:

$$B_e = \begin{bmatrix} 0.5 & 0 & 0 & 0 \\ 0 & 0.5 & 0 & 0 \\ 0 & 0 & 0.5 & 0 \\ 0 & 0 & 0 & 0.5 \end{bmatrix} \quad (3.125)$$

The error of the controlled states from the desired value is defined by the \tilde{Y} :

$$\tilde{Y} = Y - Y_d \quad (3.126)$$

And the derivative of the \tilde{Y} is as follows:

$$\dot{\tilde{Y}} = \dot{Y} - \dot{Y}_d \quad (3.127)$$

The surface function in the sliding mode controller is formulated as:

$$s_r = \dot{Y}_d - \lambda \tilde{Y} \quad (3.128)$$

So the first order derivative of the s_r can be determined by:

$$\dot{s}_r = \ddot{Y}_d - \lambda \dot{\tilde{Y}} \quad (3.129)$$

E stands for the identity matrix:

$$E = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad (3.130)$$

K is the sliding mode control gain matrix given as:

$$K = (E - B_e)^{-1}(F_e + B_e | -f + \dot{s}_r | + \eta); \quad (3.131)$$

The boundary layer thickness b_s is implemented to remove the chattering problem of the sliding mode controller.

$$b_s = [0.8 \quad 0.8 \quad 1 \quad 1]^T \quad (3.132)$$

So \bar{K} would be the the sliding mode control gain matrix term without the issue of chattering:

$$\bar{K} = \begin{cases} \bar{K}_i = K_i s_i / b_{s_i} & \text{if } |s_i| < b_{s_i} \text{ for } i = 1, 2, 3, 4 \\ \bar{K}_i = K_i \text{sign}(s_i) & \text{if } |s_i| \geq b_{s_i} \text{ for } i = 1, 2, 3, 4 \end{cases} \quad (3.133)$$

As a result, the control input of the helicopter can be derived using the following equation:

$$U = b^{-1}(-f + \dot{s}_r - \bar{K}) \quad (3.134)$$

3.4 Environment setup

Now that we have discussed the dynamic of the helicopter, it is possible to set up the environment suitable for an RL process, which is developed in OpenAI Gym [115]. OpenAI gym is a software development kit for creating and comparing reinforcement learning algorithms. figure 3.2 shows the flowchart of the environment and the dashed line means that the agent is a system outside of the environment. While trying to implement an RL algorithm in a Gym environment, for each episode, first a reset function is called, then the step function is called until a terminal state is reached.

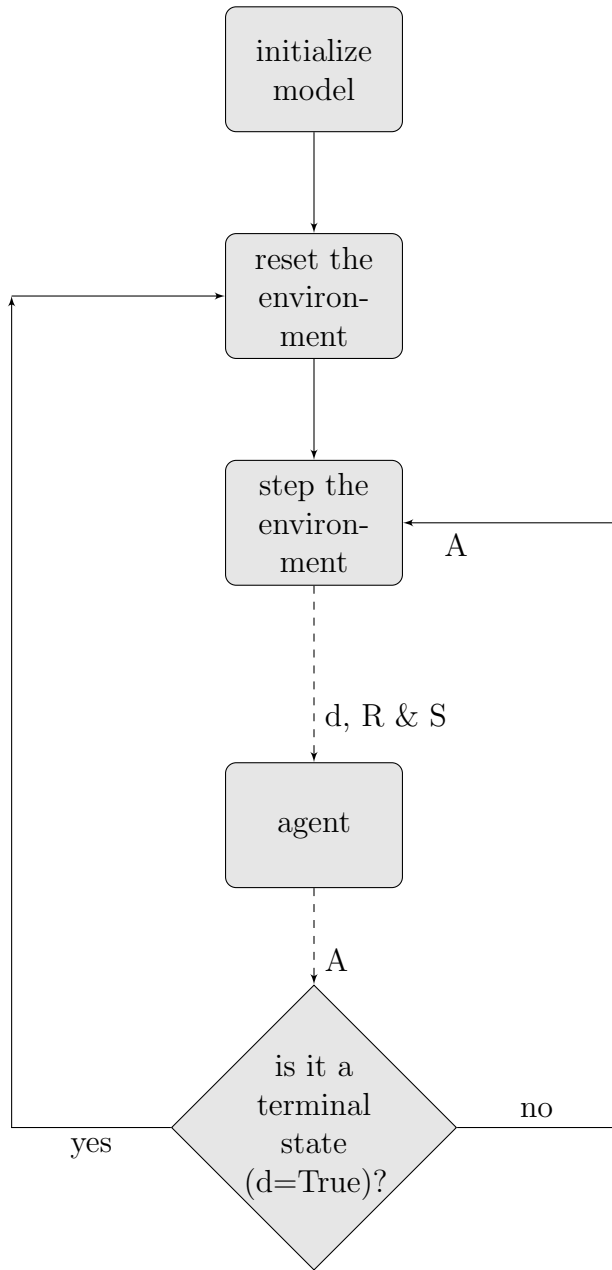


Figure 3.2: Environment Flowchart

In the following sections, the critical points in each part of this environment are discussed.

3.4.1 Environment Reset

Each time the environment is restarted, the helicopter is randomly placed in a position where x , y , and z are uniformly distributed in $[-1, 0, 1]$ so there would be 27 initial states. Other states are kept constant in this phase at hover state.

3.4.2 Step

In each step of the episode, first, the control input is generated from the actions, then the RK45 method is used for solving the set of ODEs. In addition, the reward and the condition of reaching a final state are considered. They are elaborated in the upcoming sections.

Actions

Instead of having the 4 actions as output of the agent, 16 actions are generated by the agent in each step and the control input of the helicopter is find through the following set of equations:

$$\delta_{col} = a_1z + a_2w \tag{3.135}$$

$$\delta_{lat} = a_3y + a_4v + a_5p + a_6\phi \tag{3.136}$$

$$\delta_{lon} = a_7x + a_8v + a_9q + a_{10}\theta \tag{3.137}$$

$$\delta_{ped} = a_{11}r + a_{12}\psi \tag{3.138}$$

This strategy would help the gradient ascent of the agent to find suitable actions for each step more easily.

3.4.3 Observation

The velocity, angular velocity, location, euler angles vectors, and control input of the helicopter are all considered observations in this research. control input is considered an observation since it is not directly generated by the agent.

3.4.4 Reward

The reward function is the most important part of the environment as it provides the goal of the RL algorithm. In this research it consists of 4 terms given as follows:

$$r_t(s) = r_f + r_p + r_\psi + r_u \quad (3.139)$$

In the above equation, r_f is the flying reward term, r_p represents the position reward term, r_ψ denotes the yaw angle reward term and the r_u stands for the control input reward term.

Flying term

Flying reward r_f is just a constant (18.8 in this case), assures that the algorithm is rewarded for longer episodes. The absence of this term will lead to a local minimum of reward in which the agent tries to end the episode to stop receiving negative rewards by crashing the helicopter. It also helps to stabilize the UAV in the long term.

Position term

The position error r_p punishes the agent for the distance between the current position of the UAV and the origin:

$$r_p(t) = -10\|X(t)\|_2 \quad (3.140)$$

Yaw angle term

This term also similarly punishes the agent for the error of ψ :

$$r_\psi(t) = -0.25|\psi(t)| \quad (3.141)$$

Control input terms

The control input terms consist of a derivative and a norm term to reduce chattering and increase energy consumption of the UAV:

$$r_u = -0.015\|U\| - 0.08\|U'\| \quad (3.142)$$

3.4.5 Checking for a terminal state

Unless the helicopter crashes or 8 seconds have passed, it is not a terminal state. Crashing in this research is when the states are outside of the $[-100, 100]$, except for the Euler angles which the bounds are $\phi \in [-\pi, \pi]$, $\theta \in [-\pi/2, \pi/2]$ and $\psi \in [-2\pi, 2\pi]$.

3.4.6 Summary

In this section, the dynamics for 6-DOF nonlinear dynamics of a small-scale UAV is provided. It included the effect of the fuselage, main rotor, tail rotor, etc. The setup of the environment is explained, and the code is given in Appendix A. The procedure to implement the actions and rewards in this research is also explained in detail. The implementation specifics of the SAC algorithm in this context are elaborated on in the next chapter, and the results are analyzed.

Chapter 4

Result and discussion

4.1 SAC agent

In order to solve the helicopter environment, as presented in the previous chapter, we implemented the SAC algorithm. In this section, we have provided the implementation of the soft actor-critic as a controller for the helicopter which is shown in figure 4.1. We have implemented 5 other reinforcement learning methods such as D4PG [116], proximal policy optimization (PPO) [72], Trust Region Policy Optimization (TRPO) [117], deep deterministic policy gradient [118] and Twin Delayed DDPG [100] and we were unable to find an stabilized performance of the helicopter using the aforementioned algorithms.

For this study, we use garage [119] as an API for the agent. Garage implements state-of-the-art deep reinforcement learning algorithms in Python and coherently integrates with the deep learning library PyTorch [120] and Tensorflow [121]. The library provides a straightforward approach to evaluate and test different algorithms in Gym environments. The schematic of agent-environment interaction is illustrated in Figure 4.1.

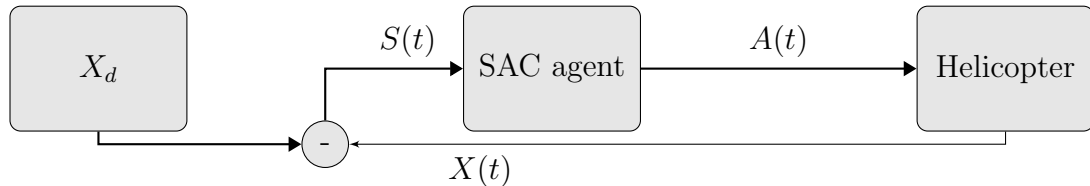


Figure 4.1: SAC controller schematic

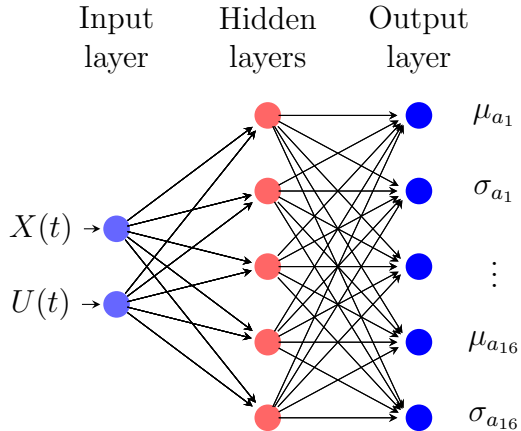


Figure 4.2: Actor diagram for SAC agent, the inputs are the 16 states of the helicopter $X(t)$ and the 4 control inputs $U(t)$, the outputs are the average μ and the standard deviation σ of the agent actions

4.1.1 Architecture

In this section, we discuss the architecture of the SAC actor and critic neural network.

