

Experiment zur Verbesserung der Robustheit von Reinforcement Learning Modellen in Drohnensimulationen anhand trainiertem Gegenspieler

Bachelorarbeit

vorgelegt am 1. Mai 2023

Fakultät Wirtschaft

Studiengang Wirtschaftsinformatik

Kurs WWI2020F

von

LEON HENNE

Betreuerin in der Ausbildungsstätte: DHBW Stuttgart:

IBM Deutschland GmbH
Sophie Lang
Senior Data Scientist

Prof. Dr. Kai Holzweißig
Studiendekan Wirtschaftsinformatik

Unterschrift der Betreuerin

Vertraulichkeitsvermerk: Der Inhalt dieser Arbeit darf weder als Ganzes noch in Auszügen Personen außerhalb des Prüfungs- und Evaluationsverfahrens zugänglich gemacht werden, sofern keine anders lautende Genehmigung des Dualen Partners vorliegt.

Inhaltsverzeichnis

Abkürzungsverzeichnis	IV
Abbildungsverzeichnis	V
Tabellenverzeichnis	VI
1 Einleitung	1
1.1 Problemstellung	1
1.2 Zielsetzung	2
1.3 Forschungsfrage	3
1.4 Forschungsmethodik	3
1.5 Aufbau der Arbeit	4
2 Diskussion des aktuellen Stands der Forschung und Praxis	5
2.1 Aufbau der Literaturrecherche	5
2.2 Verstärkendes Lernen	6
2.2.1 Methoden des verstärkenden Lernens	8
2.2.2 Algorithmen des verstärkenden Lernens	10
2.2.3 Abgrenzung zu Multi-Agent Reinforcement Learning (MARL) Algorithmen	12
2.2.4 Limitierungen und Herausforderungen von RL	13
2.3 Imitierendes Lernen	14
2.4 Simulationsumgebungen für RL	16
2.4.1 Definitionen von Simulationsumgebungen	16
2.4.2 Entwicklung von Simulationsumgebungen für RL Anwendungen	17
2.4.3 Aktuelle Physik-Engines und Simulationsanwendungen	18
2.5 Simulation der Steuerungsaufgabe von Quadroptern	19
2.5.1 Flugdynamiken eines Quadropters	20
2.5.2 Existierende Simulationen von Quadroptern	21
2.6 Robustheit und Stabilität von Strategien des verstärkenden Lernens	23
2.6.1 Definitionen von Robustheit und Stabilität	23
2.6.2 Metriken der Robustheit	24
2.6.3 Experimenteller Rahmen zur Messung der Robustheit	24
2.7 Gegnerisches verstärkendes Lernen	25
2.8 Domain Randomization	27
2.8.1 Das Prinzip hinter der Randomisierung von Simulationsumgebungen	27
2.8.2 Anwendung von Randomisierung der Simulationsumgebung	28
2.8.3 Errungenschaften unter Einsatz von Domain Randomization	29
3 Durchführung des Laborexperiments	30
3.1 Erläuterung der Forschungsmethodik	30
3.1.1 Beschreibung der Simulationsumgebung	30
3.1.2 Erläuterung der Trainingsszenarien	31
3.1.3 Erläuterung des Testszenarios	31
3.1.4 Messung der Robustheit von RL Policies	32
3.1.5 Auswertung mittels statistischer Tests	32
3.2 Programmanforderungen	34

3.2.1	Anforderungen der Simulationsumgebung	34
3.2.2	Anforderungen des Optimierungsverfahrens	36
3.2.3	Anforderungen des Laborexperiments	37
3.3	Implementierung der Simulation und des Experiments	38
3.3.1	Programmumsetzung der Simulationsumgebung	38
3.3.2	Programmumsetzung des Optimierungsverfahrens	43
3.3.3	Programmumsetzung des Laborexperiments	45
4	Ergebnisse des Laborexperiments	48
4.1	Vergleich der Robustheit der Strategien aus dem Training mit regelbasierten und RL Gegenspieler	48
4.1.1	Betrachtung der kumulierten Belohnung	49
4.1.2	Betrachtung der Belohnungsstabilität	49
4.1.3	Betrachtung der Anzahl der Misserfolge	50
4.2	Vergleich der Robustheit der Strategien aus dem Training mit regelbasierten Gegenspieler und Domain Randomization	51
4.2.1	Betrachtung der kumulierten Belohnung	51
4.2.2	Betrachtung der Belohnungsstabilität	52
4.2.3	Betrachtung der Anzahl der Misserfolge	53
4.3	Vergleich der Robustheit der Strategien aus dem Training mit RL basiertem Gegenspieler und Domain Randomization	53
4.3.1	Betrachtung der kumulierten Belohnung	54
4.3.2	Betrachtung der Belohnungsstabilität	54
4.3.3	Betrachtung der Anzahl der Misserfolge	55
4.4	Beantwortung der Forschungsfrage	55
5	Reflexion und Forschungsausblick	56
5.1	Reflexion der Forschungsmethodik	56
5.2	Ausblick für weitere Forschung	56
	Anhang	57
	Literaturverzeichnis	66

Abkürzungsverzeichnis

A2C	Advantage Actor-Critic
A3C	Asynchronous Advantage Actor-Critic
CSV	kommaseparierte Wertedatei
DART	Dynamic Animation and Robotics Toolkit
DHBW	Duale Hochschule Baden-Württemberg
DQL	Deep Q-Learning
DR	Domain Randomization
IL	Imitation Learning
KPI	Key Performance Indicator
MARL	Multi-Agent Reinforcement Learning
ODE	Open Dynamics Engine
PPO	Proximal Policy Optimization
RAEARL	Robust Adaptive Ensemble Adversarial Reinforcement Learning Framework
RL	Reinforcement Learning
ROS	Robot Operating System
SITL	Software-in-the-Loop
TRPO	Trust Region Policy Optimization

Abbildungsverzeichnis

1	vereinfachte Darstellung der Interaktion zwischen dem Agenten und seiner Umgebung	7
2	Klassifizierung von Algorithmen im Bereich des RL	8
3	Ablaufdiagramm des imitierenden Lernens	15
4	Rotationsbewegungen eines Quadropters	20
5	Aufbau des RAEARL Frameworks	27
6	Intuition hinter dem Paradigma von DR	28
7	Einfluss der Randomisierung verschiedener Effekte auf die Fehlerrate der Objekterkennung	29
8	grafisches Modell der Simulation	40
9	Episodenlänge, Belohnung und Erfolgsrate unter verschiedenen Gewichtungen	42
10	Klassendiagramm mit Funktionsabhängigkeiten	44
11	BPMN Prozess des Laborexperiments	47

Tabellenverzeichnis

1	Konzeptmatrix für Artikel zu Simulationsumgebungen und zur Robustheit RL Algorithmen nach Webster/Watson 2002. Legende: RL (Reinforcement Learning), MARL (Multi-Agent Reinforcement Learning), IL (Imitation Learning), ES (Entwicklung von Simulationsumgebungen), DS (Drohnen-simulation), KS (kompetitive Simulationsumgebungen), DR (Domain Randomization), RRLP (Robustheit von RL Policies) . . .	6
2	wichtigsten Kriterien zur Auswahl von Simulatoren ¹	18
3	Interne quantitative Metriken, dessen gemessenes Verhalten und ihre Häufigkeit ² . .	24
4	Auszug der externen quantitativen Metriken, dessen gemessenes Verhalten und ihre Häufigkeit ³	25
5	Gegenüberstellung von Auswahlkriterien und bekannten Drohnensimulationen	39
6	Beispieltabelle der erhobenen aggregierten Messdaten	46
7	Ergebnisse der statistischen Tests aus dem Leistungsvergleich der Modelle aus den Trainingsszenarien 1 und 3	48
8	Ergebnisse der statistischen Tests aus dem Leistungsvergleich der Modelle aus den Trainingsszenarien 1 und 4	51
9	Ergebnisse der statistischen Tests aus dem Leistungsvergleich der Modelle aus den Trainingsszenarien 3 und 4	54

¹Ivaldi/Padois/Nori 2014, S. 4

²Pullum 2022, S.17

³Pullum 2022, S.19

1 Einleitung

1.1 Problemstellung

Reinforcement Learning (RL) findet heutzutage bereits Anwendung in vielerlei Forschungsprojekten wie Deepmind AlphaStar oder OpenAI Five, aber auch in Produkten und Dienstleistungen wie AWSDeepRacer oder Meta's Horizon open-source RL-Plattform.⁴ RL ist im Bereich des maschinellen Lernens eine Herangehensweise zur Lösung von Entscheidungsproblemen.⁵ Ein Software-Agent leitet dabei durchzuführende Aktionen aus seiner Umgebung ab, mit dem Ziel die kumulierte erhaltene Belohnung zu maximieren, währenddessen sich seine Umgebung durch alle Aktionen verändert.⁶ Die Umgebungen beinhalten in ihrer einfachsten Form eine simulierte Welt, welche zu jedem Zeitschritt eine Aktion entgegennimmt, und den eigenen nächsten Zustand sowie einen Belohnungswert zurückgibt.⁷ Da das Sammeln von Daten in der echten Welt für den Einsatz von RL Algorithmen eine Limitierung darstellen kann, werden häufig fürs Training Simulationsumgebungen eingesetzt.⁸ Eine Limitierung können bspw. Sicherheitsaspekte sein, welche beim Training von Roboterarmen, oder sich autonom bewegenden Systemen auftreten, da die einzelnen physischen Bewegungen nicht vorhersehbar abschätzbar sind.⁹ Simulationen nehmen damit zum einen als Testumgebung eine wichtige Rolle in der Entwicklung von Kontrollalgorithmen ein.¹⁰ Zum anderen bedarf die erfolgreiche Anwendung von RL neben effizienten Algorithmen auch geeignete Simulationsumgebungen.¹¹ Besonders schwierig, und daher sehr wichtig zu erforschen, ist die Anpassung der Trainingsumgebung an die reale Welt, sodass bspw. Agenten für Roboter und autonome Fahrzeuge nach dem Training mit robusten Strategien in der Realität eingesetzt werden können.¹² In der Forschungsliteratur wird diese beschriebene Problematik als „Sim to real“-Transfer beschrieben.¹³

Ein Forschungsgebiet, bei dem die Lösung des Sim to Real Transfers betrachtet wird, ist die autonome Steuerung von unbemannten Luftfahrzeugen bzw. Drohnen.¹⁴ Das Transferproblem entsteht bspw. dabei, dass zur automatisierten Kollisionsvermeidung die Kollisionsbeispiele einer Simulationsumgebung entnommen werden, um physischen Schaden oder Drohnen Verlust zu vermeiden.¹⁵ Drohnen tragen dabei bereits in der heutigen Zeit zur Lösung vieler komplexer Aufgaben bei, wie der Katastrophenüberwachung oder der Waldbrandbekämpfung.¹⁶ Zusätzlich

⁴Vgl. Li 2019, S. 4

⁵Vgl. Schuderer/Bromuri/van Eekelen 2021, S. 3

⁶Vgl. Schuderer/Bromuri/van Eekelen 2021, S. 3

⁷Vgl. Reda/Tao, T./van de Panne 2020, S. 1

⁸Vgl. Zhao/Queralta/Westerlund 2020, S. 737

⁹Vgl. Zhao/Queralta/Westerlund 2020, S. 738

¹⁰Vgl. Cutler/Walsh/How 2014, S. 2

¹¹Vgl. Reda/Tao, T./van de Panne 2020, S. 8

¹²Vgl. Slaoui u. a. 2019, S. 1

¹³Vgl. Zhao/Queralta/Westerlund 2020, S. 738

¹⁴Vgl. Deshpande/Minai/Kumar, M. 2021, S. 1

¹⁵Vgl. Sadeghi/Levine 2016, S. 4

¹⁶Vgl. Hentati u. a. 2018, S. 1495

steigen immer weiter die komplexen Einsatzanforderungen unter dem Anstieg an Anwendungsgebieten für unbemannte Luftfahrzeuge.¹⁷

Die Simulation einer möglichst realistischen Umgebung in diesem Kontext wird in der Forschung häufig mit dem Ansatz von Domain Randomization (DR) begleitet.¹⁸ Unter dem Themenfeld der DR wird die Idee erforscht, anstelle der akkuraten Modellierung realistischer Dynamiken, diese so stark zu randomisieren, dass reale Dynamikeffekte abgedeckt sind.¹⁹ Neben den dynamischen Bedingungen unterliegt die Realität jedoch häufig auch dem Einfluss mehrerer Parteien. Diese tragen teilweise kooperierend aber auch teilweise konkurrierend zum eigenen Erfolg bei, wie z.B. im Rahmen eines dem Wettbewerb unterliegenden Marktes.²⁰ Stellt man sich ein Szenario im Kontext kooperativer oder konkurrierender Drohnen vor, ist es naheliegend, dass auch jene Einflüsse möglichst präzise in die Simulationsumgebung integriert sein müssen, um ein robustes Modell erlernen zu können. Während bereits in Produkten wie PowerTAC von Collins/Ketter 2022 die Simulation von Märkten entwickelt wurde, scheint der Einfluss des Gegenspielers in kompetitiven Drohnensimulationen auf die Robustheit von RL Algorithmen, und demnach auf die Lösung des „Sim to real“-Transfers, unerforscht.

1.2 Zielsetzung

Im Rahmen dieser Arbeit soll eine Simulationsumgebung entwickelt werden, in welcher sich RL basierter Gegenspieler und deterministische Gegenspieler integrieren lassen. Anschließend ist zu untersuchen, wie sich die Wahl des Trainingsgegenspielers auf das Leistungsverhalten der RL Modelle Verhaltensmodelle, im Kontext von RL oftmals als Policies referenziert, unter verändertem Testszenario auswirkt.

Dazu soll eine kompetitive Simulationsumgebung entwickelt werden, in welcher sich zwei konkurrierender Spieler in Form von Flugobjekten spielerisch gegenseitig bekämpfen. In der Simulation werden folgende Policies in drei verschiedenen Szenarien optimiert.

- Training mit regelbasiertem Gegenspieler unter gleichbleibenden Dynamikparametern
- Training mit RL basiertem Gegenspieler unter gleichbleibenden Dynamikparametern
- Training mit regelbasiertem Gegenspieler unter sich verändernden Dynamikparametern

Anschließend werden alle trainierten Policies in einem Testszenario untersucht. Das Testszenario verfügt dabei über festgelegte, sich vom Training unterscheidende Eigenschaften, wie Flugbahnen, oder Start- und Zielpositionen. Außerdem wird ein deterministischer Gegenspieler unter im Gegensatz zum Training, veränderten Handlungspräferenzen eingesetzt. Bei der Untersuchung werden jeweils die folgenden Variablen als Key Performance Indicator (KPI) betrachtet.

¹⁷Vgl. Deshpande/Kumar, R. u. a. 2020, S. 1

¹⁸Vgl. Sadeghi/Levine 2016, S. 1

¹⁹Vgl. Zhao/Queralt/Westerlund 2020, S. 4f.

²⁰Vgl. Collins/Ketter 2022, S. 2

- kumulierte erzielte Belohnung
- Varianz der Belohnungen
- Anzahl an Misserfolgen

Durch die Auswertung des Testszenarios kann der Effekt des RL basierten Gegenspielers auf die Robustheit evaluiert werden, indem die Leistungsdiskrepanz im Testszenario zwischen den Strategien im Training mit RL- und regelbasierten Gegenspielers- und der Domain Randomization verglichen wird.

1.3 Forschungsfrage

Aus der beschriebenen Problemstellung und der für den Rahmen dieser Arbeit festgelegten Zielsetzung ergibt sich folgende Forschungsfrage:

Inwiefern kann durch den Einsatz eines mittels RL trainierten Gegenspielers die Robustheit der gelernten Policy verbessert werden?

Zur Beantwortung der Forschungsfrage werden folgende Hypothesen aufgestellt und im Rahmen der Arbeit untersucht:

Hypothese 1: *Die im Testszenario erzielte kumulierte Belohnung ist unter Verwendung der Policy aus dem Training mit RL basiertem Gegenspieler signifikant und zuverlässig höher als die Policy aus dem Training mit regelbasiertem Gegenspieler.*

Hypothese 2: *Die Varianz der im Testszenario erzielten kumulierten Belohnung ist unter Verwendung der Policy aus dem Training mit RL basiertem Gegenspieler signifikant und zuverlässig geringer als die Policy aus dem Training mit regelbasiertem Gegenspieler.*

Hypothese 3: *Die im Testszenario erreichte Anzahl von Misserfolgen ist unter Verwendung der Policy aus dem Training mit RL basiertem Gegenspieler signifikant und zuverlässig geringer als die Policy aus dem Training mit regelbasiertem Gegenspieler.*

1.4 Forschungsmethodik

Als Forschungsmethodik soll im Rahmen dieser Arbeit ein quantitatives Laborexperiment nach Recker 2021 durchgeführt werden. Hierbei wird häufig nach dem hypothetisch-deduktiven Modell vorgegangen, in welchem Hypothesen formuliert, empirische Studien entwickelt, Daten gesammelt, Hypothesen anhand dieser evaluiert und gewonnene Erkenntnisse berichtet werden.²¹ Damit stellt das Laborexperiment eine Möglichkeit der Untersuchung der Ursache- und Wirkungsbeziehung dar.²² Dabei wird die kontrollierte Umgebung der Simulation erschaffen, deren

²¹Vgl. Recker 2021, S. S.89f.

²²Vgl. Recker 2021, S. 106

Aufbau die unabhängige Variable des Laborexperiments darstellt. Die Metriken, anhand der die Leistung und die Robustheit der trainierten Policies gemessen werden, bilden im Experiment die abhängigen Variablen, zu dessen Grundlage die Forschungsfrage beantwortet wird.

1.5 Aufbau der Arbeit

Insgesamt gliedert sich die Arbeit nach einem Schema von Holzweißig 2022. Die Arbeit beginnt mit einem einleitenden Kapitel, in welchem Motivation, Problemstellung, Zielsetzung und Forschungsmethodik erläutert sind. Anschließend wird im zweiten Kapitel der aktuelle Stand der Forschung zu den relevanten Konzepten beschrieben. Im dritten Kapitel wird die Forschungsmethodik dargestellt, indem die Simulationsumgebung als Messinstrument entwickelt wird, verschiedene Messszenarien erläutert und entsprechende Daten gesammelt werden. Daraufhin erfolgt im vierten Kapitel die Auswertung der Messdaten sowie die Beantwortung der Forschungsfrage durch Annahme oder Ablehnung der eigens aufgestellten Hypothesen. Im Zuge dessen kann ebenso die Forschungsfrage anhand der Annahme oder Ablehnung der Hypothesen beantwortet werden. Abschließend wird im letzten Kapitel ein Fazit zu den erzielten Forschungsergebnissen dargelegt und ein Ausblick auf nachfolgende Forschungsansätze gegeben.

2 Diskussion des aktuellen Stands der Forschung und Praxis

2.1 Aufbau der Literaturrecherche

In Anlehnung an die Literaturrecherche nach Webster/Watson 2002 wurden alle voraussichtlich benötigten Konzepte für die Durchführung der beschriebenen Forschungsmethodik in Tabelle 1 festgehalten. Alle angeführten Konzepte wurden mittels verschiedener Suchbegriffe in Suchmaschinen, Datenbanken und Bibliotheken wie *Google Scholar* 2/28/2023, *IEEE Xplore* 2/28/2023 oder die digitale Bibliothek der Association for Computing Machinery (ACM) ACM Digital Library 2/28/2023 recherchiert. In der daraus gefunden Literatur wurden zitierte Werke ebenfalls nach den beschriebenen Konzepten durchsucht und insgesamt jede Literaturquelle in Tabelle 1 den in ihnen enthaltenen Konzepten zugeordnet.

Artikel	Konzepte							
	RL	MARL	IL	ES	DS	KS	DR	RRLP
Sutton/Barto 2018	X							
Li 2019	X							
Zhao/Queralta/Westerlund 2020	X			X			X	
Wang/Hong 2020	X							
Zhang/Wu/Pineau 2018	X			X				X
Cutler/Walsh/How 2014	X			X				
Canese u. a. 2021	X	X						
Reda/Tao, T./van de Panne 2020	X			X			X	
Ningombam 2022	X							
Arulkumaran u. a. 2017	X							
Huang u. a. 2017	X							
Mnih/Kavukcuoglu u. a. 2013	X							
Wong u. a. 2022	X	X						
Hussein u. a. 2017	X		X					
Attia/Dayan 2018			X					
Gao u. a. 2014			X					
Fang u. a. 2019			X					
Balakrishna u. a. 2020			X					
Schuderer/Bromuri/van Eekelen 2021	X	X		X				
Körber u. a. 2021				X				
Bharadhwaj u. a. 2019				X			X	
Foronda 2021				X				
Maria 1997				X				

Artikel	Konzepte							
Brockman u. a. 2016	X			X				
Yan Duan u. a. 2016	X			X				X
Ivaldi/Padois/Nori 2014				X				
Ayala u. a. 2020				X				
Todorov/Erez/Tassa 2012				X				
Koch u. a. 2018	X				X			
Deshpande/Kumar, R. u. a. 2020	X				X			
Deshpande/Minai/Kumar, M. 2021					X		X	X
Hentati u. a. 2018					X			
Molchanov u. a. 2019					X		X	X
Furrer u. a. 2016					X			
Silano/Iannelli 2019					X			
Shah u. a. 2017					X			
Panerati u. a. 2021				X	X			
Moos u. a. 2022								X
Pullum 2022								X
Liu u. a. 2023								X
Yan Duan u. a. 2016								X
Schott/Hajri/Lamprier 2022	X					X		
Pinto u. a. 2017	X					X		X
Pan u. a. 2021						X		
Zhai u. a. 2022						X		X
Tobin u. a. 2017							X	
Chen u. a. 2021							X	
Hsu u. a. 2023							X	
Alghonaim/Johns 2021							X	
Sadeghi/Levine 2016							X	

Tab. 1: Konzeptmatrix für Artikel zu Simulationsumgebungen und zur Robustheit RL Algorithmen nach Webster/Watson 2002. Legende: RL (Reinforcement Learning), MARL (Multi-Agent Reinforcement Learning), IL (Imitation Learning), ES (Entwicklung von Simulationsumgebungen), DS (Drohensimulation), KS (kompetitive Simulationsumgebungen), DR (Domain Randomization), RRLP (Robustheit von RL Policies)

2.2 Verstärkendes Lernen

Verstärkendes Lernen, oder auch als Reinforcement Learning (RL) in der Fachsprache bezeichnet, definiert einen konzeptionellen Ansatz, zielorientiertes Lernen von Entscheidungen zu verstehen

und zu automatisieren.²³ Dabei besteht der Fokus darauf, dass ein Agent aus der direkten Interaktion mit seiner Umgebung lernt, ohne dass explizite Überwachung notwendig ist.²⁴ Der Agent lernt über die Zeit eine optimale Strategie zur Lösung des Entscheidungsproblems durch Ausprobieren und Scheitern die gewünschte Veränderung in seiner Umwelt herzustellen.²⁵ Notwendig dabei ist es, dass der Agent den Zustand seiner Umgebung wahrnehmen und durch entsprechende Aktionen beeinflussen kann, um das Erreichen des Zielzustandes zu ermöglichen.²⁶ Dazu muss der Agent diejenigen Aktionen finden, die ihm die größtmögliche kumulierte Belohnung liefern, wobei Aktionen nicht nur die unmittelbaren, sondern auch zukünftige Belohnungen beeinflussen.²⁷ Zusammengefasst lässt sich die beschriebene Interaktion des Agenten mit seiner Umgebung wie in Abbildung 1 darstellen.