Actor

The actor diagram and schematic of this SAC agent is given in figures 4.2 and 4.4 which includes 2 hidden layer of size 128 and 128 fully connected layers with a Rectified Linear Unity (ReLU) activation function and a \tanh activation function at last in order to narrow the result to $[-1,1]$, the actions are then linearly mapped to the action range based on the environment*. In order to constraint the standard deviation of the policy, it is set to be between $[e^{-20}, e^1]$.

Critic

The diagram for the Critic neural network is depicted in figs. 4.3 and 4.4. The diagram includes 2 hidden layer of size 256 and 256 fully connected layers with ReLU activation function after each hidden layers.

4.1.2 Hyper-parameters

Optuna package is utilized for optimization of all the hyperparameters of the SAC agent [122]. The SAC agent’s primary hyper-parameters are given in table 4.1 as part

*check action wrapper at Appendix A

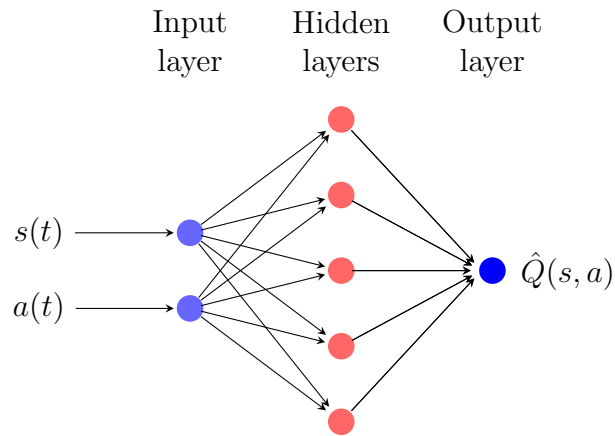


Figure 4.3: Critic diagram for SAC agent, inputs are the states and actions while the output is the Q value for the given input.

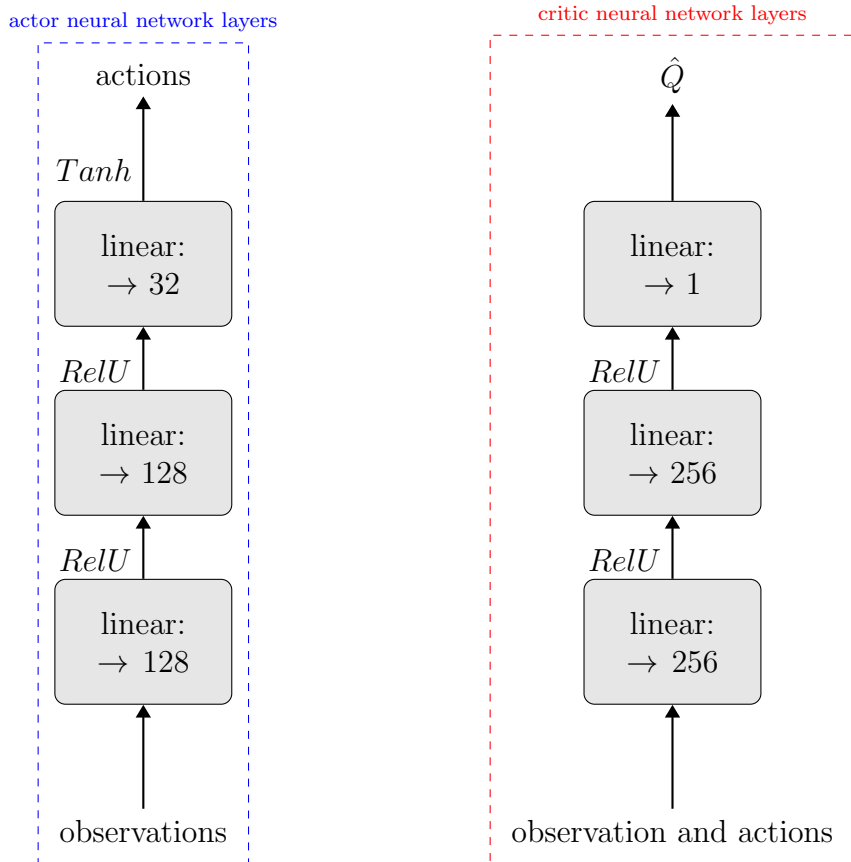


Figure 4.4: Actor and critic neural network diagram

of its implementation.

Table 4.1: Hyper parameters of SAC agent.

Hyper parameter	Value	description
\mathcal{T}	5×10^7	Total number of environment steps.
\mathcal{B}	2048	mini batch which is the number of samples from the buffer randomly sampled for each stochastic gradient decent step update
\mathcal{D}	10^7	replay buffer size, which is the total number of steps saved in buffer (when new ones are added the oldest ones are removed.)
l_κ	3×10^{-4}	learning rate for optimization of policy by this factor.
l_δ	3×10^{-4}	learning rate for optimization of Q functions.
τ_{target}	5×10^{-3}	updating the target network linearly by this factor.
γ	0.99	factor for discounting later rewards.
α	$3 \times e^{-0.009i}$	the temperature term in SAC policy.
χ	3×10^{-3}	The learning rate of Adam optimizer.
σ_{mean}	e^{-20}	the minimum standard deviation of the actions in the stochastic policy of SAC.
σ_{mean}	e^1	the maximum standard deviation of the actions in the stochastic policy of SAC
i_{max}	10^4	the maximum number of iteration

The replay buffer is designed to contain 10^7 transition tuples before starting the SAC algorithm. The experience replay buffer enables learning from prior policy experiences while avoiding correlated samples in the gradient step. Furthermore, we implement an upgraded target network with a target factor of 5×10^{-3} . This is mo-

tivated by the desire to improve the stability of the learning process. Based on line 21 of algorithm 2, if it is time to update, the gradient step per epoch is set to 2 and the gradient step per iteration is set to 8. For each time a gradient step is executed, a mini-batch of 2048 random samples is chosen from the replay buffer. The learning rate for Adam optimizer on both neural network is set to be 3×10^{-3} .

The training of the network included about 5×10^7 to 6×10^7 environment time steps of about 0.03, simulation time. The simulation is run in a 32-cores Intel(R) Xeon(R) CPU E5-2620 v4 @ 2.10GHz which took about 3–4 days.

4.2 Results

4.2.1 Training and evaluation

The training result is given in 4.5. The iteration is continued for 10,000 iterations; however, no significant improvement is found after 6000. The standard deviation increases as the number of iterations increases. This is acceptable because at first, in all the episodes, the simulated UAV crashes. However, as the training progresses, the agent improves its action to stabilize the helicopter and hence the difference between the return of points closer to the origin and those placed at a more distant point from the origin grows.

4.2.2 Comparison of Controllability and stability to the SMC

In order to be able to compare the results obtained by an RL method to a traditional control method, the results of controlling the initial point set to $[-1,-1,-1]$ are shown in figs. 4.6 to 4.11 for both the sliding mode controller (SMC) and the SAC policy obtained. As seen, the resulting policy achieves good stabilization capability.

The settling time is considered to be the time when the states reach 0.10 m of the origin which is the desired state. The rise time in this study is considered to be the time for the response to rise from the absolute value of 0.9 m to 0.1 m in the vicinity of origin for each x, y and z states. Considering the aforementioned definitions the comparison between the SAC and SMC is given in table 4.2.

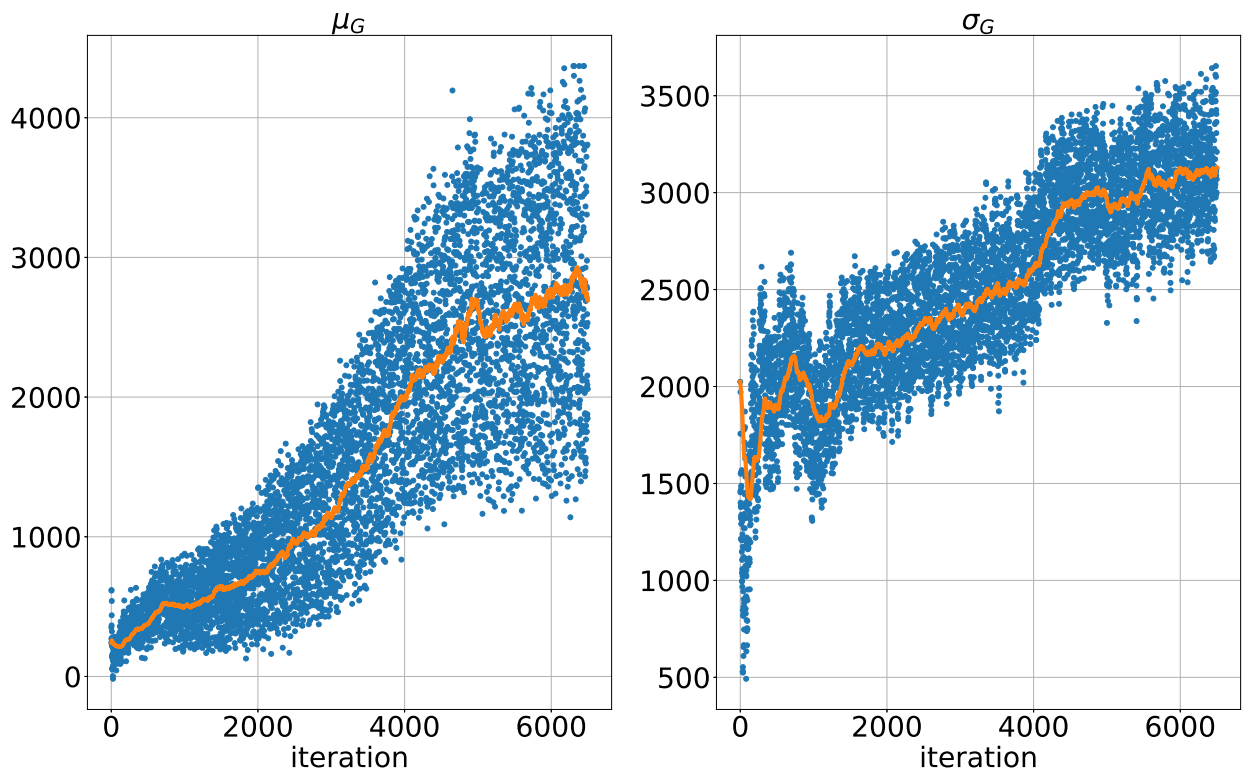


Figure 4.5: Averaged discounted return μ_G and standard deviation σ_G of each iteration using the random actions of the policy.