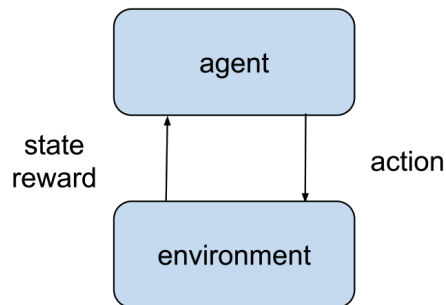


Abb. 1: vereinfachte Darstellung der Interaktion zwischen dem Agenten und seiner Umgebung²⁸

Ein Standardaufbau einer Aufgabe für verstärkendes Lernen kann demnach als sequentielles Entscheidungsproblem verstanden werden, zu dessen Lösung ein Agent zu jedem diskreten Zeitschritt eine Aktion ausführt, welche den Zustand der Umgebung verändert.²⁹ Betrachtet man die technische Umsetzung einer solchen Interaktion zwischen dem Agenten und dessen Umgebung, wird häufig zur Modellierung ein Markov Entscheidungsprozess verwendet. Im Kontext von RL ist der Entscheidungsprozess definiert als ein Tupel aus folgenden Elementen:³⁰

- Zustandsraum S
- Aktionsraum A
- initiale Zustandsverteilung $p_0(S)$
- Übergangswahrscheinlichkeit $T(S_{t+1}|S_t, A_t)$
- Belohnungswahrscheinlichkeit $R(r_{t+1}|S_t, A_t)$

²³Vgl. Sutton/Barto 2018, S. 13

²⁴Vgl. Sutton/Barto 2018, S. 13

²⁵Vgl. Li 2019, S. 4

²⁶Vgl. Sutton/Barto 2018, S. 2

²⁷Vgl. Sutton/Barto 2018, S. 1

²⁸Enthalten in: Li 2019, S. 5

²⁹Vgl. Zhao/Queralt/Westerlund 2020, S. 2

³⁰Vgl. Zhang/Wu/Pineau 2018, S. 2

Zum Finden der optimalen Strategie existieren modellbasierende und modellfreie Algorithmen des verstärkenden Lernens.³¹ Bei modellbasierenden Algorithmen wird das Umgebungsverhalten, also die Übergangs- und Belohnungswahrscheinlichkeiten, als bekannt vorausgesetzt.³² Als modellbasierender Algorithmus wird z. B. dynamische Programmierung eingesetzt, um mittels Strategieevaluation und Strategieiteration die optimale Strategie zu finden.³³ Unter modellfreien Algorithmen wird zwischen einem wertbasierten-, einen strategiebasierten und einem Akteur-Kritik-Ansatz unterschieden.³⁴ Der Agent im Kontext von modellfreien RL Methoden kennt dabei nur die Zustände S und die Aktionen A , jedoch nicht das Umgebungsverhalten T oder die Belohnungswahrscheinlichkeit R .³⁵ Abbildung 2 stellt eine Klassifizierung der Algorithmen und Methoden von RL dar.

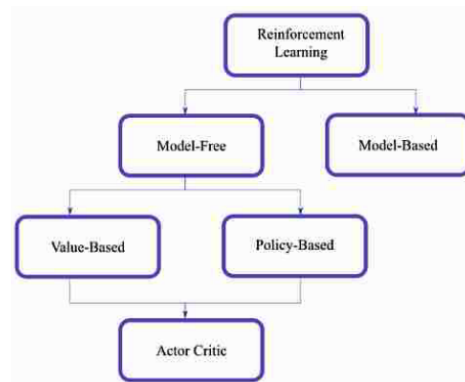


Abb. 2: Klassifizierung von Algorithmen im Bereich des RL³⁶

2.2.1 Methoden des verstärkenden Lernens

Der Agent sucht im Kontext von **wertbasierenden Methoden** die optimale Strategie π^* , welche allen Zuständen S die jeweilige Aktion $A(S)$ zuordnet, sodass die kumulierte Belohnungswahrscheinlichkeit $R(r_{t+1}|S_t, A_t)$ über alle Zeitschritte t maximal ist.³⁷ Neben der unmittelbaren Belohnung aus einem Zustands-Aktions Paar müssen auch die zukünftigen Belohnungen der möglichen Folgezustände betrachtet werden, wofür das Konzept der Wertigkeit eingeführt wird.³⁸ Über eine Zustands- oder Aktionswertigkeitsfunktion, oftmals als Q-Funktion referenziert, wird eine Vorhersage über die zu erwartende abgezinste Belohnung berechnet.³⁹ Durch den Abzinsungsfaktor $\gamma \in [0, 1)$ wird der Einfluss zukünftiger Belohnungen nach ihrer zeitlichen Reihenfolge priorisiert.⁴⁰ Mit der Wertigkeitsfunktion kann evaluiert werden, welche Strategie langfristig am erfolgreichsten ist, da bspw. manche Aktionen trotz geringer sofortiger Belohnung

³¹Vgl. Wang/Hong 2020, S. 3

³²Vgl. Wang/Hong 2020, S. 3

³³Vgl. Li 2019, S. 5

³⁴Vgl. Li 2019, S. 5

³⁵Vgl. Cutler/Walsh/How 2014, S. 2

³⁶Enthalten in: Canese u. a. 2021, S. 6

³⁷Vgl. Reda/Tao, T./van de Panne 2020, S. 2

³⁸Vgl. Wang/Hong 2020, S. 3

³⁹Vgl. Li 2019, S. 5

⁴⁰Vgl. Li 2019, S. 5

einen hohen Wert aufweisen können, wenn aus dem zukünftigen Zustand eine hohe Belohnung zu erwarten ist.⁴¹ Die Wertigkeitsfunktion und die daraus berechneten Wertigkeiten von Aktionen oder Zuständen werden über alle Zeitschritte neu geschätzt und stellen mit die wichtigsten Komponenten in Algorithmen des verstärkenden Lernens dar.⁴² Methoden basierend auf diesem Wertigkeitswert lernen eine Schätzfunktion der Wertigkeit für alle Zustände ($V_\pi(s) \forall s$) und alle Zustandsaktions-Paare ($Q_\pi(s_t, a_t) \forall s, a \in (S, A)$) durch das Aktualisieren der folgenden Funktionen:⁴³

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left[r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t) \right] \quad (1)$$

$$V(s_t) = \max_a Q(s_t, a | \omega) \quad (2)$$

Aus der geschätzten Wertigkeit jedes Zustandsaktions-Paares kann die optimale Strategie $\pi^*(s)$ durch $\arg \max_a Q(s, a)$ bestimmt werden.⁴⁴

Methoden, welche die Strategie durch direkte Parametrisierung anstelle einer Bewertung aller Handlungsalternativen mittels Wertigkeitsfunktion optimieren, werden als **strategiebasierend** bezeichnet.⁴⁵ Diese Methodik kann beim Trainieren deterministischer Strategien zu unerwarteten Aktionen führen, weshalb häufig das Optimieren einer Wahrscheinlichkeitsverteilung für alle Aktionen bevorzugt wird.⁴⁶ Als Subklasse der RL Methoden wird der statistische Gradientenabstieg verwendet, um die parametrisierte Strategie π_θ hinsichtlich der maximalen langfristigen kumulierten Belohnung zu optimieren.⁴⁷ Unter dem Parametervektor θ beschreibt die Strategie π_θ oder auch $\pi(a|s, \theta)$ die Wahrscheinlichkeit Aktion a im Zustand s auszuwählen.⁴⁸ Zur Optimierung der Strategie wird die Funktion der kumulierten Belohnungen J nach dem Parameter der Gewichte wie folgt in Formel drei abgeleitet und der optimierte Parametervektor anhand Formel vier aktualisiert.⁴⁹

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left[\left(\sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_t | s_t) \right) \left(\sum_{t=1}^T r(s_t, a_t) \right) \right] \quad (3)$$

$$\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta) \quad (4)$$

Zusammengefasst können die Formeln vier und fünf dabei so interpretiert werden, dass die logarithmierte Wahrscheinlichkeit Aktion a_t im Zustand s_t auszuwählen erhöht wird, wenn a_t in einer höheren kumulierten Belohnung resultiert.⁵⁰

⁴¹Vgl. Sutton/Barto 2018, S. 6

⁴²Vgl. Sutton/Barto 2018, S. 6f.

⁴³Vgl. Zhang/Wu/Pineau 2018, S. 2

⁴⁴Vgl. Zhang/Wu/Pineau 2018, S. 2

⁴⁵Vgl. Zhang/Wu/Pineau 2018, S. 2

⁴⁶Vgl. Ningombam 2022, S. 3

⁴⁷Vgl. Ningombam 2022, S. 3

⁴⁸Vgl. Sutton/Barto 2018, S. 321

⁴⁹Vgl. Wang/Hong 2020, S. 6

⁵⁰Vgl. Wang/Hong 2020, S. 6

Unter **Akteur-Kritiker Methoden** werden hybride wertebasierende und strategiebasierende Methoden verstanden, welche zugleich die Strategie optimieren und eine Wertefunktion approximieren.⁵¹ Die strategiebasierende Methodik mit der lernenden Strategie agiert dabei als Akteur, wohingegen die Wertefunktion, welche jeder Aktion und jedem Zustand einen Belohnungswert zuweist, als Kritiker handelt.⁵² Der Akteur wählt somit aus seiner Wahrscheinlichkeitsverteilung die auszuführende Aktionen aus, während der Kritiker diese anhand seiner Wertigkeit bewertet.⁵³ Betrachtet man den Trainingsprozess von Akteur-Kritiker basierten Methoden, ist dieser wie folgt aufgebaut:⁵⁴

1. Aktueller Zustand der Umgebung als Eingabe dem Akteur und Kritiker übergeben
2. Akteur liefert eine auszuführende Aktion basierend auf dem Umgebungszustand
3. Der Kritiker bekommt die Aktion als Eingabe und berechnet deren Wertigkeit mittels Q-Funktion
4. Durch die Wertigkeit seiner Aktion kann der Akteur seine Strategie anpassen
5. Mit der neuen Strategie führt der Akteur die nächste Aktion im folgenden Zustand aus
6. Die Q-Funktion des Kritikers wird mit den neuen Informationen aus der erhaltenen Belohnung angepasst

2.2.2 Algorithmen des verstärkenden Lernens

Im vorherigen Kapitel sind die Methoden des verstärkenden Lernens klassifiziert und deren Unterschiede beschrieben worden. Anschließend daran wird in diesem Abschnitt der Arbeit auf eine Auswahl konkreter Algorithmen und Implementierungen, die verschiedenen Methoden sowie deren Funktionsweise eingegangen. Die Auswahl an Algorithmen beinhaltet die fundamentalen Algorithmen des verstärkenden Lernens Deep Q-Learning (DQL), Trust Region Policy Optimization (TRPO) und Asynchronous Advantage Actor-Critic (A3C).⁵⁵

Mit der Entwicklung von DQL konnte erstmals ein Algorithmus entwickelt werden, welcher eine Reihe von Atari 2600 Spielen auf der Fähigkeitsebene eines professionellen Videospieltesters spielen konnte.⁵⁶ Als wertbasierte Methode wird für jeden Zustand die Wertigkeit berechnet, was unter DQL mittels eines neuronalen Netzes erreicht wird, welches die Q-Funktion approximiert.⁵⁷ Durch Auswählen der Aktion mit der höchsten Wertigkeit in jedem Zustand, kann die deterministische Strategie abgeleitet werden.⁵⁸ DQL adressiert das Instabilitätsproblem von

⁵¹Vgl. Zhang/Wu/Pineau 2018, S. 2f.

⁵²Vgl. Sutton/Barto 2018, S. 321

⁵³Vgl. Ningombam 2022, S. 3

⁵⁴Vgl. Ningombam 2022, S. 4

⁵⁵Vgl. Arulkumaran u. a. 2017, S. 1

⁵⁶Vgl. Arulkumaran u. a. 2017, S. 6

⁵⁷Vgl. Huang u. a. 2017, S. 4

⁵⁸Vgl. Huang u. a. 2017, S. 4

Funktionsapproximation unter dem Einsatz von Erfahrungswiederholung und Zielnetzwerken.⁵⁹ Erfahrungswiederholung beschreibt die Technik, die Beobachtungen des Agenten, also in welchem Zustand welche Aktion zu welcher Belohnung und welchem Folgezustand führt, in jedem Zeitschritt zu speichern und zufällig anhand Übergangserfahrungen aus anderen Trainingsläufen die Gewichte der Wertigkeitsfunktion anzupassen.⁶⁰ Mittels Erfahrungswiederholung wird eine bessere Dateneffizienz zum einen durch die Wiederholung, und zum anderen durch das Auflösen von Korrelationen zwischen aufeinanderfolgenden Zuständen erzielt.⁶¹ Das Zielnetzwerk beinhaltet die zunächst fixen Gewichte der Strategie, welche lediglich nach einer festen Anzahl an Schritten angepasst werden, um die Fluktuation der approximierten Q-Funktion auszugleichen.⁶² Die Stärke von DQL liegt in der kompakten Repräsentation der Q-Funktion und der hochgradig dimensionierten Zustandsbeobachtungen durch neuronale Netze.⁶³

Trust Region Policy Optimization (TRPO) ist ein strategiebasierender Gradientenalgorithmus zur effektiven Optimierung großer nicht linearer Strategien wie z. B. neuronale Netze.⁶⁴ Der Algorithmus weist dabei die zwei Varianten single-path, welche im modellfreien Kontext angewendet werden kann, und *vine* auf, welche sich nur in Simulationen eignet, da das System zu bestimmten Zuständen gespeichert wird.⁶⁵ Die Varianten unterscheiden sich im Schätzverfahren des Gradienten, wobei die single-path Methodik diesen anhand einer Aktions-Zustandskette der initialen Verteilung bestimmt, während unter *vine* eine Teilmenge verschiedener Aktion-Zustandspaare verwendet wird.⁶⁶ Die daraus entstehenden Anpassungen der Strategieparameter θ zwischen der alten und neuen Strategie werden mittels Kullback-Leibler-Divergenz bemessen und kontrolliert.⁶⁷ Eine Weiterentwicklung des TRPO Algorithmus ist Proximal Policy Optimization (PPO), durch dessen Verwendung die Implementierung erleichtert und die Datenprobenkomplexität verbessert wird.⁶⁸ PPO passt die Beschränkung der KL-Divergenz an, indem die Wahrscheinlichkeitsmasse der Optimierung außerhalb der gesetzten Hyperparametergrenzen bearbeitet, oder die Größe der KL-Divergenz bestraft wird.⁶⁹ Durch diese Anpassung der KL-Divergenz Beschränkung ermöglicht der PPO Algorithmus mehrfache Berechnungen des Gradienten anhand eines Datenobjektes.⁷⁰

Zusätzlich zu den bisherigen Algorithmen kann durch jeweilige asynchrone Varianten, welche den Trainingsprozess des Agenten parallelisieren, die Stabilität des Trainings verbessert werden.⁷¹ Einer dieser asynchronen Algorithmen im Bereich der Akteur-Kritiker Methodik stellt der asynchronous advantage actor-critic (A3C) Algorithmus dar, welcher anstelle von Erfahrungswiederholung

⁵⁹Vgl. Arulkumaran u. a. 2017, S. 7

⁶⁰Vgl. Mnih/Kavukcuoglu u. a. 2013, S. 4

⁶¹Vgl. Mnih/Kavukcuoglu u. a. 2013, S. 4f.

⁶²Vgl. Arulkumaran u. a. 2017, S. 7

⁶³Vgl. Arulkumaran u. a. 2017, S. 7

⁶⁴Vgl. Schulman/Levine u. a. 2015, S. 1

⁶⁵Vgl. Schulman/Levine u. a. 2015, S. 1

⁶⁶Vgl. Schulman/Levine u. a. 2015, S. 4

⁶⁷Vgl. Huang u. a. 2017, S. 4

⁶⁸Vgl. Schulman/Wolski u. a. 2017, S. 1

⁶⁹Vgl. Schulman/Wolski u. a. 2017, S. 3f.

⁷⁰Vgl. Schulman/Wolski u. a. 2017, S. 4

⁷¹Vgl. Mnih/Badia u. a. 2016, S. 1

mehrere Agenten parallel in unterschiedlichen Umgebungsinstanzen trainiert.⁷² A3C kombiniert die Berechnung der relativen Wertigkeit einer Aktion, anstelle der absoluten Wertigkeit durch die Q-Funktion, mit der Akteur-Kritiker Struktur und lässt sich auf einzelnen sowie verteilten Systemen einsetzen.⁷³ Wird mit dem Algorithmus lediglich ein Agent trainiert, wird dies häufig auch als advantage actor-critic (A2C) referenziert.⁷⁴ Neben A3C ist der soft actor-critic Algorithmus ein weiterer Algorithmus der Akteur-Kritiker Methodik, welcher die kumulierte Belohnung, aber auch die Entropie maximiert.⁷⁵ Durch dieses Konzept der Maximierung der Belohnung und der Entropie ergeben sich die Vorteile eines stärker erkundenden Algorithmus, welcher gleichzeitig mehrere nahezu optimale Lösungen erfassen kann.⁷⁶ Zur Evaluation wird die auf der Belohnung und Entropie basierte Wertigkeit iterativ anhand des Bellman Operators nach Bellman 1966 bestimmt, und anschließend die Strategie hinsichtlich der exponentiellen Q-Funktion angepasst.⁷⁷ Durch die abwechselnde Approximation der Q-Funktion des Kritikers und der Strategie des Akteurs, anstelle der Berechnung dieser bis zur Konvergenz, kann der soft actor-critic Algorithmus besonders mit großen kontinuierlichen Handlungsbereichen umgehen.⁷⁸

2.2.3 Abgrenzung zu Multi-Agent Reinforcement Learning (MARL) Algorithmen

Innerhalb dieses Unterkapitels soll der beschriebene Aufbau von RL Algorithmen und deren Optimierungsproblem zu den von MARL Systemen abgegrenzt werden. Bei MARL Systemen wird anstatt eines Agenten eine Menge von Agenten eingesetzt, welche alle mit ihrer Umgebung interagieren, um den Weg der Zielerreichung zu lernen.⁷⁹ Dieser Ansatz dient dazu, die Vielzahl an Problemstellungen in der echten Welt, welche nicht vollständig durch einen einzelnen Agenten lösbar sind, zu bearbeiten.⁸⁰ Einsatzgebiete von MARL sind dabei unter anderem das Routing von Netzwerkpaketen, Wirtschaftsmodellierung oder zusammenhängende Robotersysteme.⁸¹ Je nach Ziel und der demnach definierten Belohnungsfunktion können die Agenten auf die drei unterschiedlichen Arten vollständig kooperativ, vollständig kompetitiv und der Mischung aus beiden miteinander interagieren.⁸² Aus den einzelnen Interaktion jedes Agenten mit derselben Umgebung ergibt sich der Unterschied, dass die Umgebungsdynamik aus der Kombination aller Aktionen der Agenten beeinflusst wird anstatt aus der Aktion des einzelnen Agenten.⁸³ Da dieser Effekt die Annahme der Stationarität von Markov Entscheidungsprozessen verletzt, bedarf die Umgebung einer anderen Representation.⁸⁴ Ein Konzept, das dafür häufig verwendet wird,

⁷²Vgl. Mnih/Badia u. a. 2016, S. 1

⁷³Vgl. Arulkumaran u. a. 2017, S. 9

⁷⁴Vgl. Arulkumaran u. a. 2017, S. 9

⁷⁵Vgl. Haarnoja u. a. 2018, S. 1

⁷⁶Vgl. Haarnoja u. a. 2018, S. 3

⁷⁷Vgl. Haarnoja u. a. 2018, S. 4

⁷⁸Vgl. Haarnoja u. a. 2018, S. 4

⁷⁹Vgl. Wong u. a. 2022, S. 6

⁸⁰Vgl. Canese u. a. 2021, S. 1

⁸¹Vgl. Canese u. a. 2021, S. 1

⁸²Vgl. Canese u. a. 2021, S. 8f.

⁸³Vgl. Wong u. a. 2022, S. 2

⁸⁴Vgl. Wong u. a. 2022, S. 6

ist das Markov Spiel, welches sich anders als der Entscheidungsprozess durch einen mehrdimensionalen Aktions- und Belohnungsraum aus der Kombination aller N Agenten auszeichnet.⁸⁵ Betrachtet man die Limitierungen von MARL, erkennt man aus den beschriebenen Punkten die Herausforderungen der nicht vorhandenen Stationarität und der Skalierbarkeit, welchen sich die Herausforderung der teilweisen Beobachtbarkeit der Umgebung anschließt.⁸⁶

2.2.4 Limitierungen und Herausforderungen von RL

Trotz signifikanter Errungenschaften birgt der Einsatz der besprochenen RL Algorithmen weiterhin Limitierungen und Risiken für ungewolltes Verhalten.⁸⁷

Eine der Herausforderungen zeigt sich bei der Representation der Agentenumwelt, da RL stark auf diesem Konzept basiert.⁸⁸ Daraus ergibt sich die Aufgabe, die Umwelt und deren Verhalten sowie die Wahrnehmung durch den Agenten realitätsgetreu und präzise zu gestalten.⁸⁹ Neben der Definition und Wahrnehmung des Umweltverhaltens ist die Spezifikation des Ziels des Agenten ein ebenso kritischer Teil, da unerwartete Intentionen aus der Zielstellung abgeleitet werden könnten.⁹⁰ Zusätzlich teilen RL Algorithmen auch Herausforderungen aus anderen Gebieten des maschinellen Lernens wie Genauigkeit, Interpretierbarkeit und die im Rahmen dieser Arbeit untersuchte Robustheit von Modellen.⁹¹

Eine weitere Limitierung stellt der große Suchraum an Aktionen und das unbekannte Verhalten der Umgebung dar. Dies sorgt dafür, dass die Effizienz einzelner Daten häufig sehr gering ist und die Abwägung zwischen der Exploration neuer Strategie und der Optimierung bekannter Verhaltensmuster ein wichtiger Bestandteil ist.⁹² Aufgrund der geringen Effizienz der Daten, aber des dennoch hohen Bedarfs an bewerteter Agentenerfahrung wird häufig auf simulierte Daten zurückgegriffen.⁹³ Simulierte Daten werden dabei häufig von möglichst hochqualitativen Simulationsumgebungen bereitgestellt, da zu dem hohen Bedarf der Methodik häufig Limitierungen in der Sammlung von Daten in der echten Welt bestehen.⁹⁴

Aufgrund der Bedeutung der Simulationsumgebung und dessen Handlungsraum für RL Algorithmen und deren Transfer in die echte Welt, werden in den beiden nachfolgenden Kapiteln imitierendes Lernen und Simulationen genauer betrachtet.

⁸⁵Vgl. Canese u. a. 2021, S. 4

⁸⁶Vgl. Canese u. a. 2021, S. 9ff.

⁸⁷Vgl. Li 2019, S. 7

⁸⁸Vgl. Sutton/Barto 2018, S. 8

⁸⁹Vgl. Sutton/Barto 2018, S. 7

⁹⁰Vgl. Li 2019, S. 7

⁹¹Vgl. Li 2019, S. 7

⁹²Vgl. Li 2019, S. 7

⁹³Vgl. Zhao/Queralta/Westerlund 2020, S. 7

⁹⁴Vgl. Li 2019, S. 8

2.3 Imitierendes Lernen

Für Algorithmen des maschinellen Lernens wie z. B. RL, erweist die Skalierung zu höherdimensionalen Handlungsräumen ein großes Optimierungsproblem auf.⁹⁵ Hinzu kommen Schwächen in der Anwendung von verstärkendem Lernen wie Handlungsstrategien in unerfahrenen Umgebungszuständen.⁹⁶ Imitierendes Lernen (IL) kann hier als ein Vorschrift zur Optimierung von RL Modellen dienen, womit sich signifikante Belohnungsverbesserungen erzielen lassen.⁹⁷ Ein Anwendungsgebiet des imitierenden Lernens ist die Robotik, in welchem ein chirurgischer Roboter durch Teleoperation bereits erfolgreich chirurgische Eingriffe demonstrieren konnte.⁹⁸ Allgemeiner wird IL häufig eingesetzt um menschliches Verhalten nachzuahmen und zu deren Aktionen zu beschleunigen.⁹⁹ Das Ziel von IL ist hierbei eine Zuordnung zu erlernen zwischen den Beobachtungen und Aktionen von Expertendemonstrationen.¹⁰⁰ Dadurch soll ein Roboter das ihm demonstrierte Verhalten reproduzieren und auch in unbekannten Szenarien generalisieren können.¹⁰¹ Problemstellungen können somit auf das Demonstrieren und Anlernen der Fähigkeit beschränkt werden, ohne dass explizite Anweisungen oder Belohnungsfunktionen programmiert werden müssen.¹⁰² IL stellt dabei einen passiven Ansatz dar, die Zielstrategie durch Beobachten der vollständigen Ausführungsabläufe zu erlernen.¹⁰³ Dabei wird eine Ersatz-Verlustfunktion optimiert, welche die Diskrepanz zwischen der gewählten Aktion der parametrisierten Strategie π_θ und der des Experten ψ misst.¹⁰⁴

Der typische Einsatzprozess von IL beinhaltet das Akquirieren von Demonstrationsbeispielen, deren Kodierung als Zustands-Aktions Paare und das Optimieren einer Strategie.¹⁰⁵ Einzelne Prozessschritte können sich dabei je nach Anwendungsgebiet unterscheiden, oder Nachfolgeprozesse wie die Feinabstimmung mittels RL enthalten, was an Abbildung 3 deutlich wird.

Algorithmen des imitierenden Lernens unterscheiden sich in der Literatur nach den Begriffen *on-policy* und *off-policy*.¹⁰⁷ Dabei werden zum einen erlernte Strategien in der Umgebung ausgeführt und mittels Experten evaluiert, zum anderen werden lediglich die Expertendemonstrationen verwendet.¹⁰⁸ DAgger stellt ein *on-policy* Algorithmus dar, welcher zu jeder Iteration die Optimierung aller bisher beobachteten Aktions-Zustands Paare vornimmt.¹⁰⁹ Behavioral Cloning als

⁹⁵Vgl. Hussein u. a. 2017, S. 3

⁹⁶Vgl. Attia/Dayan 2018, S. 1

⁹⁷Vgl. Hussein u. a. 2017, S. 4

⁹⁸Vgl. Gao u. a. 2014

⁹⁹Vgl. Attia/Dayan 2018, S. 1

¹⁰⁰Vgl. Hussein u. a. 2017, S. 1

¹⁰¹Vgl. Fang u. a. 2019, S. 365

¹⁰²Vgl. Hussein u. a. 2017, S. 1

¹⁰³Vgl. Attia/Dayan 2018, S. 2

¹⁰⁴Vgl. Balakrishna u. a. 2020, S. 2f.

¹⁰⁵Vgl. Hussein u. a. 2017, S. 3

¹⁰⁶Enthalten in: Hussein u. a. 2017, S. 4

¹⁰⁷Vgl. Balakrishna u. a. 2020, S. 3

¹⁰⁸Vgl. Balakrishna u. a. 2020, S. 3

¹⁰⁹Vgl. Attia/Dayan 2018, S. 5



Abb. 3: Ablaufdiagramm des imitierenden Lernens¹⁰⁶

off-policy Algorithmus erlernt hingegen die direkte Übersetzung der Zustands- und Umgebungsinformationen zu ihren Aktionen, ähnlich einem Label unter überwachtem Lernen.¹¹⁰

Insgesamt kann der Einsatz von IL unter anderem folgende Vorteile verzeichnen.¹¹¹

- Verbesserung der Anpassungsfähigkeit hinsichtlich neuer Umgebungen
- Erhöhung der Lerneffizienz aufgrund schnellerer Übertragung von Wissen
- Kompatibilität mit anderen Algorithmen des maschinellen Lernens wie z. B. RL

Aufgrund der interdisziplinären Natur von IL ergeben sich jedoch wie folgt auch eine Reihe von Herausforderungen.¹¹²

- Der IL Prozess ist sowohl während der Sammlung von Demonstrationsdaten sowie während des Trainingsverlaufs anfällig für Rauschen und Fehler der Sensorik.
- Fähigkeiten, Aufbau und Freiheitsgrade des lernenden Modells und des Experten müssen bestmöglich übereinstimmen.
- Anwendungen lassen sich aufgrund ihres Echtzeitbezugs und Limitierungen der Rechenkapazitäten nur schwierig hinsichtlich hoher Freiheitsgrade skalieren.

¹¹⁰Vgl. Fang u. a. 2019, S. 4

¹¹¹Vgl. Fang u. a. 2019, S. 1

¹¹²Vgl. Hussein u. a. 2017, S. 4f.

2.4 Simulationsumgebungen für RL

Anders als im klassischen Bereich des maschinellen Lernens wie dem überwachten- und unüberwachten Lernen, werden beim verstärkenden Lernen viele der Testdatensätze nicht aus der echten Welt akquiriert.¹¹³ Um entsprechend realistische Daten für das Training bereitzustellen, werden Simulationsumgebungen in Abhängigkeit von ihrer RL Anwendung ausgewählt.¹¹⁴ Dennoch bleibt nahezu immer eine gewisse Diskrepanz zwischen der Dynamik in der Simulation und der Dynamik in der echten Welt.¹¹⁵ Möglichkeiten, diese Diskrepanz zu minimieren, sind zum einen das Fehlverhalten von Sensoren einzubinden oder ein reales Signal mit der virtuellen Umgebung zu verknüpfen.¹¹⁶ Dennoch lässt sich kaum garantieren, dass erlernte Strategien der Agenten sich auf nur leicht veränderte Umgebungen übertragen lassen.¹¹⁷ Anders als in der Simulation von Flüssigkeiten und deren Dynamik, bedarf RL eine reaktive Umgebung, dessen Verhältnis der Simulationszeit zur echten Zeit mindestens eins oder darüber liegt.¹¹⁸ RL kann von einem erhöhten Echtzeitfaktor profitieren, trotz der damit einhergehenden verringerten Präzision.¹¹⁹

2.4.1 Definitionen von Simulationsumgebungen

Ausgehend von der Literaturrecherche zeigt sich, dass in der Forschungsliteratur die allgemeine Definition von Simulationen kaum aufgegriffen wird. Eine mögliche Definition nach Maria 1997 wird wie folgt dargelegt:

Eine Simulation eines existierenden Systems stellt die Anwendung eines Modells dar, welches konfigurierbar zu experimentellen Zwecken das eigentliche System vertritt, um wirtschaftliche oder systematische Herausforderungen des existierenden Systems zu umgehen. Das Model wird in diesem Kontext definiert als Repräsentation des Aufbaus und der Verhaltensweise des existierenden Systems.

Innerhalb bestimmter Anwendungsgebiete, wie der Medizin und der Pflege, werden zusätzlich virtuelle Simulationen wie nachstehend definiert.

Unter virtuellen Simulationen versteht man eine digitale Lernumgebung, welche durch teilweiser Immersion eine wahrnehmbare Erfahrung bereitstellt.¹²⁰

¹¹³Vgl. Zhang/Wu/Pineau 2018, S. 1

¹¹⁴Vgl. Körber u. a. 2021, S. 7

¹¹⁵Vgl. Bharadhwaj u. a. 2019, S. 1

¹¹⁶Vgl. Zhang/Wu/Pineau 2018, S. 1

¹¹⁷Vgl. Bharadhwaj u. a. 2019, S. 1

¹¹⁸Vgl. Körber u. a. 2021, S. 3

¹¹⁹Vgl. Körber u. a. 2021, S. 3

¹²⁰Vgl. Foronda 2021, S. 1

2.4.2 Entwicklung von Simulationsumgebungen für RL Anwendungen

Im weiteren Teil dieses Kapitels wird aufbauend auf den zuvor angeführten Definitionen, die Entwicklung von Simulationen betrachtet. Allgemein lässt sich dieser Entwicklungsprozess in die folgenden Teilschritte gliedern:¹²¹

1. Identifikation der Herausforderungen im existierenden System und Ableitung von Anforderungen für die Simulation.
2. Zielgruppe, Funktionsrahmen und quantitative Bewertungskriterien der Simulation definieren.
3. Analyse des zu simulierenden Systemverhaltens durch Sammeln und Verarbeiten von realen Daten des existierenden Systems.
4. Entwicklung einer schematischen Darstellung des Modells und dessen Überführung in nutzbare Software.
5. Validierung des Modells durch bspw. den Vergleich mit dem existierenden System.
6. Dokumentierung des Modells, dessen Variablen, Metriken und getroffene Annahmen.

Die Entwicklung von Simulationen wurde in der Forschungsliteratur besonders durch den Fortschritt im Bereich des verstärkenden Lernens vorangetrieben, da der Vergleich von RL-Algorithmen zuverlässige Benchmarks in Form von Simulationsumgebungen benötigt.¹²² Aus dieser Motivation wurde 2016 durch die *OpenAI* der *OpenAI Gym* Werkzeugkasten entwickelt, welcher eine Sammlung an Benchmarksimulationen mit einer einheitlichen Schnittstelle für RL Algorithmen enthält.¹²³ Seither wurde diese definierte Schnittstelle vielfach verwendet, um RL Umgebungen mit dem Ziel zu entwickeln, diese zu publizieren und dessen Wiederverwendung zu ermöglichen.¹²⁴ Die Entwicklung dieses Softwareframeworks wurde nach der Version 0.26.0 im Jahr 2022 durch ein neues Team unter dem Namen Gymnasium weitergeführt.¹²⁵ Die Schnittstelle ist definiert als Python Klasse *gym.Env*, von welcher weitere Klassen erben und die vorgeschriebenen Funktionen zum Zeitschritt und zum Zurücksetzen der Simulation implementieren.¹²⁶ Der Werkzeugkasten von OpenAI fokussiert sich auf einen „Episoden ähnlichen Rahmen“, in welchem der Agent durch zunächst zufälliges Auswählen von Interaktionen lernt.¹²⁷ Weitere Entwicklungsentscheidungen des OpenAI Gym Werkzeugkastens umfassen z. B. die bewusst fehlende Schnittstelle des Agenten, die strikte Versionierung der Umgebung oder die standardmäßige Simulationsüberwachung.¹²⁸

¹²¹Vgl. Maria 1997, S. 8f.

¹²²Vgl. Brockman u. a. 2016, S. 1

¹²³Vgl. Brockman u. a. 2016, S. 1

¹²⁴Vgl. Schuderer/Bromuri/van Eekelen 2021, S. 4

¹²⁵Vgl. GitHub 4/4/2023

¹²⁶Vgl. Schuderer/Bromuri/van Eekelen 2021, S. 4

¹²⁷Vgl. Brockman u. a. 2016, S. 1

¹²⁸Vgl. Brockman u. a. 2016, S. 2f.

Werden Lernumgebungen nach der Gym Schnittstelle oder nach eigener Definition für RL Anwendungen eingesetzt, kann sich deren Gestaltung unterschiedlich auf die Leistung der Anwendung auswirken.¹²⁹ Eine enge initiale Wahrscheinlichkeitsverteilung des Umgebungszustandes kann die Lerneffizienz erhöhen, wohingegen eine weite Wahrscheinlichkeitsverteilung positiv die Robustheit der erlernten Strategie beeinflusst.¹³⁰ Die Robustheit kann zusätzlich durch die Einbindung von Fehlverhalten in der Wahrnehmung der Umgebung beeinflusst werden, da auch in realen Szenarien ein Risiko für Fehlverhalten besteht.¹³¹ Im Bereich der Robotik bzw. in der Simulation von Bewegungen, kann auch durch die Gestaltung des Aktionsraumes, basierend auf elektrischer Regelungstechnik mittels PID-Regler, anstatt basierend auf Drehmomenten ein effizienterer Lernprozess stattfinden.¹³²

Neben den beschriebenen Eigenschaften von Umgebungen für verstärkendes Lernen unterliegen auch die verwendeten Simulationen bestimmten Merkmalen, welche in der Entwicklung zu berücksichtigen sind. Laut einer Umfrage nach Ivaldi/Padois/Nori 2014 sind diese wichtigsten Eigenschaften die Stabilität, Geschwindigkeit, Präzision, Genauigkeit, Bedienbarkeit und der Ressourcenverbrauch. Die Entwicklung des Modells, welches das existierende System ersetzt, sollte sich demnach möglichst positiv auf die beschriebenen Eigenschaften auswirken. Neben den beschriebenen leistungsbezogenen Merkmalen sind die folgenden weiteren Kriterien mitunter die wichtigsten zur Auswahl einer Simulation:

Rank	Most important criteria
1	Simulation very close to reality
2	Open-source
3	Same code for both real and simulated robot
4	Light and fast
5	Customization
6	No interpenetration between bodies

Tab. 2: wichtigsten Kriterien zur Auswahl von Simulatoren¹³³

Aus Tabelle zwei lässt sich entnehmen, dass besonders die Nähe zur Realität ein wichtiges Auswahlkriterium ist. Im Kontext von verstärkendem Lernen im Robotik Bereich ist ein wichtiger Baustein die Physik-Engine zur Modellierung von Dynamiken.¹³⁴

2.4.3 Aktuelle Physik-Engines und Simulationsanwendungen

Innerhalb dieses Abschnittes wird aufgrund der Bedeutung der Physik-Engine für den Grad der Simulationsrealität eine Auswahl der aktuellen Physik-Engines und deren Simulationsanwendung betrachtet.

¹²⁹Vgl. Reda/Tao, T./van de Panne 2020, S. 1

¹³⁰Vgl. Reda/Tao, T./van de Panne 2020, S. 3

¹³¹Vgl. Yan Duan u. a. 2016, S. 2

¹³²Vgl. Reda/Tao, T./van de Panne 2020, S. 7

¹³³Enthalten in: Ivaldi/Padois/Nori 2014, s. 4

¹³⁴Vgl. Ayala u. a. 2020, S. 2

Gazebo

Gazebo ist eine durch die Open Source Robotics Foundation entwickelte Simulationsanwendung, welche mehrere Physik-Engines unterstützt.¹³⁵ Mittels Gazebo lassen sich Interaktionen zwischen Robotern in Innen- und Außenbereichen unter realistischer Sensorik simulieren.¹³⁶ Die unterstützten Physik-Engines umfassen Bullet, Dynamic Animation and Robotics Toolkit (DART), Open Dynamics Engine (ODE) und Simbody.¹³⁷

MuJoCo

MuJoCo stellt eine Physik-Engine für modellbasierte Steuerung dar, dessen Objekte durch C++ oder XML definiert und Gelenkzustände im Koordinatensystem beschrieben werden.¹³⁸ Diese beschriebene Eigenschaft lässt sich auch aus dem Namen als Abkürzung für **M**ulti-**J**oint dynamics with **C**ontact ableiten.¹³⁹ Die MuJoCo Anwendung ist lizenziert, was den Besitz einer Lizenz für die Installation oder die Virtualisierung innerhalb eines Containers voraussetzt.¹⁴⁰

PyBullet

PyBullet basiert als Simulationssoftware auf der Bullet-Engine und fokussiert sich funktional auf die Anwendung von RL im Robotikbereich.¹⁴¹ Die Bullet-Engine ist hingegen eine offene Softwarebibliothek, welche neben verstärkenden Lernen auch bei Computeranimationen angewendet wird.¹⁴² Die Handhabung von PyBullet profitiert von der ausgiebigen Dokumentation, der großen Entwicklergemeinschaft und der Unterstützung von verschiedenen Dateiformaten wie SDA, URDF und MJCF zur Einbindung von Objekten.¹⁴³

2.5 Simulation der Steuerungsaufgabe von Quadroptern

Die Popularität von unbemannten Flugzeugen im letzten Jahrzehnt nahm besonders im Bereich der Quadroptern zu, sodass durch die sinkenden Kosten von Sensorik und Minicomputern, zahlreiche zukunftssträchtige Ergebnisse und Anwendungen erforscht worden sind.¹⁴⁴ Daher gilt dieses Kapitel der theoretischen Betrachtung, dessen Eigenschaften und deren Simulation. Anwendungsgebiete beinhalten z. B. die Landwirtschaft, den Pakettransport oder die Überwachung von großflächiger Infrastruktur wie Stromnetze.¹⁴⁵ Eine Simulation von unbemannten Flugzeugen

¹³⁵Vgl. Ivaldi/Padois/Nori 2014, S. 7

¹³⁶Vgl. Ayala u. a. 2020, S. 4

¹³⁷Vgl. Körber u. a. 2021, S. 3

¹³⁸Vgl. Todorov/Erez/Tassa 2012, S. 1

¹³⁹Vgl. Todorov/Erez/Tassa 2012, S. 2

¹⁴⁰Vgl. Körber u. a. 2021, S. 3

¹⁴¹Vgl. Körber u. a. 2021, S. 3

¹⁴²Vgl. Ivaldi/Padois/Nori 2014, S. 7

¹⁴³Vgl. Körber u. a. 2021, S. 6

¹⁴⁴Vgl. Koch u. a. 2018, S. 1

¹⁴⁵Vgl. Deshpande/Kumar, R. u. a. 2020, S. 1

stellt eine Flugumgebung und vielseitige Sensorik bereit und kann je nach Anwendung Effekte wie Wind, Wolken und Niederschlag einbeziehen.¹⁴⁶

2.5.1 Flugdynamiken eines Quadropters

Die Simulation des Quadropters stellt eine Simulation eines Flugkörpers mit drei Rotations- und drei Translationsbewegungen und demnach insgesamt sechs verschiedenen Freiheitsgraden dar.¹⁴⁷ Ein Quadropter ist ein fester Körper mit vier befestigten Rotoren, welche sich ausschließlich in eine Richtung drehen und positiven Schub in die Z-Achse des Körpers ausüben können.¹⁴⁸ Die vier Rotoren werden als + oder X Konfiguration entweder direkt in Richtung der X- und Y-Achse (+), oder um 45° gedreht (X) an der Drohne befestigt.¹⁴⁹ Jede Bewegung der sechs verschiedenen Arten wird durch die unterschiedliche Ansteuerung, der im Fall des Quadropters, vier Rotoren getätigt. Aus den verschiedenen starken auftrieben resultieren die drei Rotationsbewegungen Rollen, Nicken und Gieren.¹⁵⁰



Abb. 4: Rotationsbewegungen eines Quadropters¹⁵¹

Rollen wird, wie aus Abbildung 4 hervorgeht, durch unterschiedlichen Auftrieb der zwei linken und rechten Rotatoren, das Nicken durch Unterschiede der vorderen und hinteren Rotatoren hervorgerufen. Das Gieren bzw. die Rotation um die Z-achse wird durch die stärkere Rotation der sich im Uhrzeiger drehenden oder der sich gegen den Uhrzeiger drehenden Rotoren bewerkstelligt.