Table 4.2: Comparison of SAC and SMC based on the response characteristics of the helicopter dynamic system for the case of initial point set to [-1,-1,-1].

	settling time [s]	rise time [s]	overshoot [m]	SS error [m]
x,SMC	5.07	0.79	0.13	0.0
x,SAC	4.41	3.96	0	0
y,SMC	1.36	1.0	0.09	0.05
y,SAC	2.96	2.66	0	0.05
z,SMC	1.51	1.47	0.03	0
z,SAC	2.6	2.45	0.1	0.1

Based on the given comparison result on 4.2 the sliding mode controller provides a better result in the case of x y and z. The ψ angle has a 10.8° steady-state error which is not superior to the sliding mode controller (1.8°). However, it should be mentioned that the sliding mode controller is a highly tuned controller for this system. On the other hand, the SAC agent is a model-free method in which it did not have any access to the model.

The result of the control input is given in figure 4.10. There are some vibrations in the δ_{ped} and δ_{col} , We find that it was somewhat hard to reduce these vibrations because as we increased the control derivative input term in the 3.142, the policy would alternate between getting closer to the target and achieving a stable hovering somewhere far from the origin.

robustness

In order to test the robustness of the policy, a simulated wind is blown at the UAV given the following equation:

$$V_{wind,t} = V_{wind,t-1} + W \quad (4.1)$$

in which W is a random number in [-0.1,0.1] at each time step and using the policy generated. It achieved 100% stability in all the 27 initial positions, In order to compare the results of the simulation with and without the wind, the rise time, settling time, and steady-state errors are given in table 4.3 for the [-1,-1,-1] case.

Table 4.3: Robustness response characteristics of the helicopter dynamic system for the SAC policy by a simulated wind for the case of initial point set to $[-1,-1,-1]$.

	settling time [s]	rise time [s]	overshoot [m]	SS error [m]
x	4.29	3.87	0	0.03
y	2.87	2.57	0.08	0.02
z	2.47	2.36	0.09	0.09

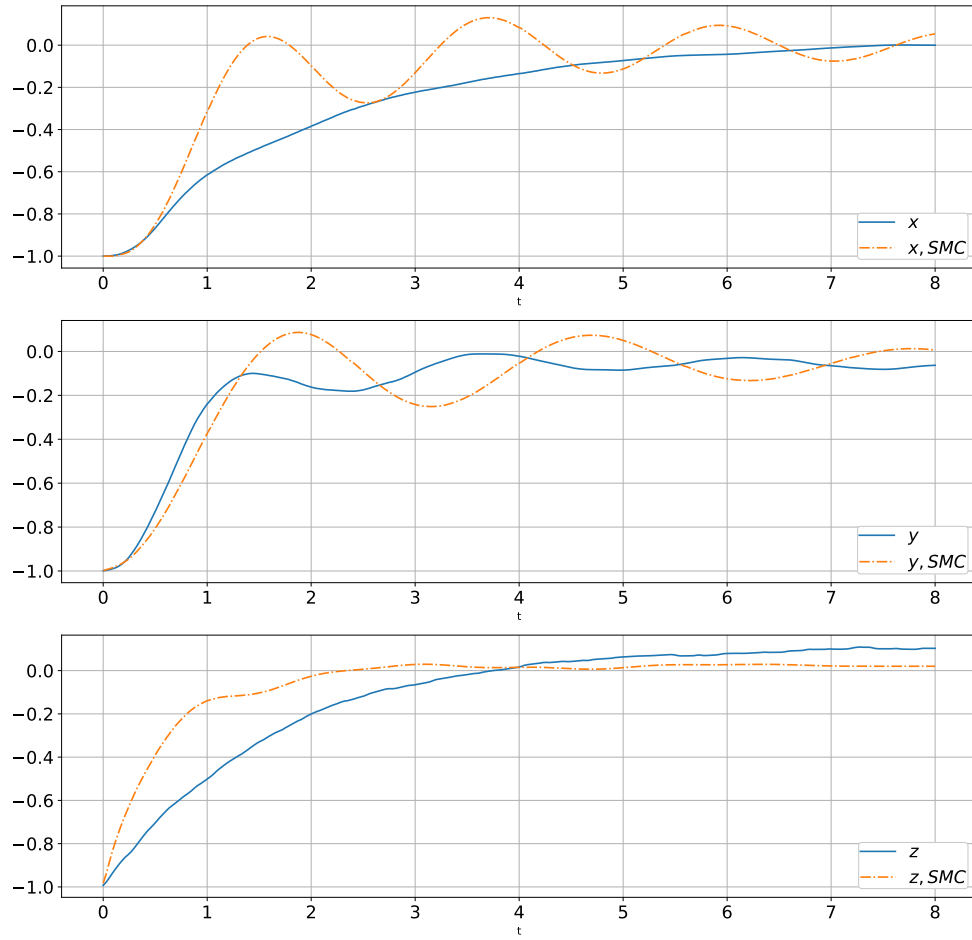


Figure 4.6: Helicopter SMC and SAC positional states by initial position $[-1,-1,-1]$.

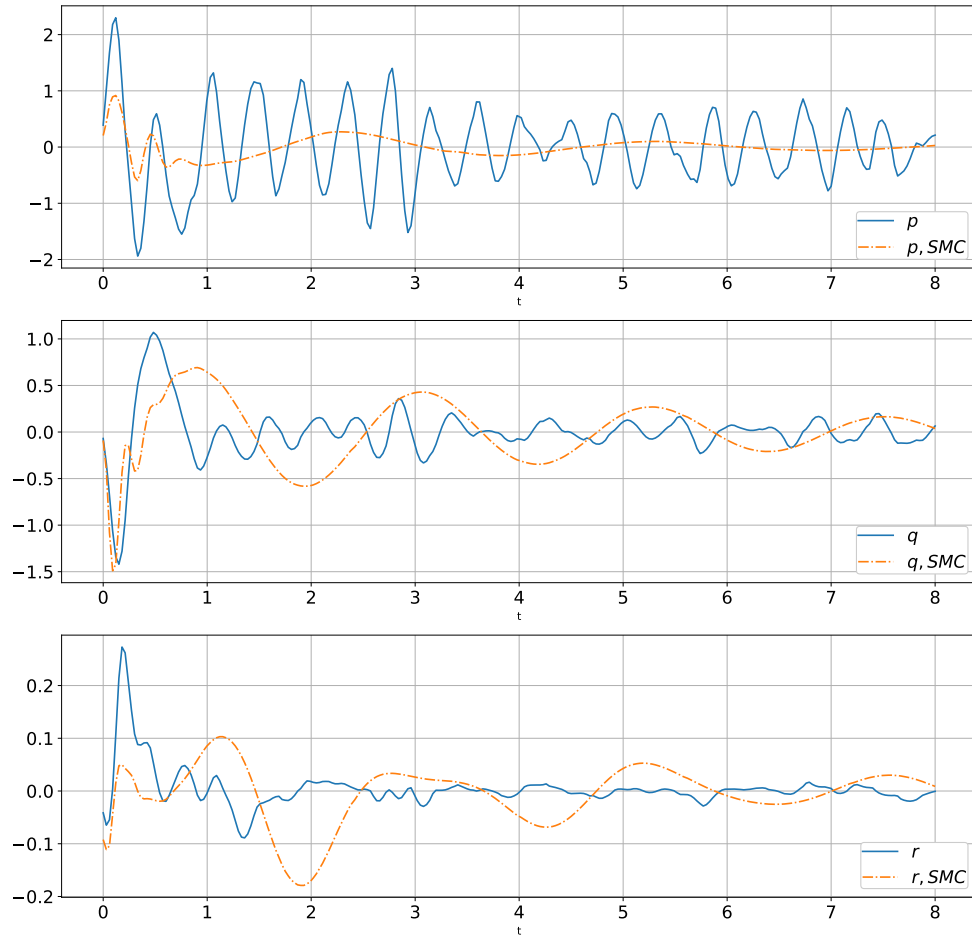


Figure 4.7: Helicopter SMC and SAC angular velocities by initial position $[-1,-1,-1]$.

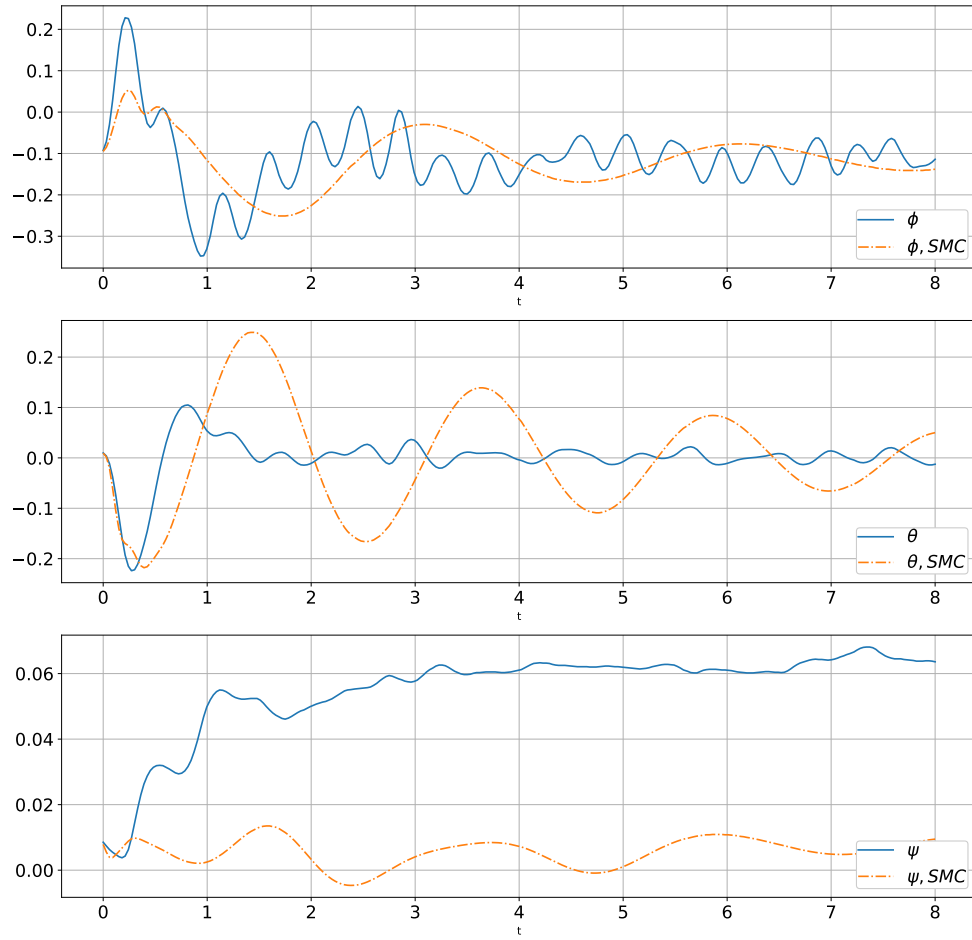


Figure 4.8: Helicopter SMC and SAC Euler angles by initial position $[-1,-1,-1]$.