Diese Rotationsbewegungen werden mathematisch als Matrix der Eulerschen Winkel repräsentiert, welche die Folge der einzelnen Drehungen entlang der X-, Y-, und Z-Achse enthält.¹⁵² Das Zusammenspiel dieser Rotationsbewegungen, dargestellt anhand Formel vier, mit der durch die Rotatoren entlang der Z-Achse der Drone erzeugte Schubkraft und den Rollraten, beschrieben in Formel fünf, sorgt für die Bewegung der Drohne in Richtung der zukünftigen Koordinaten nach Formel sechs.¹⁵³

¹⁴⁶Vgl. Hentati u. a. 2018, S. 1496

¹⁴⁷Vgl. Koch u. a. 2018, S. 2

¹⁴⁸Vgl. Molchanov u. a. 2019, S. 3

¹⁴⁹Vgl. Koch u. a. 2018, S. 2

¹⁵⁰Vgl. Koch u. a. 2018, S. 2

¹⁵¹Enthalten in: Koch u. a. 2018, S. 2

¹⁵²Vgl. Deshpande/Kumar, R. u. a. 2020, S. 3

¹⁵³Vgl. Deshpande/Minai/Kumar, M. 2021, S. 2

$$R = \begin{pmatrix} C_\psi C_p & S_\xi S_p C_\psi - S_\psi C_\xi & S_\xi S_\psi + S_p C_\xi C_\psi \\ S_\psi C_p & S_\xi S_\psi S_p + C_\xi C_\psi & -S_\xi C_\psi + S_\psi S_p C_\xi \\ -S_p & S_\xi C_p & C_\xi C_p \end{pmatrix} \quad (4)$$

In Formel vier wird der Rotationszustand der Drohne zu den Weltachsen durch die Sinus- und Cosinuswinkel S_a und C_a repräsentiert, wobei für a die Art der Rotation also Rollen (ξ), Nicken (p) und Gieren (ψ) eingesetzt wird.¹⁵⁴

$$I \begin{pmatrix} \dot{p} \\ \dot{q} \\ \dot{r} \end{pmatrix} = \begin{pmatrix} l(F_1 + F_2 - F_3 - F_4) \\ l(-F_1 + F_2 + F_3 - F_4) \\ -M_1 + M_2 - M_3 + M_4 \end{pmatrix} - \begin{pmatrix} p \\ q \\ r \end{pmatrix} \times I \begin{pmatrix} p \\ q \\ r \end{pmatrix} \quad (5)$$

Formel fünf inkludiert in der Matrix I die Trägheitsmomente um die X-, Y- und Z-Achsen in Abhängigkeit der Rollrate p , Nickrate q und Gierrate r .¹⁵⁵ Der Schub der in der Länge l vom Schwerpunkt entfernten Rotatoren wird mittels F_i und das Drehmoment mit $M_i, \forall i \in \{1, 2, 3, 4\}$ gekennzeichnet.¹⁵⁶

$$m \begin{pmatrix} \ddot{x} \\ \ddot{y} \\ \ddot{z} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ -mg \end{pmatrix} + R \begin{pmatrix} 0 \\ 0 \\ \sum_{i=1}^4 F_i \end{pmatrix} \quad (6)$$

Weiterhin ist die Beschleunigung \ddot{x} , \ddot{y} und \ddot{z} in Richtung der drei Achsen durch die Masse m und die Gravitation g beeinflusst.¹⁵⁷

2.5.2 Existierende Simulationen von Quadroptern

RotorS

Die Simulationsumgebung RotorS wurde auf Basis des Robot Operating Systems (ROS) entwickelt, um Programmtestzeiten zu verkürzen, Fehlersuche zu vereinfachen und Unfälle mit echten Mikroflugzeugen zu vermindern.¹⁵⁸ Dabei wurde die Simulation modular entworfen, sodass Komponenten wie Steuerung oder Zustandsschätzung austauschbar sind und das Hinzufügen von neuen Drohnen erleichtert wird.¹⁵⁹ Die Komponenten der Drohne stellen dabei Plug-ins der verwendeten Gazebo Physik-Engine dar, wodurch ein Mikroflugzeug aus den Teilen des Körpers, der Anzahl der Rotoren und Sensorik an fixen Position zusammengesetzt wird.¹⁶⁰ Mittels der standardmäßigen Sensorik können Informationen über die direkte und visuell gemessene Trägheit sowie über die Wegbestimmung erzielt werden.¹⁶¹ Anstelle des Wegbestimmungssensors kann

¹⁵⁴Vgl. Deshpande/Minai/Kumar, M. 2021, S. 2

¹⁵⁵Vgl. Deshpande/Minai/Kumar, M. 2021, S. 2

¹⁵⁶Vgl. Deshpande/Kumar, R. u. a. 2020, S. 3

¹⁵⁷Vgl. Deshpande/Kumar, R. u. a. 2020, S. 3

¹⁵⁸Vgl. Furrer u. a. 2016, S. 596

¹⁵⁹Vgl. Furrer u. a. 2016, S. 595

¹⁶⁰Vgl. Furrer u. a. 2016, S. 597

¹⁶¹Vgl. Furrer u. a. 2016, S. 597

auch eine Komponente zur Zustandsschätzung für hochfrequente Abfragen implementiert werden.¹⁶² Die Steuerungskomponente wird durch eine einfache Schnittstelle einer geometrischen Steuerung bedient, welche Aktionen in Form von Rotationswinkeln, Höhen oder Positionen entgegennimmt.¹⁶³

CrazyS

CrazyS stellt eine Erweiterung der Simulation RotorS auf Basis desselben ROS, um die Modellierung des Nano-Quadropters Crazyflie samt ihrer Dynamik, ihres Kontrollsystems und ihrer Sensorik dar.¹⁶⁴ Mit der Modellierung der Nanodrohne wurde gleichzeitig ein Konzept zur Erweiterung der RotorS Fähigkeiten dargelegt sowie die Entwicklung von Software-in-the-Loop (SITL) als nahezu Echtzeitüberwachung vorangetrieben.¹⁶⁵

AirSim

AirSim ist eine open-source Simulationsplattform, mit der das Ziel verfolgt wird, durch eine detaillierte Simulation die Entwicklung von RL und anderen Methoden des maschinellen Lernens voranzutreiben.¹⁶⁶ Zur Modellierung der Simulationsphysik wird die Unreal Engine 4 aufgrund ihres hohen Grades an physikalischer und visueller Realität eingesetzt.¹⁶⁷ Der Aufbau der Simulation folgt einem modularen Entwurf, welcher unter anderem die einzelnen Komponenten Fahrzeug, Umgebung, Physik-Engine, Sensorik und Darstellungsschnittstelle beinhaltet.¹⁶⁸ Die Schnittstelle des Fahrzeugs erlaubt eine Steuerung über viele Betätigungselemente und deren Eigenschaftsparameter wie Masse, Trägheit, Widerstand oder Reibung.¹⁶⁹ Ein Fahrzeug ist dabei durch die Umgebungskomponente beeinflusst, welche physikalische Effekte wie Gravitation, Luftwiderstand, Luftdruck und magnetische Felder simuliert.¹⁷⁰ Die Umgebung wird für das Fahrzeug wahrnehmbar durch die Modellierung von Sensoren wie GPS, Beschleunigungsmesser, Gyroskop, Barometer und Magnetometer.¹⁷¹

gym-pybullet-drones

Eine weitere nach der Gym Schnittstelle definierte Simulationsumgebung ist gym-pybullet-drones.¹⁷² Basierend auf der Bullet Physik-Engine ermöglicht die Simulation unter anderem das visuell basierte Training von mehreren Agenten mittels RL unter realistischer Modellierung von Kollisionen und aerodynamischen Effekten.¹⁷³ Die Wahl der Physik-Engine wurde aufgrund des CPU und

¹⁶²Vgl. Furrer u. a. 2016, S. 598

¹⁶³Vgl. Furrer u. a. 2016, S. 598

¹⁶⁴Vgl. Silano/Iannelli 2019, S. 81

¹⁶⁵Vgl. Silano/Iannelli 2019, S. 82

¹⁶⁶Vgl. Shah u. a. 2017, S. 2

¹⁶⁷Vgl. Shah u. a. 2017, S. 1

¹⁶⁸Vgl. Shah u. a. 2017, S. 3

¹⁶⁹Vgl. Shah u. a. 2017, S. 5

¹⁷⁰Vgl. Shah u. a. 2017, S. 6

¹⁷¹Vgl. Shah u. a. 2017, S. 9

¹⁷²Vgl. Panerati u. a. 2021, S. 1

¹⁷³Vgl. Panerati u. a. 2021, S. 1

GPU basierenden Renderings, des Kollisionsmanagements und der Kompatibilität mit dem Unified Robot Description Format (URDF) getroffen¹⁷⁴ Durch die Kompatibilität mit dem URDF Format kann die standardmäßige Simulation des Drohnenmodells Bitcraze Crazyflie 2.x um weitere Nanoquadropters erweitert werden.¹⁷⁵

2.6 Robustheit und Stabilität von Strategien des verstärkenden Lernens

Anders als in Simulationen, lassen sich in der echten Welt häufig Unsicherheiten, Störeinflüsse und grundlegende Veränderungen der Umgebung wahrnehmen, für welche die Methoden des RL standardmäßig nicht robust genug sind.¹⁷⁶ Im nachfolgenden Kapitel wird daher genauer dargestellt, was unter der Robustheit von Algorithmen des verstärkenden Lernens verstanden wird, was Kenngrößen sind und in welchem Kontext diese erfasst werden können.

2.6.1 Definitionen von Robustheit und Stabilität

In der aktuellen Forschungsliteratur findet sich nur wenig Gemeinsamkeit innerhalb der unterschiedlichen Definitionen von Stabilität und Robustheit.¹⁷⁷ Die Definition der Robustheit im Kontext von verstärkendem Lernen wird verschieden interpretiert, wie z. B. als Robustheit gegen Störeinflüsse, Beeinflussung der Belohnung, oder Umgebungsunterschiede.¹⁷⁸ Pullum 2022 definiert Stabilität und Robustheit im Kontext der Literaturanalyse wie folgt:

*Stabilität ist eine Eigenschaft des lernenden Algorithmus, die sich auf dessen Leistungsvarianz bezieht und bei geringer Varianz auf ein stabiles Modell hinweist.*¹⁷⁹

*Robustheit im Kontext von Software, referenziert eine Eigenschaft eines Systems, welches nicht nur ausschließlich unter normalen, sondern auch unter außergewöhnlichen Bedingungen, die die Annahmen des Entwicklers übersteigen, gut funktioniert.*¹⁸⁰

Moos u. a. 2022 beschreibt die Robustheit in seiner Literaturanalyse als Fähigkeit mit Variationen und Unsicherheiten in der Umgebung umgehen zu können, wobei Unsicherheiten häufig variierende physische Parameter darstellen.¹⁸¹

¹⁷⁴Vgl. Panerati u. a. 2021, S. 3

¹⁷⁵Vgl. Panerati u. a. 2021, S. 3

¹⁷⁶Vgl. Moos u. a. 2022, S. 1

¹⁷⁷Vgl. Pullum 2022, S. 5

¹⁷⁸Vgl. Liu u. a. 2023, S. 2

¹⁷⁹Vgl. Pullum 2022, S. 5

¹⁸⁰Vgl. Pullum 2022, S. 5

¹⁸¹Vgl. Moos u. a. 2022, S. 1

2.6.2 Metriken der Robustheit

Werden Algorithmen und der Erfolg von Veränderungen dieser Experimente betrachtet, kommt es oftmals zu einem Vergleich zwischen Leistung und Stabilität, womit die Belohnung die einzige Kenngröße bildet.¹⁸² Wird diese Metrik im Kontext unterschiedlicher Umgebungen überprüft, kann allerdings so auch die Robustheit von RL Algorithmen betrachtet werden.¹⁸³ Neben der Betrachtung der Belohnung lassen sich auch weitere quantitative sowie qualitative Kenngrößen untersuchen, welche die Robustheit und Stabilität eines Algorithmus innerhalb der Umgebung widerspiegeln.¹⁸⁴ In der Literaturanalyse nach Pullum 2022 werden quantitative Metriken zusätzlich nach internen Kenngrößen, welche den Trainingsprozess beschreiben, und externen Kenngrößen, welche die Modellqualität repräsentieren, klassifiziert sowie folgende Tabellen drei und vier zu dessen Metriken angeführt.¹⁸⁵

Internal Quantitative Metric	Behavior	Total citations
Reward or Score – magnitude, mean/ variance, variation in average reward, time to threshold, episode duration	Stability, Robustness	75
Policy entropy	Stability	2
Variations in control strategy approximation weights	Stability, Robustness	2
Convergence rate	Stability	2
Lyapunov stability criteria calculated	Stability	1
Policy weight	Robustness	1
Regret	Robustness	1
Wasserstein function bounds calculated	Robustness	1
		85

Tab. 3: Interne quantitative Metriken, dessen gemessenes Verhalten und ihre Häufigkeit¹⁸⁶

2.6.3 Experimenteller Rahmen zur Messung der Robustheit

Die Tabellen drei und vier zeigen auf, dass trotz der Existenz weiterer Metriken, die Belohnung, deren Durchschnitt, Varianz und Entwicklung am häufigsten eingesetzt werden, um die Robustheit von RL Algorithmen zu messen. Dabei werden Experimente unter festgelegten oder optimierten Hyperparametern, in mehreren Simulationsumgebungen durchgeführt, um aus dem Vergleich derselben Strategie in unterschiedlichen Umfeldern Rückschlüsse auf die Robustheit zu

¹⁸²Vgl. Yan Duan u. a. 2016, S. 6

¹⁸³Vgl. Pinto u. a. 2017, S. 6

¹⁸⁴Vgl. Pullum 2022, S. 15

¹⁸⁵Vgl. Pullum 2022, S. 16

¹⁸⁶Ähnlich enthalten in: Pullum 2022, S.17

¹⁸⁷Ähnlich enthalten in: Pullum 2022, S.19

External Quantitative Metric	Behavior	Total citations
Deviations/variation in other (than precision, accuracy and recall) performance-related metrics	Stability, Robustness, Resilience	39
Error and failure rates/success rate	Stability, Robustness	28
Performance of tracking/trajectories estimation error; mean absolute deviation, mean square error, mean absolute percentage error, margins and magnitude of correlation coefficient	Stability, Robustness	23
Network-related timing/delay, path and link metrics, connectivity, delivery ratio, routing loops, path optimality, visitation distribution, structural Hamming distance, Small base station-serving ratio, sum-rate and 5th percentile rate	Stability, Robustness	15
Mean/average and variation inaccuracy, precision and recall, area under the receiver operating characteristic (ROC) curve (AUC)	Stability, Robustness, Resilience	12
Variance of the estimation of loss, regret	Robustness	5
		122

Tab. 4: Auszug der externen quantitativen Metriken, dessen gemessenes Verhalten und ihre Häufigkeit¹⁸⁷

ziehen.¹⁸⁸ Unterschiede in den Umgebungen können bspw. durch fixe Dynamiken wie z. B. Reibwerte während des Trainingsprozesses und unterschiedlicher Reibwerte während der Testphase realisiert werden.¹⁸⁹ Ein weiterer Ansatz kann die Sim-to-Sim Verifikation sein, bei der die optimierten Strategien in einer nicht während des Trainings verwendeten Simulationen untersucht werden.¹⁹⁰

2.7 Gegnerisches verstärkendes Lernen

Methoden des gegnerischen verstärkenden Lernens verfolgen das Ziel, Strategien zu lernen die robuster gegenüber Risiken wie beeinflusste Wahrnehmung, unbekannte Situationen oder ansteigende Umgebungskomplexität agieren.¹⁹¹ Die Robustheit gegenüber fehlerhafter Umgebungsbe-

¹⁸⁸Vgl. Pinto u. a. 2017, S. 5

¹⁸⁹Vgl. Pinto u. a. 2017, S. 6

¹⁹⁰Vgl. Molchanov u. a. 2019, S. 5

¹⁹¹Vgl. Schott/Hajri/Lamprier 2022, S. 2

trachtung kann durch Störung des Wahrnehmungszustands des trainierenden Agenten erzielt werden.¹⁹² Das Ziel ist es dabei durch gegnerische Aktionen eine veränderte Umgebungswahrnehmung herzustellen, auf deren Basis sich der lernende Agent verbessert.¹⁹³

Um die lernende Strategie besser auf unbekannte Situationen und steigende Komplexität vorzubereiten zu können, werden destabilisierende Kräfte in der Dynamik eingeführt.¹⁹⁴ Anders als beim ersten Ansatz werden dafür nicht nur die Wahrnehmung des lernenden Agenten beeinflusst, sondern direkte Einflüsse auf die Umgebung ausgeübt.¹⁹⁵ Hierbei wird der gegnerische Agent dafür belohnt, mittels Kräfteinfluss die Umgebungsdynamik zu verändern, sodass der lernende Agent an seiner Aufgabe scheitert.¹⁹⁶ Dazu kann ein zusätzlicher Agent mit z. B. gleichem Aktionsraum den gemeinsamen Umgebungszustand beeinträchtigen.¹⁹⁷ Pan u. a. 2021 zeigt eine solche Anwendung im Rahmen der Kontrolle eines Stromnetzes, bei welcher ein gegnerischer Agent für das Trennen von Verbindungen im Netz belohnt wird.¹⁹⁸ Zum Generieren von Störeinflüssen auf die Umgebungsdynamik kann das *Robust Adversarial Reinforcement Learning* (RARL) Framework verwendet werden.¹⁹⁹ Dabei wird der gegnerische Agent selbst durch RL daran angelernt, die möglichst effektivsten Destabilisierungsmaßnahmen zu finden.²⁰⁰ Formal dargestellt folgt dieses gegnerische Spiel der in Formel sieben angeführten Minimax Optimierung.²⁰¹

$$R^{1*} = \min_{\nu} \max_{\mu} E_{s_0 \sim p, a^1 \sim \mu(s), a^2 \sim \nu(s)} [\sum_{t=0}^{T-1} r^1(s, a^1, a^2)] \quad (7)$$

Zu jedem Zeitschritt t in der gegnerischen Simulation wird von beiden Spielern eine Aktion $a_t^N \sim \mu(s_t) \forall N \in \{1, 2\}$ nach der Wahrnehmung des Umgebungszustands s ausgeübt, was zum Erhalt der Belohnung $r_t^1 = r_t$ und $r_t^2 = -r_t$ führt.²⁰² Das Training erfolgt aus der abwechselnden Optimierung einer der beiden Strategien bis zu deren Konvergenz, während die jeweils andere nicht verändert wird.²⁰³

Im Kontext der Steuerung von Quadroptern greifen Zhai u. a. 2022 das *Robust Adaptive Ensemble Adversarial Reinforcement Learning Framework* (RAEARL) auf. Unter dessen Einsatz wird der Trainingsprozess des normalen und gegnerischen Agenten zu Gunsten der Kontinuität und Stabilität getrennt, und das System mit einem PID-Regler erweitert, um die Stärke des Gegners über den Trainingsverlauf anzupassen.²⁰⁴ Das getrennte Training des normalen und gegnerischen Agenten verwendet kopierte Agenten des jeweiligen Gegenparts, welche zwar dieselben

¹⁹²Vgl. Schott/Hajri/Lamprier 2022, S. 2

¹⁹³Vgl. Schott/Hajri/Lamprier 2022, S. 3

¹⁹⁴Vgl. Pinto u. a. 2017, S. 1

¹⁹⁵Vgl. Schott/Hajri/Lamprier 2022, S. 2

¹⁹⁶Vgl. Pinto u. a. 2017, S. 2

¹⁹⁷Vgl. Pinto u. a. 2017, S. 2

¹⁹⁸Vgl. Pan u. a. 2021, S. 2

¹⁹⁹Vgl. Schott/Hajri/Lamprier 2022, S. 2

²⁰⁰Vgl. Pinto u. a. 2017, S. 1

²⁰¹Vgl. Pinto u. a. 2017, S. 3

²⁰²Vgl. Pinto u. a. 2017, S. 3f.

²⁰³Vgl. Pinto u. a. 2017, S. 4

²⁰⁴Vgl. Zhai u. a. 2022, S. 2

neuronalen Netzstrukturen und initialen Parameter besitzen, jedoch schwächer sind, da deren Lernfortschritt ein Zeitschritt verschoben ist.²⁰⁵

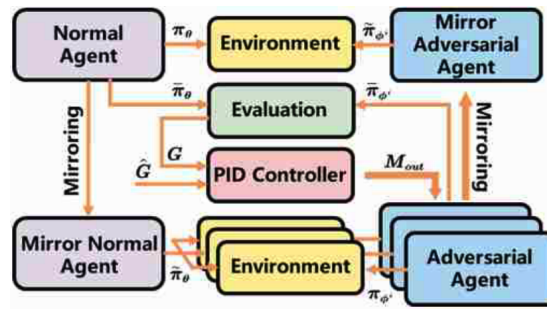


Abb. 5: Aufbau des RAEARL Frameworks²⁰⁶

Abbildung 5 zeigt den beschriebenen getrennten Aufbau des RAEARL Frameworks, in dem die Strategien der originalen Agenten mit π_θ bzw. π_{ϕ^i} und jene kopierte Strategien mit $\tilde{\pi}_\theta$ bzw. $\tilde{\pi}_{\phi^i}$ notiert sind.²⁰⁷ Zur Evaluation jeder Epoche und der Bestimmung der kumulierten Belohnung G werden die deterministischen Strategien $\bar{\pi}_\theta$ bzw. $\bar{\pi}_{\phi^i}$ einbezogen.²⁰⁸

2.8 Domain Randomization

Domain Randomization (DR) ist eine weitere Methodik, die Realitätslücke zwischen der Simulation und der echten Welt zu schließen, in dem die Simulation vielfach variiert wird während des Trainingsprozesses.²⁰⁹ Im Gegensatz zur Systemidentifikation als Prozess die Parameter bestmöglich an ein physisches System anzupassen, ermöglicht DR einen geringeren Zeitaufwand und geringere Fehleranfälligkeit.²¹⁰ Teilweise kann durch diesen Ansatz die Diskrepanz so weit geschlossen werden, dass eine Datensammlung in der echten Welt nicht erforderlich ist.²¹¹ Anhand von Komponenten der Simulationsumgebung lässt sich DR in visuelle und dynamische Randomisierung unterteilen.²¹² Weitere Methoden die Realitätslücke zu verkleinern sind progressive Netzwerke, inverse Dynamikmodelle und bayessche Methoden.²¹³

2.8.1 Das Prinzip hinter der Randomisierung von Simulationsumgebungen

DR verfolgt den Ansatz eine hohe Anzahl an Simulationsumgebungen mit randomisierten Eigenschaften zu erzeugen, in welchen die Strategie dahingehend optimiert wird, in allen Umgebungen

²⁰⁵Vgl. Zhai u. a. 2022, S. 2f.

²⁰⁶Enthalten in: Zhai u. a. 2022, S. 3

²⁰⁷Vgl. Zhai u. a. 2022, S. 3

²⁰⁸Vgl. Zhai u. a. 2022, S.3

²⁰⁹Vgl. Bharadhwaj u. a. 2019, S. 3

²¹⁰Vgl. Tobin u. a. 2017, S. 1

²¹¹Vgl. Molchanov u. a. 2019, S. 2

²¹²Vgl. Zhao/Queralt/Westerlund 2020, S. 5

²¹³Vgl. Chen u. a. 2021, S. 2

ihr Ziel bestmöglich zu erfüllen.²¹⁴ Durch das Aussetzen des Agenten gegenüber einer Vielzahl von Umgebungen, wird unter der Erwartung einer guten Leistung das Ziel verfolgt, den Agenten direkt in die physische Welt übertragen zu können.²¹⁵ Diese Erwartung resultiert aus der Hypothese, dass unter genug Variabilität die echte Welt lediglich eine bereits trainierte Variation darstellt.²¹⁶ Dieses Prinzip kann auch der Abbildung 6 entnommen werden, in welcher das physische System im Bereich der Randomisierung der Simulation liegt.



Abb. 6: Intuition hinter dem Paradigma von DR²¹⁷

2.8.2 Anwendung von Randomisierung der Simulationsumgebung

Die Anwendung von DR beginnt nach der Erstellung einer möglichst realen Simulation zum Sicherstellen, dass spätere Variation die wirkliche Welt abbilden können.²¹⁸ Anschließend kann die Randomisierung verschiedener Aspekte den Agenten darin unterstützen, die Strategie, welche z. B. als rekurrentes neuronales Netz dargestellt werden kann, so zu optimieren, dass diese für die Wirklichkeit gut genug generalisiert.²¹⁹ Aspekte, welche im Rahmen von dynamischer DR randomisiert werden können, sind z. B. die Masse von Objekten, die Dauer eines Zeitschrittes oder Schubkoeffizienten von Rotatoren.²²⁰ Verwendet man visuelle Randomisierung, so sind bspw. Position, Textur, Ausrichtung und Lichteinwirkung jedes Objektes sowie Kamerawinkel und Zufallsrauschen Eigenschaften, die während des Trainingsprozesses verändert werden können.²²¹ Alghonaim/Johns 2021 untersuchten dazu detailliert den Einfluss der Randomisierung von Hintergrundfarben und -texturen sowie Ablenkungsobjekte auf die Modellleistung und erhielten folgende Ergebnisse.

Aus der Untersuchung in Abbildung 7 wird deutlich, dass die Fehleranfälligkeit der Erkennung der X-Y Position, der Z-Position und der Orientierung, durch das Hinzufügen von Texturen und Ablenkungsobjekten mitunter am besten verringert werden kann.