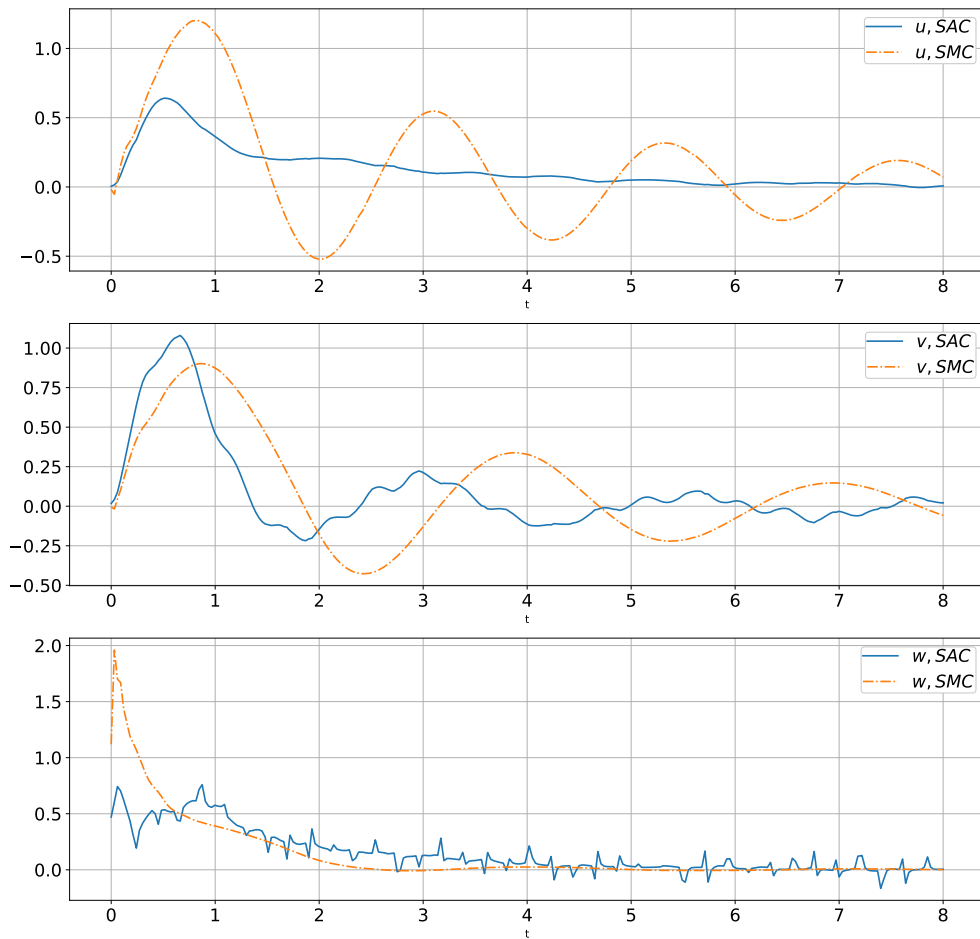


Figure 4.9: Helicopter SMC and SAC velocities by initial position $[-1,-1,-1]$.

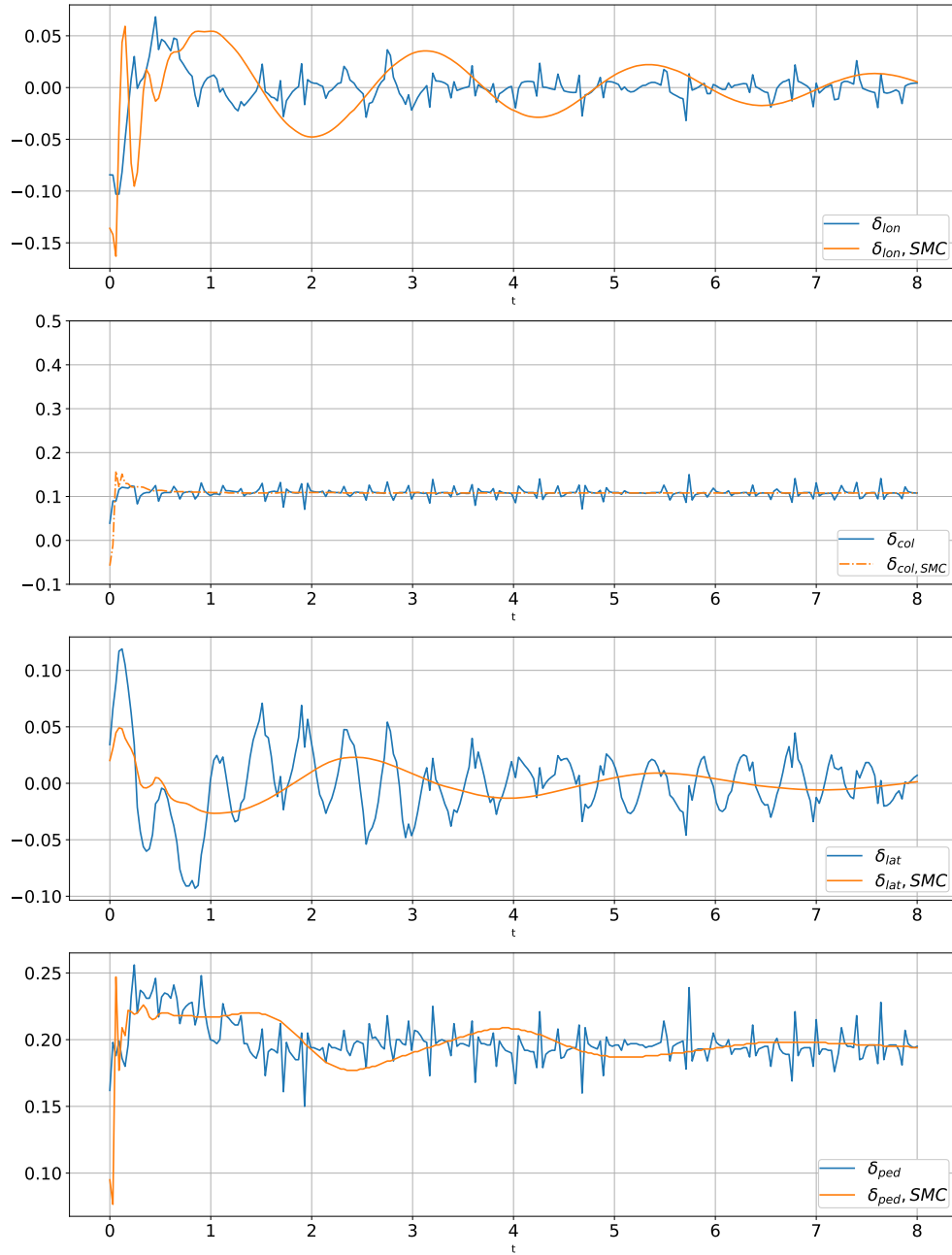


Figure 4.10: Helicopter SMC and SAC control input by initial position $[-1,-1,-1]$.

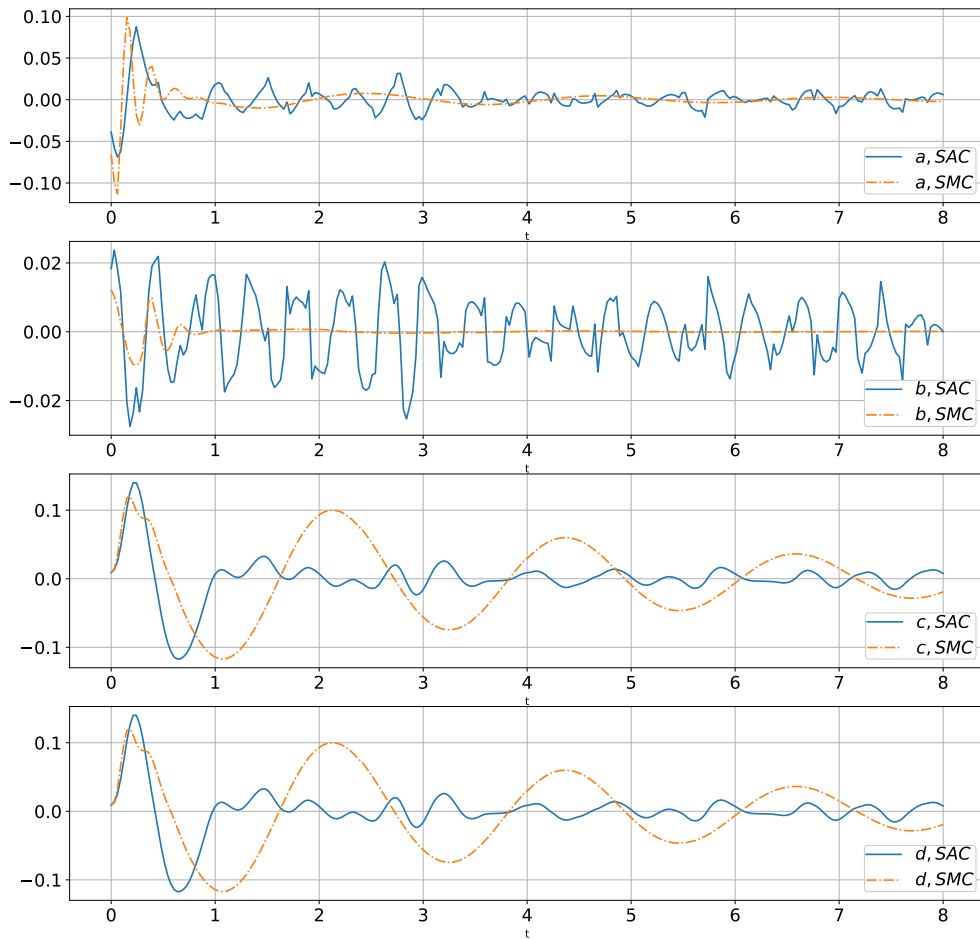


Figure 4.11: Helicopter SMC and SAC flapping states by initial position $[-1,-1,-1]$.