²¹⁴Vgl. Hsu u. a. 2023, S. 1

²¹⁵Vgl. Chen u. a. 2021, S. 2

²¹⁶Vgl. Tobin u. a. 2017, S. 1

²¹⁷Enthalten in: Zhao/Queralta/Westerlund 2020, S. 6

²¹⁸Vgl. Chen u. a. 2021, S. 4

²¹⁹Vgl. Chen u. a. 2021, S. 4

²²⁰Vgl. Molchanov u. a. 2019, S. 4

²²¹Vgl. Tobin u. a. 2017, S. 3

²²²Enthalten in: Alghonaim/Johns 2021, S. 6



Abb. 7: Einfluss der Randomisierung verschiedener Effekte auf die Fehlerrate der Objekterkennung²²²

2.8.3 Errungenschaften unter Einsatz von Domain Randomization

Mit dem Einsatz von DR ist es bereits nach aktuellem Stand der Forschung und Praxis ermöglicht worden, z. B. Strategien ohne reale Bilder in ausschließlichen Simulationen mittelmäßiger Qualität zu trainieren und für reale Drohnenflüge im Innenbereich einzusetzen.²²³ Weiterhin konnte erreicht werden, ausschließlich aus dem Training in Simulationen eine Objekterkennung und Lokalisation zu einer Genauigkeit von 1,5 cm zu realisieren.²²⁴ Neben der Anwendung von DR wird in der aktuellen Literatur sehr stark thematisiert, welche Randomisierung positive Effekte erzielen kann, jedoch wenig zu deren Erklärungen und Funktionsweise detailliert betrachtet.²²⁵ Dies erschwert die Entwicklung effizienter Simulationen und die Bestimmung der Zufallsverteilungen zur Randomisierung.²²⁶

²²³Vgl. Sadeghi/Levine 2016, S. 1

²²⁴Vgl. Tobin u. a. 2017, S. 1

²²⁵Vgl. Zhao/Queralt/Westerlund 2020, S. 6

²²⁶Vgl. Zhao/Queralt/Westerlund 2020, S. 6

3 Durchführung des Laborexperiments

Anschließend an die Diskussion des aktuellen Stands der Forschung und der Praxis wird im Rahmen dieses Kapitels die Forschungsmethodik angewandt, um die folgende Forschungsfrage zu beantworten:

Inwiefern kann durch den Einsatz eines mittels RL trainierten Gegenspielers die Robustheit der gelernten Strategie verbessert werden?

In der ersten Sektion dieses Kapitels wird auf den Aufbau des Laborexperiments eingegangen. Dabei werden grundlegende Trainingsdaten sowie zu untersuchende Hypothesen, ausgeübte Einflüsse und erfasste Metriken beschrieben. Anschließend werden daraus in der nächsten Sektion resultierende Anforderungen für die Entwicklung der Simulation und der Testumgebung abgeleitet. Die dritte Sektion beinhaltet daraufhin die Umsetzung der beschriebenen Anforderungen und erläutert die endgültige Implementierung der Simulation, der Testdatenerhebung und deren Auswertung.

3.1 Erläuterung der Forschungsmethodik

Innerhalb des Laborexperiments soll die Beziehung zwischen dem Trainingsszenario als Ursache und der Leistungsrobustheit als Wirkung betrachtet werden. Dabei sind zunächst Strategien durch verstärkendes Lernen unter unterschiedlichen Szenarien zu optimieren. Anschließend werden die trainierten Policies in einer veränderten Simulation als Testszenario ausgeübt, und währenddessen verschiedene Metriken der Strategieleistung betrachtet. Ziel ist es, die Leistungsdiskrepanzen zwischen den Strategien des Trainings mit deterministischen und mit optimiertem Gegenspieler zu betrachten. Kein Teil der umzusetzenden Forschungsmethodik ist die Anwendung der simulationsbasierten Strategien zur Steuerung von Quadroptern in der echten Welt. Weiterhin wird nicht die Ergebnisabhängigkeit zu Faktoren wie der Wahl des Simulationsframeworks, der Physik-Engine oder des Abstraktionsniveaus untersucht. Die Auswahl jener Aspekte wird unter den Gesichtspunkten der Softwarearchitektur getroffen, dessen Anforderungen und Implementierung in nachfolgenden Sektionen behandelt wird.

3.1.1 Beschreibung der Simulationsumgebung

Beginnend mit dem Training der Policies durch verstärkendes Lernen, wird hierfür anders als bei überwachtem und unüberwachtem Lernen kein unmittelbar vorliegender Datensatz benötigt. Anstelle dessen basiert das Training auf der vollständigen oder teilweisen Wahrnehmung einer Lernumgebung, welche den lernenden Agenten für ausgeführte Aktionen positiv oder negativ belohnt. Wie auch in der Einleitung erwähnt, wird im Kontext dieser Arbeit die Simulation von Quadroptern dafür verwendet, die Lernumgebung und damit die Trainingsdaten für die

Algorithmen des RL zur Verfügung zu stellen. Die Simulation von Quadrokoptern stellt ein hochdynamisches Anwendungsgebiet dar, bei dem von einer hohen Diskrepanz zwischen der echten und simulierten Welt ausgegangen werden kann. Die Simulationsumgebung stellt grundsätzlich ein Szenario dar, in welchem zwei verschiedene Drohnen kompetitiv gegeneinander agieren. Das Ziel einer Drohne ist es, zu einem festgelegten Punkt hinter der zweiten Drohne zu gelangen, währenddessen die zweite Drohne versucht die angreifende Drohne auf seinem Weg abzufangen. Die verteidigende Drohne verwendet dafür das Mittel einer bewussten Kollision. Beide Drohnen werden unter der Einschränkung ihrer Nähe zum Zielpunkt zufällig in einem drei dimensionalen Raum initialisiert, welcher nur horizontal einseitig beschränkt wird. Durch die einseitige Beschränkung des Flugraums wird der natürliche Untergrund dargestellt. Zu Beginn soll die verteidigende Drohne näher am Zielpunkt als die angreifende Drohne starten, sodass stets ein Abfangen möglich ist. Während der Simulation wird die angreifende Drohne sich entweder regelbasiert auf einer geraden Linie, oder sich durch optimierte Aktionen, zum Zielpunkt bewegen. Die Bewegungsaktionen der verteidigenden Drohne werden durch ein sich optimierendes neuronales Netz bestimmt. Ist das Ziel einer der beiden Drohnen erreicht, oder besteht Kontakt mit dem Untergrund, so wird die Simulation beendet.

3.1.2 Erläuterung der Trainingsszenarien

Durch die erläuterte Simulation werden für das Laborexperiment Strategien in vier unterschiedlichen Szenarien mit einem Algorithmus des verstärkenden Lernens optimiert. Jede der Strategien wird dabei in einer der nachfolgenden Simulationsszenarien optimiert, dessen Auswahl die unabhängige Variable des Laborexperiments darstellt.

1. Das erste Szenario beinhaltet das Training der zu verteidigenden Drohne gegen eine angreifende Drohne, welche eine deterministische regelbasierte Strategie ausführt.
2. Im zweiten Szenario wird eine deterministische Strategie für die anzugreifende Drohne mittels RL optimiert, während die verteidigende Drohne anhand einer regelbasierten Strategie agiert.
3. Anschließend enthält das dritte Szenario das Training einer Verteidigungsstrategie im kompetitiven Spiel mit der anzugreifenden Drohne unter zuvor optimierter Policy.
4. Zum Abschluss wird das erste Szenario mit regelbasiertem Gegenspieler unter DR, also unter der Randomisierung dynamischer Parameter, als viertes Szenario wiederholt. zur Randomisierung der Trainingsdomäne wird ein konstanter Windeffekt simuliert, dessen Richtung und Stärke zu jeder Trainingsepisode variiert.

3.1.3 Erläuterung des Testszenarios

Nach den Trainingsphasen werden einzelne Policies in mehreren Episoden, eines vom Training abweichenden Testszenarios, ausgeführt und deren Leistungsverhalten gemessen. Eine Episode

entspricht dabei einer Simulation des Testszenarios bis zu ihrem erfolgreichen oder nicht erfolgreichem Ende durch Drohnen-, Ziel- oder Bodenkontakt. Zusätzlich zu den Trainingsszenarien ist neben den initialen Startpositionen der Drohnen auch der von der Angreiferdrohne zu erreichende Zielpunkt zufällig zu verändern. Neben der zusätzlichen Randomisierung des Zielpunktes wird außerdem ein zum Trainingsszenario verbesserter regelbasierter Gegenspieler verwendet.

3.1.4 Messung der Robustheit von RL Policies

Die in den Testszenarien betrachteten abhängigen Variablen sind die erzielte kumulierte Belohnung, dessen Varianz sowie die Anzahl an unbeabsichtigten Misserfolgen. Die kumulierte Belohnung und dessen Varianz stellen wie bereits beschrieben in der Forschung wichtige Kenngrößen dar, um die Robustheit von RL Policies zu bestimmen. Als Robustheit wird im Rahmen dieser Arbeit die Signifikanz und Stabilität einer Leistungsdiskrepanz von Modellen zwischen Trainings- und Testszenarien definiert. Die Metrik der Anzahl von unbeabsichtigten Misserfolgen spiegelt im behandelten Anwendungsfall das Fehlschlagen der trainierten Strategie wider die gegnerische Drohne nicht abfangen zu können. Mit der Auswahl der Metriken werden die Strategieeigenschaften der maximalen Strategieleistung, der Leistungsabweichung und der Strategiesicherheit deutlich. Aus der Abweichung zwischen dem Leistungsverhalten während des Trainings und während des Tests soll so die Robustheit von RL Policies erkennbar und messbar gestaltet werden. Der Fokus liegt dabei auf der Messung der Robustheit von den optimierten Strategien aus Szenario eins, drei und vier. Aus dem Vergleich des Leistungsverhaltens zwischen Strategie eins und drei kann eine mögliche Verbesserung der Robustheit gegenüber unbekannten Szenarien abgeleitet werden. Anschließend kann der erzielte Effekt mit der Abweichung zwischen eins und vier verglichen, und so aktuell verwendete Methoden zur Erhöhung der Robustheit einbezogen werden.

3.1.5 Auswertung mittels statistischer Tests

Durch den Vergleich einzelner Strategien mittels dieser Vorgehensweise werden im Laborexperiment die zu Beginn der Arbeit aufgestellten Hypothesen auf ihre Gültigkeit untersucht. Die nachfolgenden bereits zu Beginn angeführten Hypothesen beinhalten, wie im vorherigen Kapitel diskutiert, verschiedene Aspekte der Robustheit von RL Strategien, welche in späteren Abschnitten einzeln ausgewertet werden.

1. *Die im Testszenario erzielte kumulierte Belohnung ist unter Verwendung der Policy aus dem Training mit RL basiertem Gegenspieler signifikant und zuverlässig höher als die Policy aus dem Training mit regelbasiertem Gegenspieler.*
2. *Die Varianz der im Testszenario erzielten Belohnung ist unter Verwendung der Policy aus dem Training mit RL basiertem Gegenspieler signifikant und zuverlässig geringer als die Policy aus dem Training mit regelbasiertem Gegenspieler.*

3. *Die im Testszenario erreichte Anzahl von unbeabsichtigten Misserfolgen ist unter Verwendung der Policy aus dem Training mit RL basiertem Gegenspieler signifikant und zuverlässig geringer als die Policy aus dem Training mit regelbasiertem Gegenspieler.*

Jede der drei oben genannten Behauptungen wird in zwei statistischen Tests ausgewertet. Dabei wird zum einen ein Signifikanztest zur Gleichheit der Verteilungen und zum anderen ein Signifikanztest zur Verschlechterung der Metrik durch RL basierten Gegenspieler durchgeführt. Die Thesen zur gleichen oder schlechteren Verteilung der Metrik durch trainierten Gegenspieler werden in den Signifikanztests als H_0 Hypothese eingesetzt. Die entsprechenden Gegenhypothesen formulieren jeweils die gegensätzliche Annahme einer Ungleichheit der Verteilungen oder einer Verbesserung der Metrik. Fasst man diesen Aufbau der statistischen Signifikanztests zusammen, können folgende Testhypothesen festgelegt werden.

Gleichheit des ersten Messwerts

- **Hypothese H_0 :** *Die erzielte kumulierte Belohnung im dritten Testszenario ist signifikant gleich zu der des ersten Testszenarios.*
- **Hypothese H_1 :** *Die erzielte kumulierte Belohnung im dritten Testszenario ist signifikant unterschiedlich zu der des ersten Testszenarios.*

Verschlechterung des ersten Messwerts

- **Hypothese H_0 :** *Die erzielte kumulierte Belohnung im dritten Testszenario ist signifikant geringer zu der des ersten Testszenarios.*
- **Hypothese H_1 :** *Die erzielte kumulierte Belohnung im dritten Testszenario ist signifikant höher zu der des ersten Testszenarios.*

Gleichheit des zweiten Messwerts

- **Hypothese H_0 :** *Die Varianz der erzielten Belohnung im dritten Testszenario ist signifikant gleich zu der des ersten Testszenarios.*
- **Hypothese H_1 :** *Die Varianz der erzielten Belohnung im dritten Testszenario ist signifikant unterschiedlich zu der des ersten Testszenarios.*

Verschlechterung des zweiten Messwerts

- **Hypothese H_0 :** *Die Varianz der erzielten Belohnung im dritten Testszenario ist signifikant höher zu der des ersten Testszenarios.*
- **Hypothese H_1 :** *Die Varianz der erzielten Belohnung im dritten Testszenario ist signifikant geringer zu der des ersten Testszenarios.*

Gleichheit des dritten Messwerts

- **Hypothese H_0 :** *Die erreichte Anzahl von unbeabsichtigten Misserfolgen ist im dritten Testszenario ist signifikant gleich zu der des ersten Testszenarios.*

- **Hypothese H1:** *Die erreichte Anzahl von unbeabsichtigten Misserfolgen ist im dritten Testszenario ist signifikant unterschiedlich zu der des ersten Testszenarios.*

Verschlechterung des dritten Messwerts

- **Hypothese H0:** *Die erreichte Anzahl von unbeabsichtigten Misserfolgen ist im dritten Testszenario ist signifikant höher zu der des ersten Testszenarios.*
- **Hypothese H1:** *Die erreichte Anzahl von unbeabsichtigten Misserfolgen ist im dritten Testszenario ist signifikant geringer zu der des ersten Testszenarios.*

Die Auswahl des Signifikanztests wird anhand der Wahrscheinlichkeit einer vorliegenden Normalverteilung vorgenommen. Kann nach einem Kolmogorov- oder Shapiro- Test zu einem Signifikanzniveau von 10 % eine Normalverteilung angenommen werden, lässt sich zur Überprüfung der Hypothesen ein T-Test einsetzen. Ist keine Normalverteilung gegeben, wird ein Mann-Whitney U Test verwendet. Liegt schlussendlich der P-Wert des T- oder Mann-Whitney U Signifikanztests zur Prüfung eines Unterschieds unter 10 %, so wird die H0 Hypothese abgelehnt und die Abweichung in den Robustheitsdaten als signifikant betrachtet. Zeigt der anschließende statistische Signifikanztest einen P-Wert unter dem Signifikanzniveau, wird H0 abgelehnt und die Verbesserung mittels RL Gegenspieler als signifikant wahrgenommen. Die Beantwortung der Forschungsfrage erfolgt abschließend mit der Betrachtung der angenommenen oder abgelehnten Forschungshypothesen. So werden die verschiedenen Merkmale der RL Policies untersucht und es kann festgestellt werden, welche Effekte auf Robustheit erzielt worden sind.

3.2 Programmanforderungen

Für die Durchführung des Laborexperiments wird ein Softwareprogramm implementiert, welches in der Lage ist, die zuvor beschriebene Methodik umzusetzen. In diesem Kapitel werden die dafür bestehende technische und funktionale Anforderungen genauer thematisiert.

Eine allgemeine Anforderung besteht in der Wahl einer Entwicklungssprache, welche die Entwicklung der Simulation anhand von Softwareframeworks, und zugleich die Integration mit Algorithmen des verstärkenden Lernens erlaubt. Die Entwicklungssprache und verwendete Softwarebibliotheken sollten einer möglichst aktuellen Version entsprechen, und deren Einsatz mit ihrer Version dokumentiert sein. Durch die Dokumentierung kann ein gleiches Abbild der Softwareumgebung auf anderen Instanzen eines Betriebssystems eingerichtet werden.

3.2.1 Anforderungen der Simulationsumgebung

Die übergeordnete Hauptaufgabe der Simulation ist die Generierung von Trainingsdaten für die Optimierung von Strategien mittels RL. Die Qualität der Simulation und besonders der Grad

des Realismus spielen dabei besonders für die spätere Auswertung der Robustheit einzelner Policies eine wichtige Rolle. In diesem Abschnitt wird die Hauptaufgabe in mehrere Anforderungen unterteilt, welche im Entwicklungsprozess umzusetzen sind.

Aus der Betrachtung des aktuellen Stands der Forschung und der Praxis für die Entwicklung von Simulationen geht hervor, dass besonders das OpenAI Gym Framework bzw. dessen Nachfolger Gymnasium in diesem Kontext ein wichtiger Bestandteil ist. Das Gymnasium Framework definiert die wichtigsten Funktionsschnittstellen, welche zur Entwicklung der Simulation im Rahmen dieser Arbeit zu implementieren sind.

Eine der Funktionen beschreibt die **Initialisierung** der eigenen Umgebungsklasse, welche von der Umgebungsklasse der Gymnasium-Bibliothek erbt. Initial sind darin für die Simulation die Startpositionen der verteidigenden Drohne und der angreifenden Drohne sowie der anzugreifende Zielpunkt anzugeben. Die Simulation soll eine zufällige Platzierung der Drohnen und des Zielpunktes im Raum ermöglichen, jedoch ausschließlich unter der Bedingung, dass die direkte Strecke zum Zielpunkt für die verteidigende Drohne kürzer ist, als für die angreifende Drohne. Zur Erfüllung dieser Anforderung muss eine Wahrscheinlichkeitsfunktion für die Zielpunktkoordinaten sowie je Drohne, für die drei Startkoordinaten X, Y und Z bestimmt werden. Die Randomisierung dieser Koordinaten soll über eine Variable an- und abgeschaltet werden können, sodass einfach zwischen Trainings- und Testszenario gewechselt werden kann. Als Spielgebiet wird eine flache feste Ebene auf der Höhe Null eingesetzt, über dessen der drei dimensionale Raum unbegrenzt in X-, Y- und Z-Richtung zur Verfügung steht.

Eine weitere Funktion beinhaltet die Definition des **Beobachtungsraums**, also welche Informationen für den Agenten wahrnehmbar sind. In der beschriebenen Simulation soll der Agent die Position, die Ausrichtung, die Richtungsgeschwindigkeit der Drohne und die Geschwindigkeit der Rotatoren wahrnehmen. Jede diese Eigenschaften wird von der verteidigenden Drohne und zusätzlich von der angreifenden Drohne wahrgenommen. Insgesamt stelle der Beobachtungsraum somit ein Datenobjekt dar, welches alle Informationen wie Koordinaten, Geschwindigkeiten und Drehzahlen beinhaltet.

Neben dem Beobachtungsraum muss auch ein **Handlungsraum** festgelegt werden, welcher die Steuerung der verteidigenden Drohne umfassen soll. Hierbei soll durch den Agenten die Geschwindigkeit und ihre Richtung in Form einer vier elementigen Liste vorgegeben werden. Auch die regelbasiert gesteuerten Drohnen folgen damit dem Ziel, eine dieser Form entsprechende Aktion zu erzeugen. Anschließend besteht die Anforderung, diese Aktionen in einen konkreten Einfluss auf den Zustand des Quadropters zu übersetzen.

Führt der Agent eine Aktion seines Handlungsraums aus, wird die Implementierung eines **Zeitschritts** der Simulation benötigt. In dieser muss deklariert sein, wie sich die gewählte Aktion auf den Zustand der Drohnen auswirkt. Während jedes Zeitschritts der Simulation sind dabei die Kriterien zum Zurücksetzen der Simulation zu prüfen.

Im Anwendungsbereich dieser Arbeit ist die Simulation zurücksetzen, sobald das Abfangen der angreifenden Drohne erfolgt ist, oder eine der Drohnen in Kontakt mit dem Untergrund kommt. Das Gymnasium Framework definiert die **Rücksetzfunktion**, welche den initialen Simulationszustand herstellt, sobald ein Kriterium erfüllt ist. Durch die Schrittfunktion ist abschließend die neue Wahrnehmung der Simulation, die Erfüllung des Rücksetzkriteriums und der Belohnungswert zurückzugeben.

Zur Berechnung einer entsprechenden Belohnung, für die Ausführung der gewählten Aktion, ist eine **Belohnungsfunktion** einzusetzen. Diese muss anhand des neuen Simulationszustandes Kennzahlen ermitteln und gewichten, und dadurch eine numerische Bewertung der neuen Situation vornehmen. Hinsichtlich der Maximierung des Belohnungswertes optimiert der Agent mittels des RL Algorithmus seine Aktion in Abhängigkeit der Simulation. Die Anforderungen dieser Optimierung der Aktionen werden in der nachfolgenden Sektion detailliert betrachtet.

3.2.2 Anforderungen des Optimierungsverfahrens

Das Optimierungsverfahren trägt die Aufgabe zur Bestimmung der bestmöglichen Aktion zu jedem Zeitschritt der Simulation. Dabei ist die Strategie zur Steuerung des Quadropters mittels RL-Algorithmus zu optimieren. Da die eigene Implementierung des RL-Algorithmus nicht im Fokus dieser Arbeit liegt, werden frei verfügbare Softwarebibliotheken verwendet. Eine Bibliothek, welche eine Vielzahl an RL-Algorithmen beinhaltet, ist Stable-Baselines3 nach Raffin u. a. 2021. Die Softwarebibliothek integriert weitere Bibliotheken wie Tensorboard zum Messen von Trainings- und Evaluationskennzahlen, was die Durchführung der Forschungsmethodik unterstützt. Weiterhin ist eine umfangreiche Dokumentation zur Implementierung von RL Algorithmen in Verbindung mit dem Gymnasium Framework vorhanden, was die Umsetzung des Trainingsprozesses in Verbindung mit der Simulationsumgebung erleichtert. Aus der Softwarebibliothek Stable-Baselines3 soll der PPO-Algorithmus zum Optimieren je einer Policy in jedem Trainingsszenario eingesetzt werden.

Der Trainingsprozess beinhaltet im ersten Schritt die Bestimmung der unabhängigen Variablen. In Kontext unserer Arbeit ist dies die Wahl des Trainingsszenarios, dessen Eigenschaften in der Simulation einzustellen sind. Alle Eigenschaften der Szenarios sollen über die Veränderung weniger Parameter als zentrale Schnittstelle im Programm auswählbar sein. Anschließend sind die gewählten Metriken anzugeben, welche im Trainingsprozess erfasst werden sollen. Diese beinhalten die durchschnittliche Episodenlänge und kumulierte Belohnung sowie die Erfolgsrate des Zusammenstoßes beider Drohnen, über alle im Trainingsprozess evaluierten Strategien. Zur Erfassung dieser Trainingsmetriken gilt es, die Integration mit der Tensorboard Bibliothek zu nutzen. Daraufhin ist das Training über die dafür vorgesehene Funktion des RL-Modells zu starten und nach dessen Ablauf das optimierte Modell zu speichern. Das Speicherformat ist so zu wählen, dass gespeicherte Modelle zur Auswertung in das Testszenario geladen werden können.

Aufgrund des hohen dreidimensionalen Aktionsraums, welcher zu jedem Zeitschritt zu optimieren ist, soll vor dem beschriebenen Trainingsprozess ein Training durch IL eingesetzt werden. Im IL Prozess ist ein Modell dahingehend zu optimieren, eine vorgeschriebene regel-basierte Strategie zur Steuerung des Quadropters nachzuahmen. Die Policies der später mittels RL zu optimierenden Quadropters werden demnach zunächst regelbasiert als Experte definiert. Daraufhin ist es das Ziel, die Gewichte eines IL-Modells wie z. B. Behavioral Cloning so anzupassen, dass diese bestmöglich die regelbasierte Strategie widerspiegeln. Anschließend sind die in diesem Prozess erlernten Gewichte zu einem RL-Modell zu transferieren und mittels RL-Algorithmus zu verfeinern. Durch diesen Transfer soll eine hohe Qualität der RL-Modelle gewährleistet werden, sodass die Messung der Robustheit valide Ergebnisse erzielt.

3.2.3 Anforderungen des Laborexperiments

Anhand der Implementierung des Laborexperiments sollen die notwendigen Messdaten erhoben, die Forschungshypothesen evaluiert, und die Forschungsfrage beantwortet werden. Dazu ist es notwendig, die beschriebenen Metriken, also die kumulierte Belohnung, dessen Varianz und die Anzahl an Misserfolgen, zu allen Policies im Testszenario zu erfassen. Das Testszenario wird durch eine Instanz der Simulationsumgebung repräsentiert, die eine zum Training veränderte Situation beinhaltet.

Die Abweichung der Domäne zu den Trainingsszenarien wird durch die zusätzliche Randomisierung des Zielpunktes und durch die Verbesserung des regelbasierten Gegenspielers erzielt. Die Verteilung des Zielpunktes soll dabei ähnlich zu den Startkoordinaten so gewählt sein, dass sich die verteidigende Drohne näher an ihr befindet. Das Abfangen eines direkten Anflugs der anzugreifenden Drohne bleibt dadurch möglich. Im Testszenario soll die Verbesserung der anzugreifenden Drohne aus einer parabelförmigen anstatt geraden Fluglinie zum Zielpunkt bestehen. Damit die Validität der Messdaten gegeben ist, muss sichergestellt sein, dass alle Strategien den gleichen Testbedingungen und demnach denselben zufälligen Start- und Zielpunkten unterliegen. Hierfür ist ein Zufallswert zu übergeben, worunter bei jedem Programmstart die gleichen zufälligen Testbedingungen hergestellt werden. Anschließend ist jede Policy in mehrfachen verschiedenen Episoden auszuführen, und deren Leistungsverhalten zu messen.

Innerhalb eines Testszenarios sind zu jedem Zeitschritt die aktuelle Belohnung zu speichern und zu messen, ob ein Misserfolg der Drohne vorliegt. Beide Variablen folgen dem Format einer Zeitreihe über den Verlauf der Testsimulation. Das Programm zur Auswertung muss nach 100-facher Ausführung der Policy im Testszenario die Kennzahlen der durchschnittlichen Belohnung und dessen Varianz berechnen. Zur Dokumentation und für die weitere Auswertung mittels statistischer Tests sind die Messdaten zu speichern.

Durch das Laden der Messdaten aus Dateien in entsprechende Datenobjekte liegen alle Metriken zur Bestimmung des statistischen Tests vor. Anschließend sind die Daten durch das Auswertungsprogramm auf ihre Verteilung hin zu untersuchen. Hierfür ist ein Kolmogorov- oder

Shapiro- Test, anhand verfügbarer Softwarebibliotheken zu verwenden. Können die Messdaten als normalverteilt angenommen werden, gilt es ein T-Test, andernfalls einen Mann-Whitney U Test zur Überprüfung der Hypothesen einzusetzen.

Die genaue Implementierung dieser Tests soll ebenfalls durch zusätzliche Softwarebibliotheken bestimmt sein. Anhand der Testergebnisse sind durch das Programm die H_0 und H_1 Hypothesen zu evaluieren und darauf aufbauend die Forschungshypothesen zu bestätigen oder zu verwerfen. Zusätzlich zu den Messdaten der Testszenarien sind auch die Ergebnisse der Signifikanztests zu speichern. Schlussendlich sind in einer Tabelle zu jeder Forschungshypothese die Ergebnisse der statistischen Tests anzuführen. Durch die gesamte einheitliche Betrachtung der Robustheitsmerkmale kann final die Forschungsfrage beantwortet werden.

3.3 Implementierung der Simulation und des Experiments

Fort folgend an die in der letzten Sektion beschriebenen Anforderungen, wird in dieser Sektion die finale Implementierung der einzelnen Programmabschnitte erläutert. In Anlehnung an den Aufbau der letzten Sektion, wird die Umsetzung der Anforderungen zu den Programmabschnitten Simulationsumgebung, Trainingsverfahrens und Laborexperiment dargelegt.

Allgemeiner Natur ist die Wahl der Entwicklungssprache Python. Die Entwicklungssprache ermöglicht eine hohe Kompatibilität mit bekannten Simulationsframeworks, RL-Softwarebibliotheken und Bibliotheken zur statischen Auswertung von Messdaten. Im Rahmen dieser Arbeit wird konkret die Python Version 3.10 eingesetzt, um von den Möglichkeiten neuer Funktionen zu profitieren und zugleich eine möglichst hohe Kompatibilität zu existieren frei verfügbaren Bibliotheken aufzuweisen. Alle in dieser Arbeit verwendeten Bibliotheken sowie deren Abhängigkeiten sind in einer virtuellen Umgebung installiert, sodass bestmöglich Störeinflüsse durch bereits installierte Pakete vermieden werden. Ebenso sind alle eingesetzten Softwarepakete in einer Textdatei inklusive ihrer Version dokumentiert.

3.3.1 Programmumsetzung der Simulationsumgebung

Die Umsetzung der Anforderungen der Simulationsumgebung beinhaltet im Fokus die Entwicklung einer Simulationsumgebung anhand des Gymnasium-Programmiergerüsts. Die entwickelte Simulation baut dabei auf bestehende Programmbibliotheken auf. Dies fördert vor allem die im Zuge dieser Arbeit umsetzbare Qualität und ermöglicht einen hohen Grad an Realismus, welcher wie in den Anforderungen beschrieben unmittelbar das Sim2Real Problem und damit die Robustheit beeinflusst. Weiterhin stellt die Neuentwicklung einer Simulationsumgebung kein Entwicklungsaufwand dar, welcher einen Einfluss auf die auszuwertende Forschungsfrage ausübt. Wird die im zweiten Kapitel angeführte Auswahl an in der Literatur beschriebenen Simulationsumgebungen betrachtet, kann aus diesen eine Basis zur Entwicklung der eigenen Simulationen gewählt werden.

Zur Auswahl einer Basissimulation, in welcher die eigenen Trainings- und Testszenarien abgebildet werden können, sind unterschiedliche Kriterien zu beachten. Ein Kriterium ist die Kompatibilität der Basissimulation mit der gewählten Entwicklungssprache. Hierbei sollte entweder die Simulation selbst in Python, oder eine entsprechende Schnittstelle entwickelt sein. Die gewählte Simulation sollte eine Physik-Engine einsetzen, welche zum einen einen möglichst hohen Grad an Realismus erlaubt, um die Forschungsfrage möglichst valide zu beantworten. Zum Anderen sollte die Simulation einen beherrschbaren Rechenleistungsaufwand verursachen, da die Zielsetzung die Optimierung mehrerer Policies beinhaltet. Weiterhin, zur Erfüllung der zuvor beschriebenen Anforderungen, muss die Basissimulation dem Framework Gym oder dessen Nachfolger Gymnasium entsprechen. Um die Entwicklung des Optimierungsverfahrens zu unterstützen, wird eine RL Integration oder zumindest eine umfangreiche Dokumentation zu dessen Einsatz vorausgesetzt. Die nachfolgende Tabelle stellt die in Betracht gezogenen Drohnensimulationen den für die Auswahl getroffenen Kriterien gegenüber.

	RotorS	CrazyS	AirSim	gym-pybullet-drones
Python API	X	X	X	X
Integration des Gymnasium Frameworks	X	X	X	X
höchster Realismusgrad			X	
kontrollierbarer Rechenleistungsaufwand	X	X		X
umfangreiche RL Integration und Dokumentation			X	X
Simulation mehrerer Quadrokopter enthalten				X

Tab. 5: Gegenüberstellung von Auswahlkriterien und bekannten Drohnensimulationen

Aus der Gegenüberstellung von Auswahlkriterien und aktuellen Simulation geht hervor, dass zunächst alle Simulationen eine Schnittstelle für die Entwicklungssprache bereitstellen. Dabei wurden native Schnittstellen sowie Schnittstellen aus zusätzlichen Softwarebibliotheken miteinbezogen. Auch die Integration des Gym oder Gymnasium Frameworks wird durch alle Simulationen eigens oder durch zusätzliche Bibliotheken sichergestellt. Bei der Betrachtung der Realitätsnähe zeigt die AirSim Umgebung aufgrund der verwendeten Unreal Physik-Engine den höchsten Grad an Realismus auf. Im Gegenzug werden durch diese Simulation Rechenkapazitäten erwartet, welche im Rahmen dieser Arbeit nicht zur Verfügung stehen. Eine umfangreiche Dokumentation und bereits vorhandene Integration von RL wird durch AirSim und gym-pybullet-drones gewährleistet. Die Simulationen RotorS und CrazyS bieten eine Integration mit RL nur auf Basis der wenig dokumentierten zusätzlichen Softwarebibliotheken GymFC und gym_multirotor. Wird das letzte Kriterium der Simulation mehrerer Drohnen betrachtet, tritt dies nur in Beispielen der gym-pybullet-drones Simulation auf. Fasst man die Erfüllung der Auswahlkriterien

²²⁶Eigene Darstellung

zusammen wird erkenntlich, dass gym-pybullet-drones als einzige Drohnensimulationen die zusätzlichen Anforderungen der RL Integration und der Dokumentation erfüllt. Demnach wird für die nachfolgende Entwicklung der eigenen Simulation auf der gym-pybullet-drones Bibliothek aufgebaut.

Die gym-pybullet-drones Simulation ist in einzelne Simulationszenarien gegliedert, welche je eine Python Datei umfassen. Als Aviaries gekennzeichnet, bilden sie je eine Trainingsumgebung nach der Schnittstellendefinition des Gymnasium-Frameworks ab. In der Struktur der Simulationsumgebung wird das Programmierkonzept der Vererbung verwendet, um ähnliche Simulationen auf denselben übergeordneten Klassen zu basieren. Jede erbende Klasse enthält damit alle Funktionen und Variablen der übergeordneten Elternklasse. Dies ermöglicht die verschiedenen starke Abstraktion der Simulation auf den unterschiedlichen Vererbungsebenen. Weiterhin ermöglicht die Bibliothek verschiedene Steuerungsarten der Quadrokopter. Eine Option ist die Steuerung über die direkte Vorgabe der Drehzahlen aller Rotatoren. Zusätzlich kann auch ein PID-Kontroller eingesetzt werden um die Steuerungssignale entgegenzunehmen. Zur Entwicklung der eigenen Simulation wird die Steuerung über einen Richtungsvektor und einen Geschwindigkeitswert gewählt. Unabhängig des gewählten Aktionstyps werden alle Steuerungssignale durch Kontrollerklassen in konkrete Drehzahlen übersetzt. Dessen genaue Übersetzung wird im Rahmen dieses Kapitels nicht erläutert und ist in der entsprechenden Dokumentation nach Panerati u. a. 2021 nachzulesen. Auch der Beobachtungsraum kann zum einen visuell und zum anderen kinetisch erfolgen, was im Rahmen dieser Arbeit verwendet wird. Der kinetische Beobachtungsraum stellt eine Reihe von wahrnehmbaren Eigenschaften wie Position, Ausrichtung, Geschwindigkeit und Drehzahl dar. Unabhängig der Überwachungsart kann durch die pybullet Physik-Engine ein grafisches Modell der Simulation erzeugt werden, welches in Abbildung 8 dargestellt wird.



Abb. 8: grafisches Modell der Simulation²²⁷

Die Umsetzung, der zu Beginn dieses Kapitels beschriebenen Simulationsumgebung, basiert auf der Verwendung der Basissimulation und der Kontrollerklassen. Zusätzlich sind teilweise Funktionen aus ähnlichen, bereits nativ vorhandenen Szenarien mit Änderungen übernommen und im Programm entsprechend gekennzeichnet. Die beschriebene Simulation ist innerhalb einer Klasse nach dem Gymnasium Framework aufgebaut. Zusätzlich wird die Simulation durch eine er-

²²⁷Eigene Darstellung

bende Klasse umhüllt, um für den Algorithmus eine verarbeitbare Abstraktionsschicht zu erzeugen. In dieser Schicht wird die Steuerung der angreifenden Drohne abstrahiert, sodass die Trainingsumgebung lediglich die Bestimmung einer Aktion für die verteidigende Drohne erwartet. Nachfolgend wird die Umsetzung dieser beiden Klassen nach den einzelnen Funktionen des Gymnasium-Frameworks erläutert.

Die Simulation beider Quadrokopter ist in der Klasse **DualDroneAviary** definiert. Zur Initialisierung einer Simulationsinstanz wird die **Init-Funktion** aufgerufen, welche die Übergabe aller simulations relevanten Parameter erwartet. Die Parameter umfassen unter anderem Informationen zur Anzahl und Art und Position der Drohnen, gewählte Handlungs- und Beobachtungsarten sowie die unabhängigen Variablen der Trainingsszenarien. Außerdem wird die Gewichtung der Belohnungsfunktion als Parameter entgegengenommen. Kernaufgabe dieser Funktion ist die Umsetzung der Eigenschaften eines Trainings- und Testszenarios. Dies beinhaltet die Erzeugung und Übermittlung der zufälligen Start- und Zielpunkte und des zufälligen Windeffektes mittels unterstützender Funktionen aus der eigenen Hilfsklasse. Die verteidigende Drohne wird zufällig innerhalb von 1,5 Meter in X- und Y-Richtung und 0.1 bis 3 Meter in Z-Richtung zum Ursprung initialisiert. Hingegen die angreifende Drohne wird im Rahmen von zwei bis sechs Metern uniform um den Ursprung herum auf einer möglichen Höhe von 0.1 bis fünf Meter platziert. Sind die Eigenschaften bestimmt, können diese zur Erzeugung der darunter liegenden Abstraktionsschicht der BaseAviary Instanz genutzt werden.

Der **Beobachtungsraum** in dieser Simulation wird durch einen Informationsvektor definiert. Der für die Experimente eingesetzte Informationsvektor je Drohne enthält die Position, Ausrichtung, Fortbewegungs- und Drehgeschwindigkeit sowie die Rotationsdrehzahlen. Die Position wird durch dreidimensionale Koordinaten, die Ausrichtung durch Quaternion und Drehwinkel kodiert. Positionen können dabei von $-\infty$ bis $+\infty$ in X- und Y-Richtung und von 0 bis $+\infty$ in Z-Richtung betrachtet werden. Die Quaternion erlauben Werte zwischen -1 und $+1$, die Drehwinkel von $-\pi$ bis $+\pi$. Drohnen-, Dreh- und Rotatoren Geschwindigkeit werden als Geschwindigkeitsvektor von minus bis plus unendlich zu den Koordinatenachsen bzw. als vier einzelne Werte von null bis zu maximalen Drehzahl definiert. Insgesamt ist für jede Drohne ein Beobachtungsvektor mit 20 Elementen innerhalb eines Wörterbuchobjektes definiert.

Die Definition des **Aktionsraums** ist in der ActionSpace Funktion bestimmt. Je nach Aktionstyp ist ein drei oder vierdimensionaler Vektorbereich definiert. Jeder Wert des Vektors kann besitzt eine untere und obere Grenze, aus dessen Bereich die endgültige Aktion gewählt wird. Zur Steuerung des Quadrokopters wird im Rahmen dieser Arbeit der Aktionstyp auf Basis der Geschwindigkeit verwendet. Der Wertebereich des Richtungsvektors ist dabei von $[-1, -1, -1]$ bis $[1, 1, 1]$ definiert. Das vierte Element spiegelt den Anteil der maximalen Geschwindigkeit von null bis eins wider.

Ist eine Aktion regelbasiert oder durch ein RL-Modell bestimmt, wird diese als Parameter an die **Step-Funktion** weitergegeben. Auf der Ebene der DualDroneAviary Klasse kann die übergebene Aktion in Abhängigkeit vom anfangs initialisierten Wind abgewandelt werden. Der Windvektor

wird in Abhängigkeit des Winkels zwischen Bewegungsrichtung der Drohne und des Windes zur Aktion hinzugefügt. Der Drohnenbewegungsrichtung wird die Windrichtung addiert und die Geschwindigkeit in Abhängigkeit des Eintreffwinkels vollständig positiv, negativ oder gar nicht beeinflusst. Anschließend wird unter veränderter oder unveränderter Aktion die Schrittfunktion der Elternklasse aufgerufen. Die Schrittfunktion der Elternklasse übersetzt die Aktion aus dem Aktionsraum in konkrete Rotordrehzahlen und übt diese unter der simulierten Physik aus. Zur Erhöhung des Realismusgrades wird zusätzlich zu den Effekten der pybullet Physik-Engine ein Boden- und Trägheitseffekt eingesetzt. Anschließend wird die neue Zustandsbeobachtung, die entstandene Belohnung, das Rücksetzkriterium und ein Datenobjekt für Zusatzinformationen zurückgegeben.

Die Bestimmung der Belohnung ist unter der **reward-Funktion** der DualDroneAviary Klasse implementiert. Innerhalb dieser Funktion werden aus dem aktuellen Zustand der Simulation Kennzahlen zu dessen Evaluation berechnet und gewichtet. Durch die Gewichtung und dessen Optimierung kann das Strategieziel formuliert und der Einfluss des Wertebereichs der Kennzahlen angepasst werden. Die Aufsummierung der Produkte aus den folgenden Kennzahlen und deren Gewichtung spiegelt die Belohnung zu jedem Zeitschritt wider.

- Distanz zwischen der angreifenden und verteidigenden Drohne
- Belohnung bei erfolgreichem Abfangen der angreifenden Drohne
- Negative Belohnung, sofern der Zielpunkt erreicht wurde
- Belohnung für eine möglichst direkte Flugrichtung zur gegnerischen Drohne
- Lineare Belohnung der Geschwindigkeit

Eine Optimierung der Gewichte wurde durch Ausprobieren unterschiedlicher Gewichtungskombinationen und Messen der entsprechenden Erfolgsrate erzielt. Zu jedem Zeitschritt eintretende Kennzahlen wurden aufgrund ihrer Häufigkeit in Fünfer Schritten, außerordentliche Belohnungen in 10000er-Schritten untersucht. Insgesamt sind daraus unter anderem folgende in Abbildung 9 dargestellte Log-Dateien entstanden, welche den Optimierungsprozess widerspiegeln.



Abb. 9: Episodenlänge, Belohnung und Erfolgsrate unter verschiedenen Gewichtungen²²⁸

²²⁸Diagramme durch die Tensorboard Bibliothek erstellt

Die Optimierung der Belohnungsfunktion brachte $[-10, 50000, 50000, 25, 20]$ als optimalen Gewichtsvektor hervor, da hierunter mehrere Algorithmen wie PPO und A2C, eine positive Erfolgsrate verzeichneten.

Ist das Rücksetzkriterium erfüllt, wird die Simulationsumgebung mittels **reset-Funktion** neu gestartet. Dabei sind neue Zufallspositionen der Drohnen und des Ziels sowie im vierten Trainingsszenario neue Windeigenschaften zu generieren. Sind diese Variablen verändert, kann die Rücksetzfunktion der darunterliegenden Abstraktionsschicht aufgerufen werden. Die Funktion der BaseAviary-Klasse setzt darin die Zeitschritte und die Drohneninformationen zurück und baut das grafische Modell neu auf. Über den Zugriff auf den Klienten der pybullet Physik-Engine werden die Drohnenobjekte sowie das Untergrundobjekt aus den URDF-Dateien geladen und an entsprechenden Stellen platziert. Zur Unterscheidung des verteidigenden Quadropters von der angreifenden Drohne werden zwei nahezu identische URDF-Dateien verwendet, welche sich nur in der Farbgebung zwischen rot und blau unterscheiden. Am Ende der Funktion wird der erste Beobachtungsvektor der neuen Episode sowie ein Infoobjekt zurückgegeben.

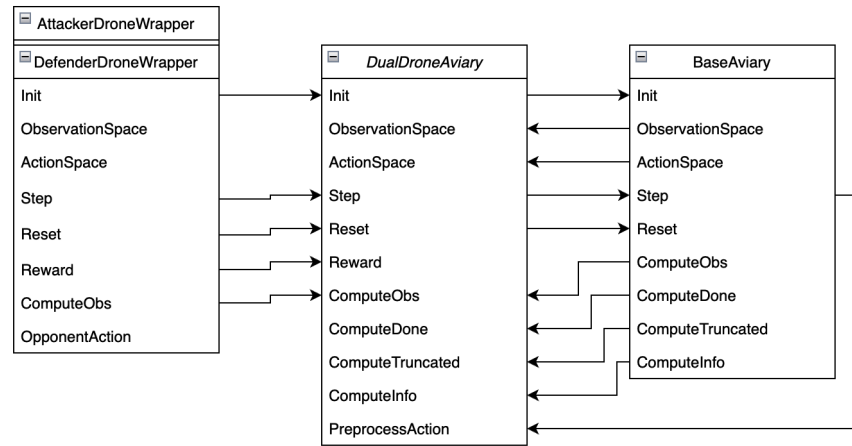
Die zuvor beschriebene Klasse vererbt ihre Variablen und Funktionen auf die höhergelegenen Abstraktionsschichten der **DefenderDroneWrapper** und **AttackerDroneWrapper** Klasse. Auf dieser Ebene bietet die Simulation ihre Schnittstelle für die Verwendung durch den IL- und RL-Algorithmus. Dafür ist es notwendig den Beobachtungsraum sowie den Aktionsraum zu neu zu definieren. Zur **Initialisierung** werden zum einen alle Parameter zur Instanziierung des DualDroneAviary Objekt erwartet. Zum Anderen sind zusätzlich die Strategie und dessen Rauschverhalten der nicht durch RL oder IL gesteuerten Drohne anzugeben, sodass diese in Klassen variablen gespeichert werden. Auf der Ebene der Wrapper-Klassen wird der **Beobachtungsraum** mit allen Informationen beider Quadropters zu einem 40-elementigen Vektor zusammengefasst. Der **Handlungsraum** hingegen umfasst auf dieser Ebene lediglich einen elementigen Vektor mit dem bekannten Wertebereich. Zur Ausführung eines Simulationsschrittes werden in der **Step-Funktion** die Aktion der jeweils gegnerischen Drohne und die durch das Modell übergebene Aktion in ein gemeinsames Datenobjekt überführt. Anschließend kann mittels dieses Datenobjektes die Schrittfunktion der Elternklasse aufgerufen, und deren Rückgabewerte an den zu optimierenden Algorithmus weitergereicht werden. Als **Belohnungswert** wird dabei lediglich auf die im Trainingsszenario relevante Komponente des Belohnungsobjektes zugegriffen.