Chapter 5

Conclusions and Future Directions

5.1 Conclusions

This study shows how to train a reinforcement learning agent using a model-free off-policy technique, specifically the Soft Actor-Critic algorithm, to produce a policy capable of performing low-level control of a simulated small-scaled helicopter. The use of this method for the same task has never been disclosed previously. The result of the SAC method is compared to an sliding mode controller, although the result was not superior, having in mind that the SAC method did not have access to the dynamics of the simulated helicopter, it provided a promising result in which the average steady state error was 0.05% for the given $[-1,-1,-1]$ case while the sliding mode controller provided an almost 0 error on the steady states. We also assessed the policy in an environment with the random wind as a disturbance to test the robustness of the method, and it is demonstrated that the SAC technique was capable of achieving stability in all trials.

5.2 Future work

The comparison of the results in the previous section showed that there is still room for improvement in case if stability criteria of the achieved policy, one way to improve it is to work in more detail about the reward function of the policy, another procedure would be to improve the RL algorithms by modifications as the RL field is improving day by day. Although it was demonstrated here that the small-sized helicopter could be stabilized using the SAC approach, trajectory tracking and recovery operations

were not conducted in this study and can be addressed in future studies. In [73], similar study was conducted on a quadcopter.

The ability to efficiently apply deep RL algorithms to the real world to address practical applications may be the most compelling motivator for future advances in the area. This study showed that RL is capable of controlling the helicopter; however, this has been done in the simulation environment, future studies could be focused on using such policies in real-world data. A review of similar approaches may be found in [77].

There is a possibility of a relatively large gap between the simulation environment and the real-world data; a possible moderator would be to take advantage of a more sophisticated model such as the one for Yamaha R-50 helicopter [16, 17] to improve the replication of the environment.

Bibliography

- [1] Kimon P Valavanis and George J Vachtsevanos. *Handbook of unmanned aerial vehicles*, volume 2077. Springer, 2015.
- [2] Eduard Semsch, Michal Jakob, Dušan Pavlicek, and Michal Pechoucek. Autonomous uav surveillance in complex urban environments. In *2009 IEEE/WIC/ACM International Joint Conference on Web Intelligence and Intelligent Agent Technology*, volume 2, pages 82–85. IEEE, 2009.
- [3] Anuj Puri. A survey of unmanned aerial vehicles (uav) for traffic surveillance. *Department of computer science and engineering, University of South Florida*, pages 1–29, 2005.
- [4] Iván Maza, Fernando Caballero, Jesús Capitán, José Ramiro Martínez-de Dios, and Aníbal Ollero. Experimental results in multi-uav coordination for disaster management and civil security applications. *Journal of intelligent & robotic systems*, 61(1):563–585, 2011.
- [5] Andreas Birk, Burkhard Wiggerich, Heiko Bülow, Max Pflingsthorn, and Sören Schwertfeger. Safety, security, and rescue missions with an unmanned aerial vehicle (uav). *Journal of Intelligent & Robotic Systems*, 64(1):57–76, 2011.
- [6] Ebtehal Turki Alotaibi, Shahad Saleh Alqefari, and Anis Koubaa. Lsar: Multi-uav collaboration for search and rescue missions. *IEEE Access*, 7:55817–55832, 2019.
- [7] Teppo Luukkonen. Modelling and control of quadcopter. *Independent research project in applied mathematics, Espoo*, 22:22, 2011.
- [8] Daniel Gheorghiuță, Ionuț Vîntu, Letiția Mirea, and Cătălin Brăescu. Quadcopter control system. In *2015 19th International Conference on System Theory, Control and Computing (ICSTCC)*, pages 421–426. IEEE, 2015.
- [9] Pengcheng Wang, Zhihong Man, Zhenwei Cao, Jinchuan Zheng, and Yong Zhao. Dynamics modelling and linear control of quadcopter. In *2016 International Conference on Advanced Mechatronic Systems (ICAMEchS)*, pages 498–503. IEEE, 2016.

- [10] Omar I Dallah Bashi, WZ Hasan, N Azis, S Shafie, and Hiroaki Wagatsuma. Unmanned aerial vehicle quadcopter: A review. *Journal of Computational and Theoretical Nanoscience*, 14(12):5663–5675, 2017.
- [11] Quan Quan. *Introduction to multicopter design and control*. Springer, 2017.
- [12] Beibei Ren, Shuzhi Sam Ge, Chang Chen, Cheng-Heng Fua, and Tong Heng Lee. *Modeling, control and coordination of helicopter systems*. Springer Science & Business Media, 2012.
- [13] JC Avila Vilchis, Bernard Brogliato, Alejandro Dzul, and Rogelio Lozano. Non-linear modelling and control of helicopters. *Automatica*, 39(9):1583–1596, 2003.
- [14] T John Koo and Shankar Sastry. Output tracking control design of a helicopter model based on approximate linearization. In *Proceedings of the 37th IEEE Conference on Decision and Control (Cat. No. 98CH36171)*, volume 4, pages 3635–3640. IEEE, 1998.
- [15] Robert Mahony, Tarek Hamel, and A Dzul. Hover control via lyapunov control for an autonomous model helicopter. In *Proceedings of the 38th IEEE Conference on Decision and Control (Cat. No. 99CH36304)*, volume 4, pages 3490–3495. IEEE, 1999.
- [16] Marco La Civita, George Papageorgiou, William C Messner, and Takeo Kanade. Integrated modeling and robust control for full-envelope flight of robotic helicopters. In *2003 IEEE International Conference on Robotics and Automation (Cat. No. 03CH37422)*, volume 1, pages 552–557. IEEE, 2003.
- [17] M La Civita, George Papageorgiou, William C Messner, and Takeo Kanade. Design and flight testing of an h00 controller for a robotic helicopter. *Journal of Guidance, Control, and Dynamics*, 29(2):485–494, 2006.
- [18] Sepehr Pourrezaei Khaligh. *Control-oriented modeling and system identification for nonlinear trajectory tracking control of a small-scale unmanned helicopter*. University of Alberta (Canada), 2014.
- [19] DJ Walker and I Postlethwaite. Advanced helicopter flight control using two-degree-of-freedom h (infinity) optimization. *Journal of Guidance, Control, and Dynamics*, 19(2):461–468, 1996.
- [20] Bernard Mettler, Mark B Tischler, and Takeo Kanade. System identification of small-size unmanned helicopter dynamics. In *Annual Forum Proceedings-American Helicopter Society*, volume 2, pages 1706–1717. Citeseer, 1999.
- [21] H Jin Kim and David H Shim. A flight control system for aerial robots: algorithms and experiments. *Control engineering practice*, 11(12):1389–1400, 2003.

- [22] Eric N Johnson and Suresh K Kannan. Adaptive trajectory control for autonomous helicopters. *Journal of Guidance, Control, and Dynamics*, 28(3):524–538, 2005.
- [23] Lorenzo Marconi and Roberto Naldi. Robust full degree-of-freedom tracking control of a helicopter. *Automatica*, 43(11):1909–1920, 2007.
- [24] Ioannis A Raptis, Kimon P Valavanis, and Wilfrido A Moreno. System identification and discrete nonlinear control of miniature helicopters using backstepping. *Journal of Intelligent and Robotic Systems*, 55(2):223–243, 2009.
- [25] Miroslav Krstic, Petar V Kokotovic, and Ioannis Kanellakopoulos. *Nonlinear and adaptive control design*. John Wiley & Sons, Inc., 1995.
- [26] Isabelle Fantoni, Rogelio Lozano, and SC Sinha. Non-linear control for under-actuated mechanical systems. *Appl. Mech. Rev.*, 55(4):B67–B68, 2002.
- [27] A Azzam and Xinhua Wang. Quad rotor arial robot dynamic modeling and configuration stabilization. In *2010 2nd International Asia Conference on Informatics in Control, Automation and Robotics (CAR 2010)*, volume 1, pages 438–444. IEEE, 2010.
- [28] Robert Mahony and Tarek Hamel. Robust trajectory tracking for a scale model autonomous helicopter. *International Journal of Robust and Nonlinear Control: IFAC-Affiliated Journal*, 14(12):1035–1059, 2004.
- [29] Bryan Godbolt. *Experimental nonlinear control of a helicopter unmanned aerial vehicle (UAV)*. University of Alberta (Canada), 2013.
- [30] Islam SM Khalil, JC Doyle, and K Glover. *Robust and optimal control*. prentice hall, new jersey, 1996.
- [31] Farid Kendoul, David Lara, Isabelle Fantoni, and R Lozano. Real-time non-linear embedded control for an autonomous quadrotor helicopter. *Journal of guidance, control, and dynamics*, 30(4):1049–1061, 2007.
- [32] H Jin Kim, David H Shim, and Shankar Sastry. Nonlinear model predictive tracking control for rotorcraft-based unmanned aerial vehicles. In *Proceedings of the 2002 American control conference (IEEE Cat. No. CH37301)*, volume 5, pages 3576–3581. IEEE, 2002.
- [33] Jyotirmay Gadewadikar, Frank Lewis, Kamesh Subbarao, and Ben M Chen. Structured h-infinity command and control-loop design for unmanned helicopters. *Journal of guidance, control, and dynamics*, 31(4):1093–1102, 2008.
- [34] Bernard Mettler. *Identification modeling and characteristics of miniature rotorcraft*. Springer Science & Business Media, 2013.