Fasst man die beschriebene Umsetzung der Simulationsumgebung zusammen, so kann folgendes Klassendiagramm und deren Funktionsabhängigkeiten in Abbildung 10 betrachtet werden.

3.3.2 Programmumsetzung des Optimierungsverfahrens

Mit der Entwicklung der Simulationsumgebung stehen Trainingsdaten zur Verfügung, um darauf aufbauend Strategien des verstärkenden Lernens zu optimieren. Der gesamte Optimierungspro-

²²⁹Eigene Darstellung

Abb. 10: Klassendiagramm mit Funktionsabhängigkeiten²²⁹

zess unterteilt sich in den Einsatz von imitierendem Lernen zu Beginn und folgender Optimierung mittels RL-Algorithmen. Die Umsetzung des Optimierungsverfahrens wird innerhalb der Python Dateien *imitation_train.py* und *train.py* beschrieben.

Das imitierende Training beginnt mit dem Herstellen des Trainingsszenarios zum Akquirieren der Demonstrationen des Experten. Innerhalb einer Trainingsiteration wird die Simulation unter den Trainingsgegebenheiten in 20000 Episoden ausgeführt. Beim imitierenden Lernen wird in der Trainingsumgebung eine gewählte Aktion für 50 Zeitschritte ausgeführt, um die Trainingsdaten aus einer Episode zu minimieren. Damit verteilen sich die daraus entstehenden Transition-Objekt der verwendeten Softwarebibliothek Imitation, auf eine höhere Anzahl von verschiedenen Start- und Zielkoordinaten. Dieses Datenobjekt beschreibt einen durch die Aktion entstehenden Übergang vom bekannten zu einem neuen beobachteten Zustand sowie den Erhalt der Belohnung. Hieraus ergeben sich eine Menge von aufeinanderfolgenden oder zeitversetzten Transition-Objekten. Die Übergänge in den Trainingsdaten liegen häufig nicht entsprechend ihrer Entstehung geordnet, um die Datenvielfalt zu erhöhen und korrelierende Effekte zu mindern. Basierend auf diesen Transaktionsdaten wird in 4000 Epochen die Parameter der zu optimierenden Strategie angepasst. Die Anpassung der Gewichte einer Akteur-Kritiker-Strategie wird durch den Behavioral Cloning Algorithmus vorgenommen. Insgesamt wird dieser abwechselnde Vorgang der Erhebung von Daten und die zur Basis dessen stattfindende Optimierung so häufig wiederholt, bis alle Iterationen die Menge an Expertendemonstrationen chargenweise abgearbeitet haben.

Als Experte wird eine regelbasierte Version der zu optimierenden Drohne eingesetzt. Alle regelbasierten Strategien sind in der *rule-based.py* Datei deklariert und umfassen die verteidigende Drohne und die Trainings- und Testversion der anzugreifenden Drohne. Die verteidigende Drohne verfolgt regelbasiert die direkte Flugbahn zu einem Vorhaltepunkt in Bewegungsrichtung der Angreiferdrohne. Der Vorhaltepunkt rückt mit geringer werdender Distanz zur generischen Drohne immer näher zu dessen aktueller Position. Die angreifende Drohne wird im ersten Trainings-szenario regelbasiert gesteuert, indem sie sich von ihrem Startpunkt in gerader Fluglinie zum Zielpunkt begibt. Dessen Erweiterung zur parabelförmigen Fluglinie wird in der nachfolgenden

Sektion der Umsetzung des Laborexperiments behandelt.

Wurde ein Akteur-Kritiker Modell durch imitierendes Lernen hin zum regelbasiertem Verhalten optimiert, wird dieses als Basis zur weiteren Optimierung mittels RL verwendet. Im RL-Prozess ist zu beachten, dass anders als beim imitierenden Lernen eine gewählte Aktion lediglich einen Zeitschritt ausgeführt wird, da hierbei keine Trainingsdaten gespeichert werden müssen. Mit dem Einsatz des *train.py* Programms lässt sich das in einer Pickle-Datei gespeicherte Modell laden. Dazu wird das grundlegende neuronale Netz des PPO-Objekts mit derselben Struktur der Modells aus dem imitierenden Lernen initialisiert. Anschließend werden die Gewichte und Vorurteile in ein Objekt des PPO-Algorithmus von Stable-Baselines3 übertragen. Mittels der *learn*-Funktion kann die erzeugte Strategie im gegebenen Trainingsszenario über eine Anzahl von Zeitschritten weiter optimiert werden. Die während des RL Prozesses verwendeten Trainingszeitschritte lagen bei 2,592,000, welche sich aus der Trainingszeit von 180 Minuten und der Simulationsfrequenz von 240 Schritten pro Sekunde zusammensetzen. Durch die Tensorboard Bibliothek werden während des Trainingsprozesses gemittelte Metriken über alle Episoden berechnet sowie alle 10000 Zeitschritte eine Evaluation durchgeführt. In der Evaluationsperiode wird die aktuelle Strategie bemessen und dessen Episodenlänge, erzielte Belohnung und Erfolg aufgezeichnet. Als erfolgreich wird die Simulation gekennzeichnet, sobald die Strategie der Verteidigerdrohne die angreifende Drohne durch Kollision abfangen konnte. Nach dem Abschluss der Optimierung werden die neuen Modelle in Pickle Dateien gespeichert.

3.3.3 Programmumsetzung des Laborexperiments

Als letzter Teil der Umsetzung wird die Durchführung des Experiments betrachtet. Kernaufgabe dieses Teils ist die Erhebung und Auswertung der Messdaten zu dem in der Einleitung des Kapitels verfassten Experimentaufbau. Die Erhebung der Testdaten erfolgt durch das *test.py*, die Auswertung durch das *experiment_evaluation.py* Skript.

Mit der Ausführung des ersten Skripts wird ein Testklasse instanziiert und darin enthalten die Simulation nach den Eigenschaften des Testszenarios erzeugt. Das Testszenario implementiert dabei die Anforderung der Randomisierung des Zielpunktes. Zur zufälligen Bestimmung des Zielpunktes wird ebenso wie zur Initialisierung der Verteidigerdrohne, eine uniforme Verteilung bis zu 1,5 Meter um den Ursprung verwendet. Eine uniforme Verteilung fördert dabei die Varianz der Testszenarien und erhöht so die Validität des Experiments. Anders als bei der Initialisierung der Drohnen wird die Höhe des Zielpunktes stets auf 0,5 Meter festgelegt, um so Bodenkontakt und Simulationserfolg abzugrenzen. Zusätzlich zur Randomisierung des Zielpunktes wird im Testszenario der regelbasierte angreifende Quadrokopter im Aspekt seiner parabelförmigen Anflugsstrategie verändert. Ausgehend von den Startkoordinaten des Angreifers werden dazu insgesamt fünf Wegpunkte entlang einer unten geöffneten Parabel berechnet, welche nacheinander angeflogen werden. Die Parameter der Parabel berechnen, sich wie in den Formeln 8 bis 10 dargestellt.

Episode	Episodenlänge	kumulierte Belohnung	Belohnungsabweichung zum Mittelwert	Episodenmisserfolg
1	1931	13512	3192	0
2	1579	8413	1873	1
3	1431	16920	6308	0
4	1944	6741	3701	1
5	1508	10173	89	0

Tab. 6: Beispieltabelle der erhobenen aggregierten Messdaten

$$a = -1/(2 * (\text{Distanz zum Zielpunkt}/6)) \quad (8)$$

$$b = 0 \quad (9)$$

$$c = \text{Z-Koordinate des Startpunktes} \quad (10)$$

Alle Messdaten werden aus der hundertfachen Durchführung der Testsimulation erhoben. Zu jedem Zeitschritt jeder Episode wird durch das RL-Modell eine bestimmte Aktion ausgeübt und eine Teilmenge dessen Resultats gemessen. Die Teilmenge umfasst die durch die Aktion unmittelbar ausgelöste Belohnung und einen Wahrheitswert, ob mit Aktion ein Misserfolg der Simulation eingegangen wurde. Beide Werte werden mit der aktuellen Episode und dem Zeitschritt einer Tabelle als Datensatz hinzugefügt. Nach dem erfolgreichen Durchlaufen aller Zeitschritte aller der 100 Episoden wird die Tabelle als Komma separierte Wertedatei (CSV) gespeichert.

Die erstellte Datei wird anknüpfend mittels *experiment_evaluation.py* ausgewertet, indem die Messdaten geladen, die Signifikanztests durchgeführt und daran die Hypothesen bestimmt werden. Für jeden dieser drei Schritte wurde eine Funktion entwickelt, welche im Skript aufgerufen wird. Zum Laden der Messdaten sind die Speicherpfade zweier Wertedateien anzugeben, woraufhin der Verlauf der Kennzahlen zu einzelnen Episoden aggregiert werden. Dabei wird zunächst die Episodenlänge, die kumulierte Belohnung und die Anzahl an Misserfolgen ermittelt. Zu dieser Basis kann die Abweichung der Belohnung je Episode zu dem Mittel aller kumulierten Belohnungen im Testprozess berechnet werden. Nach dem Laden und Vorverarbeiten entspricht der Aufbau der Messdaten zweier Tabellen ähnlich der Tabelle Sechs.

Im nächsten Schritt ist mittels statistischer Signifikanztests die Ähnlichkeit zu einer Normalverteilung sowie die ungleiche Verteilung und mögliche Verbesserung aller Metriken zu untersuchen. Die Implementierung der Signifikanztests wird im Rahmen dieser Arbeit durch die Softwarebibliothek SciPy bereitgestellt. Iterativ werden für alle Metriken beider Messtabellen die Ähnlichkeit zur Normalverteilung überprüft, um zwischen einem T-Test und einem Mann-Whitney U Test auszuwählen. Anschließend wird je Metrik der entsprechende Test zur Ungleichheit und Verbesserung der Verteilung ausgeführt, woraus je Test eine Teststatistik und ein P-Wert gespeichert wird. Liegt der P-Wert eines Tests unter dem Signifikanzniveau von 10 % wird die H_0 Hypothese abgelehnt und daraus schlussfolgernd die Ungleichheit oder Verbesserung der Metrik angenommen. Insgesamt wird diese Auswertung der Leistungsdaten zweier Strategien zweimal

²²⁹Eigene Darstellung

durchgeführt. Dabei werden die erzielten Messdaten zwischen der Strategien aus Trainingsszenario eins und drei sowie aus Szenario eins und vier betrachtet. Der vollständige Prozess des Experiments kann auch als BPMN Modell wie in Abbildung 11 dargestellt werden.



Abb. 11: BPMN Prozess des Laborexperiments²³⁰

²³⁰Eigene Darstellung

4 Ergebnisse des Laborexperiments

Anschließend an die Durchführung des Laborexperiments werden innerhalb dieses Kapitels die erzeugten Messdaten betrachtet und die zur Beantwortung der Forschungsfrage aufgestellten Hypothesen ausgewertet. Es gilt dabei einen möglichen Effekt des trainierten Gegenspielers auf die Robustheit des RL-Modells zu untersuchen. Dazu werden die Leistungen der Strategien aus den unterschiedlichen Trainingsszenarien im abweichenden Testszenario verglichen. Eine Leistungsmessung im abweichenden Testszenario zeigt, welche Lernumgebung eine stabilere und robustere Strategie hervorbringt, welche auch unter unbekannten Effekten performanter ist.

Die auszuwertenden Leistungsmetriken aggregiert zu ihrer Testepisode können im Anhang gefunden werden. Im ersten Vergleich werden die Strategie aus dem Training mit regelbasiertem und mittels RL optimierten Gegenspieler untersucht. Weiterhin wird ein Vergleich der Testmetriken zwischen dem Training ohne und mit randomisierter Umgebung sowie regelbasiertem Gegenspieler durchgeführt. Dazu werden mehrfach die selben statistischen Signifikanztests durchgeführt, um zum einen den Effekt der DR zu ermitteln. Zum anderen werden die Einflüsse der im Rahmen dieser Arbeit entwickelten Methoden wie der DR und dem RL basierten Gegenspieler gegenübergestellt. Das Ende dieses Kapitels bildet die Beantwortung der Forschungsfrage anhand der ausgewerteten Test- und Forschungshypothesen.

4.1 Vergleich der Robustheit der Strategien aus dem Training mit regelbasierten und RL Gegenspieler

Metrik	Gleichheit P-Wert	Verschlechterung P-Wert	Ungleichheit	Verbesserung
sum reward	0.6276673409642575	0.6870324210336459	False	False
reward mean difference	1.2200477617083057e-22	1.0	True	False
count of failures	0.1814596593937995	0.09072982969689974	False	True

Tab. 7: Ergebnisse der statistischen Tests aus dem Leistungsvergleich der Modelle aus den Trainingsszenarien 1 und 3

Anhand Tabelle 7 lassen sich die Ergebnisse der statistischen Signifikanztests ablesen. Die Ergebnisse basieren dabei auf dem Leistungsvergleich der Modelle aus den Trainingsszenarien 1 und 3. Den Tabellenspalten werden zum einen die P-Werte der Tests zur Gleichheit und Verschlechterung der Verteilung zugeordnet. Zum anderen wird darauf basierend die Erfüllung der H1-Hypothese anhand des gewählten Signifikanzniveaus von 10% evaluiert. Einer Zeile der Tabelle können die Ergebnissdaten zu einer der drei untersuchten Metriken entnommen werden.

4.1.1 Betrachtung der kumulierten Belohnung

Hypothese 1: *Die im Testszenario erzielte kumulierte Belohnung ist unter Verwendung der Policy aus dem Training mit RL-basiertem Gegenspieler signifikant und zuverlässig höher als die Policy aus dem Training mit regelbasiertem Gegenspieler.*

Die erste in der Einleitung aufgestellte Forschungshypothese untersucht den Trainingseinfluss auf die in einer Testepisode erzielte kumulierte Belohnung. Mit der Hypothese 1 wird die Behauptung aufgestellt, dass durch die Wahl eines Trainingsszenarios mit RL-optimierten Gegenspieler zuverlässig eine bessere kumulierte Belohnung erzielt werden kann, als mittels eines Trainings mit regelbasierten Gegenspieler. Die Auswertung der Hypothese erfolgt durch statistische Signifikanztests zur Verteilung der kumulierten Belohnungen je Testepisode. Zur Bestätigung dieser Forschungshypothese müssen die P-Werte zur Gleichheit und Verbesserung der Belohnungsverteilung je geringer als das Signifikanzniveau sein.

Der P-Wert des Signifikanztests zur Gleichheit der Verteilungen der kumulierten Belohnungen liegt in etwa bei 0.6277. Daraus kann festgestellt werden, dass dieser größer als das gewählte Signifikanzniveau von 10% bzw. 0.1 ist. Die im Test untersuchte H_0 Hypothese einer gleichen unterliegenden Verteilung ist damit anzunehmen. Eine Ungleichheit der unterliegenden Verteilung der Belohnungen, wie in Hypothese H_1 behauptet, ist folglich abzulehnen. Da demnach keine Unterschiede in der Verteilung der im Testszenario erzielten Belohnungen vorliegen ist festzuhalten, dass die Forschungshypothese 1 abgelehnt wird. Im Laborexperiment konnte somit kein Nachweis erbracht werden, dass ein Optimieren eines RL-Modells, unter Verwendung eines selbst mittels RL optimierten Gegenspielers, zuverlässig zu einer höheren kumulierten Belohnung führt.

Schließt man von diesem Ergebnis auf dessen Bedeutung für die Robustheit der Strategien aus verstärkendem Lernen zeigt sich folgende Erkenntnis. Betrachtet man lediglich den Aspekt der kumulierten Belohnung einer Episode, kann eine Lernumgebung mit RL-basiertem Gegenspieler keinen positiveren Effekt auf die Leistungsfähigkeit unter unbekannten Umgebungsverhalten ausüben. Ebenso ist auch festzuhalten, dass keine Verschlechterung der Leistungsfähigkeit im Testszenario vorliegt. Ein Einfluss eines RL-basierten Gegenspielers auf die Robustheit der kumulierten Belohnung unter unbekannten Umgebungsverhalten kann demnach als nicht gegeben betrachtet werden.

4.1.2 Betrachtung der Belohnungsstabilität

Hypothese 2: *Die Varianz der im Testszenario erzielten kumulierten Belohnung ist unter Verwendung der Policy aus dem Training mit RL-basiertem Gegenspieler signifikant und zuverlässig geringer als die Policy aus dem Training mit regelbasiertem Gegenspieler.*

Mittels der Betrachtung der zweiten Forschungshypothese wird die Stabilität der im Testszenario erzielten kumulierten Belohnung untersucht. Dazu wird die Abweichung einer Belohnung aus

einer Testepisode zum Mittelwert der Belohnungen über aller 100 Testepisoden berechnet. Forschungshypothese 2 stellt dabei die Behauptung auf, dass eine geringere Varianz der Belohnungen im Testszenario vorliegen unter Auswahl eines Trainingsszenarios mit RL basiertem Gegenspieler. Zur Bestätigung dieser Hypothese gilt es zum einen, eine Ungleichheit der unterliegenden Verteilung der Abweichungen festzustellen. Zum anderen ist es notwendig, dass die Verteilung der Abweichungen zur Basis des RL basiertem Gegenspielers signifikant geringer ist.

Zur Evaluation der Forschungshypothese 2 unterstützen die Ergebnisse der Tabelle 7. Aus der Tabelle wird erkenntlich, dass der P-Wert des Signifikanztests zur Gleichheit bei in etwa $1.2200e-22$ liegt. Da sich dieser Wert unter dem gewählten Signifikanzniveau befindet, kann daraus abgeleitet werden, dass nur sehr unwahrscheinlich eine Gleichheit der unterliegenden Verteilung vorherrscht. Die aufgestellte H_0 Hypothese wird daher verworfen, und die gegensätzliche H_1 Hypothese einer Ungleichheit der beiden Verteilungen wird angenommen. Wird hinzu der P-Wert des Signifikanztests zur Verschlechterung der Belohnungsstabilität von 1.0 betrachtet, kann die folgende Erkenntnis festgestellt werden. Aufgrund des über dem Signifikanzniveau befindlichen P-Wertes, ist die H_0 Hypothese zu bestätigen und von einer geringeren Varianz der erzielten Belohnungen auszugehen unter dem Training mit regelbasiertem Gegenspieler. Eine Verringerung der Belohnungsvarianz durch Optimierung unter RL basierten Gegenspieler, wie in der H_1 Hypothese behauptet, ist demnach auszuschließen.

Wird diese Erkenntnis im Kontext der Messung der Robustheit betrachtet werden zwei Effekte deutlich. Zum Einen zeigte sich mit der Durchführung des Laborexperiments, dass die Gestaltung des Gegenspielers in der Trainingsumgebung einen Effekt auf die Belohnungsstabilität ausübt. Zum Anderen wird deutlich, dass der Einsatz eines optimierten Gegenspielers zu einer Verschlechterung der Belohnungsstabilität führt. Eine RL-Policy, welche unter einem durch verstärkendes Lernen optimierten Gegenspieler trainiert wurde, erzielt demnach eine geringere Robustheit aus der höheren Varianz der Belohnung.

4.1.3 Betrachtung der Anzahl der Misserfolge

Hypothese 3: *Die im Testszenario erreichte Anzahl von Misserfolgen ist unter Verwendung der Policy aus dem Training mit RL basiertem Gegenspieler signifikant und zuverlässig geringer als die Policy aus dem Training mit regelbasiertem Gegenspieler.*

Als dritte Forschungshypothese zur Beantwortung der Forschungsfrage wurde die Behauptung aufgestellt, dass die erreichte Anzahl von Misserfolgen im Testszenario signifikant geringer ist durch den Einsatz des RL-basierten Gegenspielers. Zur Untersuchung dessen wird die Verteilung der erfolgreichen und erfolglosen Testepisoden mittels statistischer Signifikanztests ausgewertet. Weisen die Ergebnisse des Tests zur Gleichheit und Verschlechterung der Metrik je einen P-Wert aus, welcher unter dem Signifikanzniveau liegt, so kann die Forschungshypothese bestätigt und von einem positiven Effekt ausgegangen werden.

Betrachtet man zunächst den P-Wert zur Gleichheit der Verteilungen von erfolgreichen und erfolglosen Episoden, liegt dieser anhand Tabelle 7 bei in etwa 0.1815. Daraus, dass dieser Wert geringer als das Signifikanzniveau ist, wird die H0 Hypothese einer Gleichheit der Verteilungen angenommen. Entsprechend wird die H1 Hypothese einer ungleichen Verteilung der Misserfolge abgelehnt, sodass kein Unterschied in den Erfolgsverteilungen der RL-Policies festzustellen war. Mit der Bestätigung keiner unterschiedlichen Verteilung des Misserfolgs bei einem Training mit RL basiertem Gegenspieler ist die Forschungshypothese 3 abzulehnen.

Das Ablehnen der dritten Forschungshypothese bestätigt, dass kein Effekt durch die Wahl des Trainingsszenarios auf die Robustheit entstanden ist. Wird lediglich die Metrik des Misserfolgs betrachtet konnte festgestellt werden, dass unter unbekanntem Umgebungsverhalten ähnlich gleich häufig Misserfolge mittels der Strategie erzielt worden sind, welche den RL basierten Gegenspieler im Trainingszenario verwendete. Insgesamt kann dadurch festgehalten werden, dass keine Verbesserung der Robustheit mit der im Rahmen dieser Arbeit angewendeten Methodik erzielt wurde.

4.2 Vergleich der Robustheit der Strategien aus dem Training mit regelbasierten Gegenspieler und Domain Randomization

Metrik	Gleichheit P-Wert	Verschlechterung P-Wert	Ungleichheit	Verbesserung
sum reward	0.48390799726891964	0.7588082897727775	False	False
reward mean difference	2.2302103039562726e-05	0.9999889697144758	True	False
count of failures	0.005412977104604817	0.0027064885523024086	True	True

Tab. 8: Ergebnisse der statistischen Tests aus dem Leistungsvergleich der Modelle aus den Trainingszenarien 1 und 4

Tabelle 8 weist die Ergebnisse der statistischen Signifikanztests zum Vergleich der Strategieleistungen aus dem Training mit regelbasiertem Gegenspieler und DR aus. Dabei werden in gleicher Struktur zur Tabelle 7 die P-Werte zur Gleichheit und Verschlechterung der Verteilung sowie die Annahme der H1 Hypothesen anhand jeder Metrik dargestellt. Forschungshypothesen und Testhypothesen sind entsprechend dem Vergleich der Policies aus dem Training mit DR und ohne DR anzupassen.