- [35] Alexander Bogdanov and Eric Wan. State-dependent riccati equation control for small autonomous helicopters. *Journal of guidance, control, and dynamics*, 30(1):47–60, 2007.
- [36] H Ifassiouen, M Guisser, and H Medromi. Robust nonlinear control of a miniature autonomous helicopter using sliding mode control structure. *World Academy of Science, Engineering and Technology*, 2:1, 2007.
- [37] JK Pieper. Application of slmc: Trc control of a helicopter in hover. In *Proceedings of 1995 American Control Conference-ACC'95*, volume 2, pages 1191–1195. IEEE, 1995.
- [38] Wei Wang, Kenzo Nonami, and Yuta Ohira. Model reference sliding mode control of small helicopter xrb based on vision. *International Journal of Advanced Robotic Systems*, 5(3):26, 2008.
- [39] Jian Fu, Wen-hua Chen, and Qing-xian Wu. Chattering-free sliding mode control with unidirectional auxiliary surfaces for miniature helicopters. *International Journal of Intelligent Computing and Cybernetics*, 2012.
- [40] Farbod Fahimi. Full formation control for autonomous helicopter groups. *Robotica*, 26(2):143–156, 2008.
- [41] Yunjun Xu. Multi-timescale nonlinear robust control for a miniature helicopter. *IEEE Transactions on Aerospace and Electronic systems*, 46(2):656–671, 2010.
- [42] Alberto Isidori, Lorenzo Marconi, and Andrea Serrani. Robust nonlinear motion control of a helicopter. In *Robust Autonomous Guidance*, pages 149–192. Springer, 2003.
- [43] T John Koo and Shankar Sastry. Differential flatness based full authority helicopter control design. In *Proceedings of the 38th IEEE Conference on Decision and Control (Cat. No. 99CH36304)*, volume 2, pages 1982–1987. IEEE, 1999.
- [44] Alberto Isidori, Lorenzo Marconi, and Andrea Serrani. *Robust autonomous guidance: an internal model approach*. Springer Science & Business Media, 2012.
- [45] Bourhane Kadmiry and Dimiter Driankov. A fuzzy flight controller combining linguistic and model-based fuzzy control. *Fuzzy Sets and Systems*, 146(3):313–347, 2004.
- [46] James Downs, Ron Prentice, S Alzell, Adam Besachio, CM Ivler, Mark B Tischler, and MH Mansur. Control system development and flight test experience with the mq-8b fire scout vertical take-off unmanned aerial vehicle (vtuav). In *Annual Forum Proceedings-American Helicopter Society*, volume 63, page 566. AMERICAN HELICOPTER SOCIETY, INC, 2007.

- [47] Ioannis A Raptis and Kimon P Valavanis. *Linear and nonlinear control of small-scale unmanned helicopters*, volume 1. Springer, 2011.
- [48] Benchun Zhou, Weihong Wang, Zhenghua Liu, and Jia Wang. Vision-based navigation of uav with continuous action space using deep reinforcement learning. In *2019 Chinese Control And Decision Conference (CCDC)*, pages 5030–5035. IEEE, 2019.
- [49] Naira Hovakimyan, Nakwan Kim, Anthony J Calise, JVR Prasad, and Eric Corban. Adaptive output feedback for high-bandwidth control of an unmanned helicopter. In *AIAA Guidance, Navigation and Control Conference, AIAA-2001-4181*, 2001.
- [50] Nakwan Kim, Anthony Calise, J Eric Corban, and JVR Prasad. Adaptive output feedback for altitude control of an unmanned helicopter using rotor rpm. In *AIAA Guidance, Navigation, and Control Conference and Exhibit*, page 5323, 2004.
- [51] Russell Enns and Jennie Si. Helicopter trimming and tracking control using direct neural dynamic programming. *IEEE transactions on neural networks*, 14(4):929–939, 2003.
- [52] S Anand and G Verweij. What’s the real value of ai for your business and how can you capitalise, 2019.
- [53] Manjusha Pandey. *Machine Learning: Theoretical Foundations and Practical Applications*. Springer Nature, 2021.
- [54] Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.
- [55] Andrew Ng. Machine learning yearning. URL: [http://www. mlyearning.org/\(96\)](http://www.mlyearning.org/(96)), 139, 2017.
- [56] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [57] Rui Nian, Jinfeng Liu, and Biao Huang. A review on reinforcement learning: Introduction and applications in industrial process control. *Computers & Chemical Engineering*, 139:106886, 2020.
- [58] Warren B Powell. Ai, or and control theory: A rosetta stone for stochastic optimization. *Princeton University*, page 12, 2012.
- [59] Lukasz Kaiser, Mohammad Babaeizadeh, Piotr Milos, Blazej Osinski, Roy H Campbell, Konrad Czechowski, Dumitru Erhan, Chelsea Finn, Piotr Koza-kowski, Sergey Levine, et al. Model-based reinforcement learning for atari. *arXiv preprint arXiv:1903.00374*, 2019.

- [60] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- [61] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without human knowledge. *nature*, 550(7676):354–359, 2017.
- [62] Pieter Abbeel, Adam Coates, Morgan Quigley, and Andrew Y Ng. An application of reinforcement learning to aerobatic helicopter flight. *Advances in neural information processing systems*, 19:1, 2007.
- [63] Andrew Y Ng, Adam Coates, Mark Diel, Varun Ganapathi, Jamie Schulte, Ben Tse, Eric Berger, and Eric Liang. Autonomous inverted helicopter flight via reinforcement learning. In *Experimental robotics IX*, pages 363–372. Springer, 2006.
- [64] Andrew Y Ng, H Jin Kim, Michael I Jordan, Shankar Sastry, and Shiv Bal-lianda. Autonomous helicopter flight via reinforcement learning. In *NIPS*, volume 16. Citeseer, 2003.
- [65] Sridhar Mahadevan and Jonathan Connell. Automatic programming of behavior-based robots using reinforcement learning. *Artificial intelligence*, 55(2-3):311–365, 1992.
- [66] Vijaykumar Gullapalli, Judy A Franklin, and Hamid Benbrahim. Acquiring robot skills via reinforcement learning. *IEEE Control Systems Magazine*, 14(1):13–24, 1994.
- [67] J Andrew Bagnell and Jeff G Schneider. Autonomous helicopter control using reinforcement learning policy search methods. In *Proceedings 2001 ICRA. IEEE International Conference on Robotics and Automation (Cat. No. 01CH37164)*, volume 2, pages 1615–1620. IEEE, 2001.
- [68] Yuxi Li. Deep reinforcement learning: An overview. *arXiv preprint arXiv:1701.07274*, 2017.
- [69] William Koch, Renato Mancuso, Richard West, and Azer Bestavros. Reinforcement learning for uav attitude control. *ACM Transactions on Cyber-Physical Systems*, 3(2):1–21, 2019.
- [70] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- [71] John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *International conference on machine learning*, pages 1889–1897. PMLR, 2015.

- [72] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms, 2017.
- [73] Gabriel Moraes Barros and Esther Luna Colombini. Using soft actor-critic for low-level uav control. *arXiv preprint arXiv:2010.02293*, 2020.
- [74] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor, 2018.
- [75] Sergey Levine, Peter Pastor, Alex Krizhevsky, Julian Ibarz, and Deirdre Quillen. Learning hand-eye coordination for robotic grasping with deep learning and large-scale data collection. *The International Journal of Robotics Research*, 37(4-5):421–436, 2018.
- [76] Jens Kober, J Andrew Bagnell, and Jan Peters. Reinforcement learning in robotics: A survey. *The International Journal of Robotics Research*, 32(11):1238–1274, 2013.
- [77] Wenshuai Zhao, Jorge Peña Queraltá, and Tomi Westerlund. Sim-to-real transfer in deep reinforcement learning for robotics: a survey. In *2020 IEEE Symposium Series on Computational Intelligence (SSCI)*, pages 737–744. IEEE, 2020.
- [78] Loc D Tran, Charles D Cross, Mark A Motter, James H Neilan, Garry Qualls, Paul M Rothhaar, Anna Trujillo, and Bonnie D Allen. Reinforcement learning with autonomous small unmanned aerial vehicles in cluttered environments—”after all these years among humans, you still haven’t learned to smile.”. In *15th AIAA aviation technology, integration, and operations conference*, page 2899, 2015.
- [79] Eric Tzeng, Judy Hoffman, Trevor Darrell, and Kate Saenko. Simultaneous deep transfer across domains and tasks. In *Proceedings of the IEEE international conference on computer vision*, pages 4068–4076, 2015.
- [80] Jae Won Lee. Stock price prediction using reinforcement learning. In *ISIE 2001. 2001 IEEE International Symposium on Industrial Electronics Proceedings (Cat. No. 01TH8570)*, volume 1, pages 690–695. IEEE, 2001.
- [81] Chen-Huan Pi, Kai-Chun Hu, Stone Cheng, and I-Chen Wu. Low-level autonomous control and tracking of quadrotor using reinforcement learning. *Control Engineering Practice*, 95:104222, 2020.
- [82] Leslie Pack Kaelbling, Michael L Littman, and Andrew W Moore. Reinforcement learning: A survey. *Journal of artificial intelligence research*, 4:237–285, 1996.
- [83] Dimitri P Bertsekas. Dynamic programming and optimal control 3rd edition, volume ii. *Belmont, MA: Athena Scientific*, 2011.

- [84] Mark E Lewis and Martin L Puterman. Bias optimality. In *Handbook of Markov decision processes*, pages 89–111. Springer, 2002.
- [85] Lucian Buşoniu, Bart De Schutter, and Robert Babuška. Approximate dynamic programming and reinforcement learning. In *Interactive collaborative information systems*, pages 3–44. Springer, 2010.
- [86] Victor Gabillon, Mohammad Ghavamzadeh, and Bruno Scherrer. Approximate dynamic programming finally performs well in the game of tetris. In *Neural Information Processing Systems (NIPS) 2013*, 2013.
- [87] Gavin A Rummery and Mahesan Niranjan. *On-line Q-learning using connectionist systems*, volume 37. Citeseer, 1994.
- [88] Christopher JCH Watkins and Peter Dayan. Q-learning. *Machine learning*, 8(3-4):279–292, 1992.
- [89] Richard S Sutton. Learning to predict by the methods of temporal differences. *Machine learning*, 3(1):9–44, 1988.
- [90] Csaba Szepesvári. Algorithms for reinforcement learning. *Synthesis lectures on artificial intelligence and machine learning*, 4(1):1–103, 2010.
- [91] Leemon C Baird III. Reinforcement learning through gradient descent. Technical report, CARNEGIE-MELLON UNIV PITTSBURGH PA DEPT OF COMPUTER SCIENCE, 1999.
- [92] Hajime Kimura, Masayuki Yamamura, and Shigenobu Kobayashi. Reinforcement learning by stochastic hill climbing on discounted reward. In *Machine Learning Proceedings 1995*, pages 295–303. Elsevier, 1995.
- [93] Jan Peters and Stefan Schaal. Policy gradient methods for robotics. In *2006 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 2219–2225. IEEE, 2006.
- [94] Marc Peter Deisenroth, Peter Englert, Jan Peters, and Dieter Fox. Multi-task policy search for robotics. In *2014 IEEE International Conference on Robotics and Automation (ICRA)*, pages 3876–3881. IEEE, 2014.
- [95] Vishesh Vikas, Piyush Grover, and Barry Trimmer. Model-free control framework for multi-limb soft robots. In *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1111–1116. IEEE, 2015.
- [96] Richard S Sutton, David A McAllester, Satinder P Singh, Yishay Mansour, et al. Policy gradient methods for reinforcement learning with function approximation. In *NIPs*, volume 99, pages 1057–1063. Citeseer, 1999.
- [97] Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3):229–256, 1992.

- [98] Ozgur Simsek, Simon Algorta, and Amit Kothiyal. Why most decisions are easy in tetris—and perhaps in other sequential decision problems, as well. In *International Conference on Machine Learning*, pages 1757–1765. PMLR, 2016.
- [99] Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash Kumar, Henry Zhu, Abhishek Gupta, Pieter Abbeel, and Sergey Levine. Soft actor-critic algorithms and applications, 2019.
- [100] Scott Fujimoto, Herke Hoof, and David Meger. Addressing function approximation error in actor-critic methods. In *International Conference on Machine Learning*, pages 1587–1596. PMLR, 2018.
- [101] Tuomas Haarnoja, Sehoon Ha, Aurick Zhou, Jie Tan, George Tucker, and Sergey Levine. Learning to walk via deep reinforcement learning, 2019.
- [102] Robert M Gray. *Entropy and information theory*. Springer Science & Business Media, 2011.
- [103] Tuomas Haarnoja, Haoran Tang, Pieter Abbeel, and Sergey Levine. Reinforcement learning with deep energy-based policies. In *International Conference on Machine Learning*, pages 1352–1361. PMLR, 2017.
- [104] Piotr Doerffer and Oskar Szulc. Numerical simulation of model helicopter rotor in hover. *Task Quarterly*, 12(3):227–236, 2008.
- [105] Thomas Crittenden, Dmitry Shlyubsky, and Ari Glezer. Combustion-driven jet actuators in reversed flow configurations. In *2nd AIAA Flow Control Conference*, page 2689, 2004.
- [106] MT Patterson and PF Lorber. Computational and experimental studies of compressible dynamic stall. *Journal of fluids and structures*, 4(3):259–285, 1990.
- [107] Gareth D Padfield. *Helicopter flight dynamics: the theory and application of flying qualities and simulation modelling*. John Wiley & Sons, 2008.
- [108] Pascual Marqués and Andrea Da Ronch. *Advanced UAV Aerodynamics, Flight Stability and Control: Novel Concepts, Theory and Applications*. John Wiley & Sons, 2017.
- [109] John M Seddon and Simon Newman. *Basic helicopter aerodynamics*, volume 40. John Wiley & Sons, 2011.
- [110] Sepehr Pourrezaei Khaligh. Control-oriented modeling and system identification for nonlinear trajectory tracking control of a small-scale unmanned helicopter. 2014.
- [111] Bernard Mettler, Mark B Tischler, and Takeo Kanade. System identification modeling of a small-scale unmanned rotorcraft for flight control design. *Journal of the American helicopter society*, 47(1):50–63, 2002.

- [112] Wayne Johnson. *Helicopter theory*. Courier Corporation, 2012.
- [113] Gordon J Leishman. *Principles of helicopter aerodynamics with CD extra*. Cambridge university press, 2006.
- [114] Jean-Jacques E Slotine, Weiping Li, et al. *Applied nonlinear control*, volume 199. Prentice hall Englewood Cliffs, NJ, 1991.
- [115] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- [116] Gabriel Barth-Maron, Matthew W Hoffman, David Budden, Will Dabney, Dan Horgan, Dhruva Tb, Alistair Muldal, Nicolas Heess, and Timothy Lillicrap. Distributed distributional deterministic policy gradients. *arXiv preprint arXiv:1804.08617*, 2018.
- [117] John Schulman, Sergey Levine, Philipp Moritz, Michael I. Jordan, and Pieter Abbeel. Trust region policy optimization. *CoRR*, abs/1502.05477, 2015.
- [118] David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin Riedmiller. Deterministic policy gradient algorithms. In *International conference on machine learning*, pages 387–395. PMLR, 2014.
- [119] The garage contributors. Garage: A toolkit for reproducible reinforcement learning research. <https://github.com/rlworkgroup/garage>, 2019.
- [120] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019.
- [121] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. Software available from tensorflow.org.

- [122] Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. Optuna: A next-generation hyperparameter optimization framework. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2019.

Appendix A

The code of the helicopter environment. The full library is available at [github](https://github.com).

Algorithm 1: Helicopter environment algorithm.

```
1 import sympy as sp
2 import numpy as np
3 from numpy import concatenate as concat
4
5 import gym
6 from gym import spaces
7 from env.Helicopter import Helicopter
8 from utils_main import save_files
9
10
11 class HelicopterEnv(gym.Env):
12     def __init__(self):
13         self.U_input = [U1, U2, U3, U4] = sp.symbols("U1:5", real=True)
14         self.x_state = [
15             u_velocity,
16             v_velocity,
17             w_velocity,
18             p_angle,
19             q_angle,
20             r_angle,
21             fi_angle,
22             theta_angle,
23             si_angle,
24             xI,
25             yI,
26             zI,
27             a_flapping,
28             b_flapping,
29             c_flapping,
30             d_flapping,
31             uwind,
32             vwind,
33             wwind,
34         ] = sp.symbols("x1:20", real=True)
35         self.My_helicopter = Helicopter()
36         self.t = sp.symbols("t")
37         self.symbolic_states_math, jacobian = self.My_helicopter.
38             lambd_eq_maker(self.t, self.x_state, self.U_input)
39         self.default_range = default_range = (-2, 2)
40         self.velocity_range = velocity_range = (-100, 100)
41         self.ang_velocity_range = ang_velocity_range = (-100, 100)
42         self.ang_p_velocity_range = ang_p_velocity_range = (-100, 100)
43         self.Ti, self.Ts, self.Tf = 0, 0.03, 8
44         self.angle_range = angle_range = (-np.pi / 2, np.pi / 2)
45         self.psi_range = psi_range = (-2 * np.pi, 2 * np.pi)
46         self.observation_space_domain = {
47             "u_velocity": velocity_range,
```

```

48         "w_velocity": velocity_range ,
49         "p_angle": ang_p_velocity_range ,
50         "q_angle": ang_velocity_range ,
51         "r_angle": ang_velocity_range ,
52         "fi_angle": angle_range ,
53         "theta_angle": angle_range ,
54         "si_angle": psi_range ,
55         "xI": default_range ,
56         "yI": default_range ,
57         "zI": default_range ,
58         "a_flapping": velocity_range ,
59         "b_flapping": velocity_range ,
60         "c_flapping": velocity_range ,
61         "d_flapping": velocity_range ,
62         "delta_col": (-10, 10) ,
63         "delta_lat": (-10, 10) ,
64         "delta_lon": (-10, 10) ,
65         "delta_ped": (-10, 10) ,
66     }
67     self.states_str = list(self.observation_space_domain.keys())
68     self.low_obs_space = np.array(tuple(zip(*self.
69         observation_space_domain.values()))[0], dtype=np.float32)
70     self.high_obs_space = np.array(tuple(zip(*self.
71         observation_space_domain.values()))[1], dtype=np.float32)
72     self.observation_space = spaces.Box(low=self.low_obs_space, high=
73         =self.high_obs_space, dtype=np.float32)
74     self.default_act_range = (-0.3, 0.3)
75     def_action = (-1, 1)
76     lat_action = (-1, 1)
77     self.action_space_domain = {
78         "col_z": def_action ,
79         "col_w": def_action ,
80         "lon_x": def_action ,
81         "lon_u": def_action ,
82         "lon_q": def_action ,
83         "lon_eul_1": def_action ,
84         "lat_y": lat_action ,
85         "lat_v": lat_action ,
86         "lat_p": lat_action ,
87         "lat_eul_0": lat_action ,
88         "ped_r": def_action ,
89         "ped_eul_3": def_action ,
90     }
91     self.low_action = np.array(tuple(zip(*self.action_space_domain.
92         values()))[0], dtype=np.float32)
93     self.high_action = np.array(tuple(zip(*self.action_space_domain.
94         values()))[1], dtype=np.float32)
95     self.low_action_space = self.low_action
96     self.high_action_space = self.high_action
97     self.action_space = spaces.Box(low=self.low_action_space, high=
98         self.high_action_space, dtype=np.float32)
99     self.min_reward = -13
100
101     self.no_timesteps = int((self.Tf - self.Ti) / self.Ts)

```



```

96     self.all_t = np.linspace(self.Ti, self.Tf, self.no_timesteps)
97     self.counter = 0
98     self.best_reward = float("-inf")
99     self.longest_num_step = 0
100    self.reward_check_time = 0.7 * self.Tf
101    self.high_action_diff = 0.2
102    obs_header = str(list(self.observation_space_domain.keys()))
        [1:-1]
103    act_header = str(list(self.action_space_domain.keys()))[1:-1]
104    self.header = (
105        "time, "
106        + act_header
107        + ", "
108        + obs_header[0:130]
109        + ",a,"
110        + "b,"
111        + "c,"
112        + "d,"
113        + obs_header[189:240]
114        + ",rew,"
115        + "cont_rew,"
116        + "int_rew,"
117        + "si_rew,"
118        + "f_rew,"
119        + "dinput_rew,"
120        + "input_rew,"
121    )
122    self.saver = save_files()
123    self.reward_array = np.array((0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0), dtype=np.float32)
124    self.reward_limit = [
125        1.00e02,
126        3.40e03,
127        1.34e02,
128        1.51e03,
129        3.28e01,
130        7.78e00,
131        3.15e04,
132        3.09e01,
133        3.00e02,
134        8.46e00,
135        1.52e04,
136        9.27e01,
137    ]
138    self.constant_dict = {
139        "u": 0.0,
140        "v": 0.0,
141        "w": 0.0,
142        "p": 1.0,
143        "q": 1.0,
144        "r": 0.0,
145        "fi": 1.0,
146        "theta": 1.0,
147        "si": 0.0,

```

```

148         "x": 0.0,
149         "y": 0.0,
150         "z": 0.0,
151         "a": 0.0,
152         "b": 0.0,
153         "c": 0.0,
154         "d": 0.0,
155     }
156     self.save_counter = 0
157     self.longest_num_step = 0
158     self.best_reward = float("-inf")
159     self.diverge_counter = 0
160     self.numTimeStep = int(self.Tf / self.Ts + 1)
161     self.ifsave = 0
162     self.low_control_input = [0.01, -0.1, -0.1, 0.01]
163     self.high_control_input = [0.5, 0.1, 0.1, 0.5]
164     self.cont_inp_dom = {"col": (-2.1, 2, 1), "lat": (-3.2, 3.2), "
        lon": (-3.5, 3.5), "ped": (-1.1, 1.1)}
165     self.cont_str = list(self.cont_inp_dom.keys())
166     self.initial_states = (
167         np.array(
168             (
169                 3.70e-04, # 0u
170                 1.15e-02, # 1v
171                 4.36e-04, # 2w
172                 -5.08e-03, # 3p
173                 2.04e-04, # 4q
174                 2.66e-05, # 5r
175                 -1.08e-01, # 6fi
176                 1.01e-04, # 7theta
177                 -1.03e-03, # 8si
178                 -4.01e-05, # 9x
179                 -5.26e-02, # 10y
180                 -2.94e-04, # 11z
181                 -4.36e-06, # 12a
182                 -9.77e-07, # 13b
183                 -5.66e-05, # 14c
184                 7.81e-04,
185             ),
186             dtype=np.float32,
187         )
188         + 0.01
189     )
190
191     self.wind1 = np.array((0, 0, 0))
192     self.jk = 1
193
194     def reset(self):
195         # initialization
196         self.t = 0
197         self.all_obs = np.zeros((self.no_timesteps, len(self.
            high_obs_space)))
198         self.all_actions = np.zeros((self.no_timesteps, len(self.
            high_action_space)))

```

```

199     self.all_control = np.zeros((self.no_timesteps, 4))
200     self.all_rewards = np.zeros((self.no_timesteps, 1))
201     self.control_rewards = np.zeros((self.no_timesteps, 1))
202     self.control_rewards1 = np.zeros((self.no_timesteps, 1))
203     self.control_rewards2 = np.zeros((self.no_timesteps, 1))
204     self.control_rewards3 = np.zeros((self.no_timesteps, 1))
205     self.control_rewards4 = np.zeros((self.no_timesteps, 1))
206     self.control_rewards5 = np.zeros((self.no_timesteps, 1))
207     self.control_input = np.array((0, 0, 0, 0), dtype=np.float32)
208     self.jj = 0
209     self.counter = 0
210     self.wind = self.wind1
211     self.jk = self.jk + 0.001
212     self.current_states = concat((self.initial_states, self.wind),
        axis=0)
213     self.current_states[9] = self.initial_states[9]
214     self.current_states[10] = self.initial_states[10]
215     self.current_states[11] = self.initial_states[11]
216     self.observation = self.observation_function()
217     self.done = False
218     self.integral_error = 0
219     return np.clip(self.observation, -0.5, 0.5)
220
221 def action_wrapper(self, current_action, obs) -> np.array:
222     self.normalized_actions = current_action
223     un_act = (current_action + 1) * (self.high_action - self.
        low_action) / 2 + self.low_action
224     self.all_actions[self.counter] = self.normalized_actions #
        unnormalized_action
225     self.control_input[0] = un_act[0] * 5 * obs[11] + un_act[1] * 5
        * obs[2]
226     self.control_input[2] = (
227         un_act[2] * 5 * obs[9] + un_act[3] * 5 * obs[0] + un_act[4]
        * 5 * obs[4] + un_act[5] * obs[7]
228     )
229     self.control_input[1] = (
230         un_act[6] * 5 * obs[10] + un_act[7] * 5 * obs[1] + un_act[8]
        * 5 * obs[3] + un_act[9] * obs[6]
231     )
232     self.control_input[3] = un_act[10] * 5 * obs[5] + un_act[11] * 5
        * obs[8]
233     self.control_input[0] = 2.1167 * np.tanh(self.control_input[0])
        + 0.1
234     self.control_input[1] = 2.03125 * np.tanh(self.control_input[1])
235     self.control_input[2] = 2.02857 * np.tanh(self.control_input[2])
236     self.control_input[3] = 2.2227 * np.tanh(self.control_input[3])
        + 0.18
237
238     self.all_control[self.counter] = self.control_input
239
240 def find_next_state(self) -> list:
241     current_t = self.Ts * self.counter
242     # self.wind = self.wind + 0.4 * (np.random.random() )
243     self.current_states[16:19] = self.wind

```

```

244     self.current_states[0:19] = self.My_helicopter.RK45(
245         current_t, self.current_states[0:19], self.
                symbolic_states_math, self.Ts, self.control_input,
246     )
247
248     def observation_function(self) -> list:
249         self.observation = concat((self.current_states[0:16], self.
                control_input), axis=0)
250         self.all_obs[self.counter] = concat((self.current_states[0:16],
                self.control_input), axis=0)
251         for iii in range(20):
252             current_range = self.observation_space_domain[self.
                states_str[ iii ]]
253             self.observation[ iii ] = (
254                 2 * (self.observation[ iii ] - current_range[0]) / (
                current_range[1] - current_range[0]) - 1
255             )
256         return self.observation
257
258     def reward_function(self, observation, rew_cof=[10, 0.08, 0.015]) ->
        float:
259         error = -rew_cof[0] * (np.linalg.norm(observation[9:12].reshape
                (3), 4))
260         if all(abs(self.current_states[9:12])) < 0.1:
261             error = error + 1 - abs(observation[8])
262         reward = error.copy()
263         self.control_rewards[self.counter] = error
264
265         self.control_rewards1[self.counter] = (
266             0.025 * self.control_rewards[self.counter] + self.
                control_rewards1[self.counter - 1]
267         )
268
269         reward += self.control_rewards1[self.counter]
270         x = self.current_states[9]
271         y = self.current_states[10]
272         si = self.current_states[8]
273         z = self.current_states[11]
274         self.control_rewards2[self.counter] = -0.1 * np.tanh(
275             0.250 / ((1 + 20 * (x ** 2 + y ** 2 + z ** 2)) ** 3 / 2) *
                abs(si)
276         )
277         reward += self.control_rewards2[self.counter]
278
279         self.control_rewards3[self.counter] = 5000 / self.numTimeStep
280         reward += self.control_rewards3[self.counter]
281
282         self.control_rewards4[self.counter] = -rew_cof[1] * sum(
283             abs(self.control_input - self.all_control[self.counter - 1,
                :]))
284         )
285         reward += self.control_rewards4[self.counter]
286
287         self.control_rewards5[self.counter] = -rew_cof[2] * np.linalg.

```

```

288         norm(self.control_input , 2)
289         reward += self.control_rewards5[self.counter]
290         self.all_rewards[self.counter] = reward
291
292     return reward
293
294     def check_diverge(self , reward) -> bool:
295         bool_1 = any(np.isnan(self.current_states))
296         bool_2 = any(np.isinf(self.current_states))
297         if bool_1 or bool_2:
298             self.jj = 1
299             self.observation = self.all_obs[self.counter - 1]
300             reward = self.min_reward - 100
301             return True, reward
302         if np.isnan(reward) or np.isinf(reward):
303             reward = self.min_reward - 100
304             return True, reward
305         for i in range(12):
306             if (abs(self.all_obs[self.counter , i])) > self.
307                 high_obs_space[i]:
308                 self.saver.diverge_save(self.observation_space_domain , i
309                     )
310                 self.jj = 1
311
312         if self.jj == 1:
313             return True, reward
314         if self.counter >= self.no_timesteps - 1: # number of timesteps
315             return True, reward
316         return False , np.clip(reward , -1000, 1000)
317
318     def done_jobs(self) -> None:
319
320         self.best_reward = 0
321         counter = self.counter
322         self.save_counter += 1
323         current_total_reward = sum(self.all_rewards)
324         if self.save_counter >= 1000:
325             self.save_counter = 0
326             self.saver.reward_step_save(self.best_reward , self.
327                 longest_num_step , current_total_reward , counter)
328         if counter >= self.longest_num_step:
329             self.longest_num_step = counter
330         if current_total_reward >= self.best_reward and sum(self.
331             all_rewards) != 0:
332             self.best_reward = current_total_reward
333             ii = self.counter + 1
334             self.saver.best_reward_save(
335                 self.all_t[0:ii] ,
336                 self.all_actions[0:ii] ,
337                 self.all_obs[0:ii] ,
338                 self.all_rewards[0:ii] ,
339                 np.concatenate(

```

```

337         self.control_rewards[0:ii],
338         self.control_rewards1[0:ii],
339         self.control_rewards2[0:ii],
340         self.control_rewards3[0:ii],
341         self.control_rewards4[0:ii],
342         self.control_rewards5[0:ii],
343     ),
344     axis=1,
345 ),
346     self.header,
347 )
348
349 def step(self, current_action):
350     self.control_input = current_action
351     try:
352         self.find_next_state()
353     except OverflowError or ValueError or IndexError:
354         self.jj = 1
355     self.observation = self.observation_function()
356     reward = self.reward_function(self.observation)
357     self.done, reward = self.check_diverge(reward)
358     if self.jj == 1:
359         reward -= self.min_reward
360     if self.done:
361         self.done_jobs()
362     self.counter += 1
363     if np.isnan(reward) or any((np.isnan(self.observation))):
364         reward = -100
365         self.current_states = self.initial_states * 0 - 10
366         self.observation = self.observation_function()
367     return np.clip(self.observation, -100, 100), np.clip(reward,
368         -1000, 1000), self.done, {}
369
370 def make_constant(self, true_list):
371     for i in range(len(true_list)):
372         if i == 1:
373             self.current_states[i] = self.initial_states[i]
374
375 def close(self):
376     return None

```