4.2.1 Betrachtung der kumulierten Belohnung

Hypothese 1: *Die im Testszenario erzielte kumulierte Belohnung ist unter Verwendung der Policy aus dem Training mit DR signifikant und zuverlässig höher als die Policy aus dem Training ohne Randomisierung der Umgebung.*

Für die Ermittlung des Effekts der DR wird die Forschungshypothese 1 für den Vergleich der Strategieleistungen aus dem Training mit und ohne DR adaptiert. Die Hypothese stellt damit die Behauptung auf, dass die kumulierte Belohnung signifikant höher ist, sofern ein Training mit DR vorliegt. Dessen Auswertung erfolgt über eine ähnliche Adaptierung der H_0 und H_1 Hypothesen der Tests zur Gleichheit und Verschlechterung der Ergebnissdaten. Um die adaptierte Hypothese zu bestätigen gilt es eine Ungleichheit im ersten Signifikanztest und eine Verbesserung im zweiten Test festzustellen.

Werden zur Auswertung die P-Werte aus der Tabelle 8 betrachtet liegt unter anderem aus dem Test der Gleichheit ein P-Wert von in etwa 0.4839 vor. Da dieser Wert über dem Signifikanzniveau liegt, ist die entsprechende Null Hypothese anzunehmen. Daraus ist festzustellen, dass die H_1 Hypothese verworfen wird und folglich von einer gleichen Verteilung der Belohnungen in Abhängigkeit des Trainingsszenarios auszugehen ist.

Für das beschriebene Testszenario konnte der Einsatz von DR im Trainingsszenario demnach keine Verbesserung der Belohnungsverteilung, hervorgerufen werden. Die im Rahmen dieser Arbeit verwendete Form der Domänen Randomisierung zeigte somit auch keinen Effekt auf die Robustheit der Strategien gegenüber des unbekannten Testszenarios, wird lediglich die kumulierte Belohnung je Episode betrachtet.

4.2.2 Betrachtung der Belohnungsstabilität

Hypothese 2: *Die Varianz der im Testszenario erzielten kumulierten Belohnung ist unter Verwendung der Policy aus dem Training mit DR signifikant und zuverlässig geringer als die Policy aus dem Training ohne Randomisierung der Umgebung.*

Mittels der adaptierten zweiten Forschungshypothese wird die Verteilung der Abweichungen vom Mittelwert der Belohnungen über alle 100 Testszenarien untersucht. Dazu wird im ersten Test die Gleichheit der Verteilungen, aus der Verwendung der beiden Strategien mit und ohne DR, als H_0 Hypothese vorausgesetzt. Im zweiten Test wird die Verringerung der Verteilung aus der Verwendung der Strategie mit DR angenommen.

Tabelle 8 zeigt zum ersten Test der Gleichheit einen P-Wert von in etwa $2.2302e - 05$. Daraus resultierend können nur zu einer sehr geringen Wahrscheinlichkeit per Zufall die Daten aus einer selben unterliegenden Verteilung gezogen worden sein. Eine gleiche unterliegende Verteilung wird demnach ausgeschlossen und die H_0 Hypothese abgelehnt sowie die H_1 Hypothese angenommen. Bezogen auf den zweiten Test liegt ein P-Wert von in etwa 1.0 vor. Folgend kann abgeleitet werden, dass die H_0 Hypothese einer stärkeren Verteilung der Abweichungen anzunehmen ist. Entsprechend wird die H_1 Hypothese sowie die Hypothese 2 abgelehnt, da die angewendete Form der DR zu einen negativen Effekt auf die Varianz der Belohnungen führt.

Für die Messung der Robustheit ist somit festzuhalten, dass DR im Kontext dieser Arbeit als Teil des Trainingsszenarios, keine Verbesserung der Robustheit im Sinne einer stabileren Belohnung erzielte.

4.2.3 Betrachtung der Anzahl der Misserfolge

Hypothese 3: *Die im Testszenario erreichte Anzahl von Misserfolgen ist unter Verwendung der Policy aus dem Training mit DR signifikant und zuverlässig geringer als die Policy aus dem Training ohne Randomisierung der Umgebung.*

Mit der dritten adaptierten Forschungshypothese soll zuletzt der Effekt der Domänen Randomisierung auf die Anzahl der Misserfolge untersucht werden. Durch den Test der Gleichheit wird die Existenz eines Effekts überprüft. Sofern dieser Test einen vorhandenen Effekt bestätigt, wird dessen Auswirkung mittels des zweiten Tests zur Verschlechterung evaluiert. Damit die dritte Forschungshypothese angenommen wird, gilt es, einen verringernden Effekt auf die Häufigkeit von Misserfolgen festzustellen. Dazu sind die H_0 Hypothesen beider Tests abzulehnen.

Anhand Tabelle 8 lassen sich die Ergebnissdaten der Signifikanztests ablesen. Die Tabelle zeigt einen P-Wert von in etwa 0.0054 zum Test der Gleichheit beider Verteilungen der Misserfolge auf. Aus der erfüllten Bedingung, dass dieser Wert geringer als das Signifikanzniveau ist, kann die H_0 Hypothese abgelehnt werden. Folglich wird für den weiteren Testverlauf angenommen, dass die Verteilung der Misserfolge durch die Auswahl des Trainingsszenarios beeinflusst werden konnte. Wird anschließend der P-Wert zur Evaluation des entstandenen Effekts von in etwa 0.0027 betrachtet, können folgende Schlussfolgerungen gezogen werden. Die H_0 Hypothese einer höheren Häufigkeit von Misserfolgen kann abgelehnt werden, da der P-Wert im Ablehnungsbereich unterhalb des Signifikanzniveaus liegt. Im Gegensatz dazu kann die H_1 Hypothese und eine Verbesserung der Häufigkeit von Misserfolg angenommen werden.

Aus den Annahmen und Ablehnungen der Testhypothesen ergibt sich, dass durch die Anwendung der Domänen Randomisierung eine Verbesserung der Häufigkeit von Misserfolgen erzielen lies. Wird nur diese Metrik betrachtet, kann entsprechend ausgesagt werden, dass die Strategie durch DR robuster auf unbekannte Umgebungen reagiert, da die Testepisoden häufiger erfolgreich abgeschlossen wurden.

4.3 Vergleich der Robustheit der Strategien aus dem Training mit RL basiertem Gegenspieler und Domain Randomization

Mit Hilfe der Tabelle 9 lassen sich die Ergebniswerte der statistischen Signifikanztests zum Vergleich der Strategien aus RL basiertem Gegenspieler und DR im Training darstellen. Dazu sind die P-Werte der Tests und Wahrheitswerte der Testhypothesen je untersuchter Metrik angegeben.

Metrik	Gleichheit P-Wert	Verschlechterung P-Wert	Ungleichheit	Verbesserung
sum reward	0.21860488490723795	0.10930244245361898	False	False
reward mean difference	1.1343177758713004e-13	0.9999999999999443	True	False
count of failures	0.005412977104604817	0.09899034302740754	0.9513450789074275	True

Tab. 9: Ergebnisse der statistischen Tests aus dem Leistungsvergleich der Modelle aus den Trainingsszenarien 3 und 4

4.3.1 Betrachtung der kumulierten Belohnung

Hypothese 1: *Die im Testszenario erzielte kumulierte Belohnung ist unter Verwendung der Policy aus dem Training mit RL basiertem Gegenspieler signifikant und zuverlässig höher als die Policy aus dem Training mit randomisierter Umgebung.*

Zuerst wird die Hypothese 1 einer signifikant höheren kumulierten Belohnung, durch den Einsatz von RL zum optimieren eines Gegenspielers gegenüber der angewendeten Form von DR, untersucht. Die H0 Hypothese im ersten Test repräsentiert die Behauptung das eine die Daten der Belohnungen einer gemeinsamen Verteilung entstammen. Das die Belohnungsdaten welche durch die Strategie mit RL basiertem Gegenspieler erzeugt worden sind einer höheren Verteilung entstammen wird im zweiten Test als H0 Hypothese überprüft. Entsprechend gegensätzliche Annahmen werden in beiden Signifikanztests als H1 Hypothese verwendet.

Der erste P-wert der Tabelle 9 in der Zeile zur kumulierten Belohnung liegt mit in etwa 0.2186 über dem Signifikanzniveau von 10%. Aus diesem Grund wird die H0 Hypothese bestätigt, dass die Daten einer einheitlichen Wahrscheinlichkeitsverteilung entstammen. Entsprechend ist die H1 Hypothese abzulehnen und auch die obige adaptierte Forschungshypothese zu verwerfen. Daraus entfällt eine weitere Betrachtung des P-Wertes des zweiten Tests.

Aus der Überprüfung der Testergebnisse wird deutlich, dass RL keinen verbessernden Effekt auf die Belohnung ausüben konnte als die angewendete Form der DR. Daraus kann abgeleitet werden, dass die Strategie aus dem Training mit RL basiertem Gegenspieler nicht robuster agieren konnte, indem eine bessere Belohnung im abweichenden Testszenario erzielt wird.

4.3.2 Betrachtung der Belohnungsstabilität

Hypothese 2: *Die Varianz der im Testszenario erzielten kumulierten Belohnung ist unter Verwendung der Policy aus dem Training mit RL basiertem Gegenspieler signifikant und zuverlässig geringer als die Policy aus dem Training mit randomisierter Umgebung.*

Mittels der nächsten Forschungshypothese soll die Leistungsstabilität, der aus dem Training mit RL basiertem Gegenspieler und mit DR entstandenen Strategien, verglichen werden. Zur Untersuchung dieser Metrik werden die Daten der Abweichungen der Strategieergebnisse betrachtet.

Als erste H0 Hypothese wird die Annahme einer einheitlichen, als zweite H0 Hypothese die Annahme einer größeren Verteilung eingesetzt.

Aus Tabelle 9 geht ein P-Wert von in etwa $1.1343e - 13$ hervor. Folgend daraus ist die H0 Hypothese abzulehnen und festzuhalten, dass ein Einfluss des Trainingsszenarios auf die Verteilung der Abweichungen besteht. Wird dieser Einfluss mittels zweitem Test genauer untersucht, ist dabei der P-Wert von in etwa 1.0 zu betrachten. Daraus geht hervor, dass die Wahl eines RL basierten Gegenspieler anstelle der DR im Trainingsszenario einen negativen Effekt auf die Belohnungsstabilität ausübt.

Schlussfolgernd konnte damit festgestellt werden, dass ein RL basierter Gegenspieler im Training zu einer weniger robusten Strategie führt, als wenn DR eingesetzt würde.

Hier am Ende vlt auch die wirkliche Varianz berechnen

4.3.3 Betrachtung der Anzahl der Misserfolge

Hypothese 3: *Die im Testszenario erreichte Anzahl von Misserfolgen ist unter Verwendung der Policy aus dem Training mit RL basiertem Gegenspieler signifikant und zuverlässig geringer als die Policy aus dem Training mit randomisierter Umgebung.*

4.4 Beantwortung der Forschungsfrage

- Erläuterung der Bedeutung der einzelnen Metriken zur Robustheit
- Entstandenen Effekte im Bezug zur Auswertung mit DR setzen

5 Reflexion und Forschungsausblick

5.1 Reflexion der Forschungsmethodik

- extrem aufwendige Entwicklung notwendig - eigene Sim hilft aber sollte ausgebaut werden
- Training der Modelle sollte unter mehr Trainingsdaten erfolgen. (Hier besonders auf aggr phy steps eingehen)
- Laborexperiment besitzt beschränkte Sichtweise
- Forschungsfrage kann auf die echte Welt erweitert werden. (vlt mit anderem kompetitiven szenario)
- Standardisierung einer Methodik ist hilfreich

5.2 Ausblick für weitere Forschung

- Standardisierung kompetitive Drohnensimulation
- Erstellung eines Konzepts zur Messung der Robustheit in der echten Welt ?

Anhang

Anhangverzeichnis

Anhang 1	Messdaten des Laborexperiments	58
Anhang 1/1	Messdaten des Modells aus dem ersten Trainingszenario	58
Anhang 1/2	Messdaten des Modells aus dem dritten Trainingszenario	60
Anhang 1/3	Messdaten des Modells aus dem vierten Trainingszenario	63

Anhang 1: Messdaten des Laborexperiments

Anhang 1/1: Messdaten des Modells aus dem ersten Trainingszenario

Episode	max_timesteps	sum_reward	reward_mean_difference	count_of_failed
0	4973	-119446.26307111187	19680.69219803202	1
1	2135	-73579.65118847974	65547.30408066415	1
2	4315	-68451.57618280343	70675.37908634046	0
3	3021	-58856.40652218178	80270.5487469621	1
4	2200	-106828.42153534402	32298.533733799864	1
5	13035	-219399.8086845268	80272.85341538291	1
6	1425	-99294.08419905338	39832.87107009051	1
7	370	-20456.49584066878	118670.45942847511	1
8	2767	-66373.5605798862	72753.3946892577	1
9	394	-70932.35303882814	68194.60223031575	1
10	1377	-77663.76882989446	61463.18643924943	1
11	198	-12715.698783642225	126411.25648550167	1
12	6360	-217357.82021610026	78230.86494695637	1
13	776	-81468.55281955525	57658.40244958864	1
14	329	-18654.276777649084	120472.67849149481	1
15	160	-57678.73562351338	81448.21964563051	1
16	3646	-167379.6638554781	28252.70858633422	1
17	386	-24654.250630269184	114472.70463887471	1
18	936	-71691.87001738523	67435.08525175866	1
19	1868	-117854.37762945828	21272.57763968561	1
20	15220	-247291.44889176913	108164.49362262525	1
21	1854	-91688.46036564803	47438.494903495855	1
22	1949	-82100.30706836819	57026.6482007757	1
23	5635	-146728.2040420022	7601.248772858322	1
24	3915	-100710.63153067631	38416.323738467574	1
25	27149	-288748.7499041229	149621.794634979	1
26	17351	-773888.6594350344	634761.7041658905	1
27	18452	-266545.7001431137	127418.74487396979	1
28	4116	-128943.76280506424	10183.19246407965	1
29	4245	-146624.52291664132	7497.567647497432	1
30	693	-24905.009835843994	114221.9454332999	1
31	4204	-108705.23049853576	30421.724770608125	1
32	3839	-105973.55113709459	33153.4041320493	1
33	13880	-102975.14464771489	36151.810621429	1
34	1575	-114665.39035433184	24461.564914812043	1
35	2368	-104427.56060357616	34699.39466556773	1

Episode	max_timesteps	sum_reward	reward_mean_difference	count_of_failed
36	1223	-83398.23783177527	55728.717437368614	1
37	33802	-735697.9326871325	596570.9774179886	1
38	11678	-280712.0211739327	141585.06590478882	1
39	4613	-103340.52131554208	35786.43395360181	1
40	4939	-113005.67033766753	26121.28493147636	1
41	557	-34873.86419424162	104253.09107490227	1
42	2921	-111742.82137255678	27384.13389658711	1
43	1557	-73030.91534860493	66096.03992053895	1
44	1522	21140.63597449921	160267.5912436431	0
45	1954	-90232.52886299146	48894.42640615243	1
46	5999	7264.712146483675	146391.66741562757	0
47	281	-16075.567283123166	123051.38798602072	1
48	2070	12393.512751061113	151520.468020205	0
49	19418	-473654.72629417136	334527.7710250275	1
50	15689	-276474.9519422902	137347.9966731463	1
51	5238	-165381.53396681393	26254.578697670047	1
52	6524	-105787.99825142126	33338.95701772263	1
53	6561	-205898.10139098376	66771.14612183988	1
54	1910	-92032.12884181371	47094.826427330176	1
55	22148	-445050.6391527784	305923.6838836345	1
56	1472	-79755.73736719266	59371.217901951226	1
57	263	-15600.021604753363	123526.93366439053	1
58	2567	-61654.41882784868	77472.5364412952	1
59	1366	-51713.97716429973	87412.97810484417	1
60	4928	-95423.07511540131	43703.880153742575	1
61	304	-13998.380734446198	125128.57453469769	1
62	12421	-71413.81905101251	67713.13621813137	0
63	311	-14710.14499937111	124416.81026977277	1
64	2411	-90587.72499755406	48539.230271589826	1
65	3560	-156842.56739543795	17715.612126294058	1
66	3323	-88666.1266710266	50460.82859811728	1
67	2216	-102406.93034675243	36720.02492239146	1
68	1168	-43261.54605595228	95865.4092131916	1
69	2261	-111291.51604653947	27835.439222604415	1
70	1803	26669.781018645506	165796.7362877894	0
71	5358	-82036.28655387607	57090.66871526782	1
72	6518	-35354.203289041674	103772.75198010221	0
73	100662	-485724.2106050032	346597.25533585926	1
74	1642	-112009.36921540477	27117.586053739113	1
75	23079	-639349.6775980042	500222.7223288603	0

Episode	max_timesteps	sum_reward	reward_mean_difference	count_of_failed
76	2505	-114534.05479872528	24592.900470418608	1
77	9151	-159660.46156868307	20533.506299539178	1
78	5556	-184610.73783700686	45483.782567862974	1
79	2591	-113013.74205295392	26113.213216189964	1
80	1917	-48611.60139315391	90515.35387598997	1
81	1033	-65091.71325664009	74035.24201250379	1
82	3035	-151284.32248952115	12157.367220377258	1
83	6333	-248964.96189519798	109838.00662605409	1
84	143	-6270.742974679765	132856.21229446412	1
85	2408	35063.594815062024	174190.55008420593	0
86	3609	-120362.11771286349	18764.8375562804	1
87	2907	-95234.46438855604	43892.49088058785	1
88	6909	-254854.3644192685	115727.40915012461	1
89	606	38778.579022998085	177905.53429214197	0
90	4813	-111878.77036680779	27248.1849023361	1
91	2321	-93649.80171388677	45477.153555257115	1
92	88171	-910404.3026574675	771277.3473883236	1
93	3108	-141655.67094119178	2528.715672047896	1
94	10312	-330954.2408898612	191827.2856207173	1
95	3768	-159059.9077657962	19932.952496652317	1
96	5018	-73322.6603574858	65804.29491165809	1
97	3678	-92350.12255671938	46776.83271242451	1
98	6036	-146973.86740964398	7846.9121405000915	1
99	1573	-87413.84910489814	51713.10616424575	1

Anhang 1/2: Messdaten des Modells aus dem dritten Trainingszenario

Episode	max_timesteps	sum_reward	reward_mean_difference	count_of_failed
0	1653	-92347.12814609804	1879.280532867153	1
1	1511	-79949.83767800787	14276.571000957323	1
2	2686	-138673.8482117061	44447.4395327409	1
3	2076	-111671.73332763254	17445.324648667345	1
4	1251	-79799.73697424782	14426.671704717373	1
5	1832	-77824.71257469118	16401.696104274015	1
6	2081	-100702.72327265213	6476.3145936869405	1
7	2037	-92548.82473529059	1677.5839436746028	1
8	1546	-80243.20957471852	13983.199104246669	1
9	2100	-104374.48557198321	10148.07689301802	1
10	1739	-99047.36821871552	4820.959539750329	1
11	2437	-114813.9047312334	20587.49605226821	1

Episode	max_timesteps	sum_reward	reward_mean_difference	count_of_failed
12	2006	-108572.06517066105	14345.65649169586	1
13	1564	-89733.93105526718	4492.477623698011	1
14	2415	-125882.56288775698	31656.15420879179	1
15	2352	-123751.17987700357	29524.771198038376	1
16	2041	-107225.85252142692	12999.443842461726	1
17	1744	-89764.21221744212	4462.196461523068	1
18	2137	-119053.14423147225	24826.73555250706	1
19	1519	-87000.82563168688	7225.583047278313	1
20	1404	-81501.43041543056	12724.978263534635	1
21	1744	-103097.75428849243	8871.345609527241	1
22	1667	-89130.51424513668	5095.894433828507	1
23	1642	-6907.798450224385	87318.61022874081	0
24	1343	-83084.15449344112	11142.25418552407	1
25	1683	-90356.76164617151	3869.647032793684	1
26	1878	-111838.56754250827	17612.15886354308	1
27	1852	-91691.71347423723	2534.695204727963	1
28	1378	5312.971207702023	99539.37988666722	0
29	1950	-91346.0914973505	2880.3171816146933	1
30	1751	-14833.500644185704	79392.90803477948	0
31	2121	-122840.91833738518	28614.509658419993	1
32	2027	-107247.44465804833	13021.035979083143	1
33	2222	-100644.25859708953	6417.849918124339	1
34	1879	-99075.08534082912	4848.676661863923	1
35	1330	-89172.55540681427	5053.8532721509255	1
36	1735	-81005.49859622098	13220.910082744216	1
37	2410	-122139.76801712476	27913.35933815957	1
38	2002	-110744.31609087711	16517.907411911918	1
39	2110	-104030.87721611398	9804.468537148787	1
40	2368	-127586.09074176587	33359.68206280068	1
41	1779	-88839.11463014259	5387.294048822601	1
42	1563	-92566.27084313083	1660.1378358343645	1
43	2262	-120846.561556859	26620.15287789381	1
44	1871	-110632.23542961734	16405.82675065215	1
45	1695	-90088.13234400195	4138.276334963244	1
46	1920	-101160.33136423607	6933.922685270882	1
47	1467	-92474.14392801689	1752.2647509483068	1
48	2328	-114691.48989933486	20465.08122036967	1
49	1520	-88746.39566596689	5480.013012998301	1
50	1829	-105687.92705010656	11461.51837114137	1
51	1068	-80491.15858977714	13735.250089188048	1

Episode	max_timesteps	sum_reward	reward_mean_difference	count_of_failed
52	2549	-115674.49520694178	21448.086527976586	1
53	1460	-80996.95173668266	13229.45694228253	1
54	1927	-112770.6757002326	18544.267021267413	1
55	1422	-76259.59982803596	17966.80885092923	1
56	1930	-94052.10799082342	174.30068814176775	1
57	1949	-98282.0069838972	4055.5983049320057	1
58	1672	-77253.73157348904	16972.67710547615	1
59	1537	-105266.57754039323	11040.168861428043	1
60	2075	-113502.44137400572	19276.03269504053	1
61	1522	-95541.05480403255	1314.6461250673601	1
62	1435	-88259.02829929876	5967.380379666429	1
63	1133	-70536.33448700237	23690.07419196282	1
64	1035	-79043.81093353551	15182.59774542968	1
65	2223	-103816.50438886588	9590.095709900692	1
66	1784	-118544.62175341861	24318.213074453422	1
67	1037	28673.595925229718	122900.00460419491	0
68	1232	14926.938427057503	109153.3471060227	0
69	2659	-151453.26593478158	57226.85725581639	1
70	1688	-94274.73392669848	48.32524773328623	1
71	2299	-113496.06004772757	19269.65136876238	1
72	1693	-99589.95555226665	5363.546873301457	1
73	1959	-104666.89017587205	10440.481496906854	1
74	1338	-76191.43750774984	18034.97117121535	1
75	1799	-107730.14072747264	13503.732048507445	1
76	1726	-94312.72168004303	86.31300107784045	1
77	1637	-79745.42987692845	14480.978802036741	1
78	1896	-98852.78875571876	4626.380076753572	1
79	1644	-84666.40512403368	9560.003554931507	1
80	2133	-101018.20279461997	6791.794115654775	1
81	2230	-94581.92219002724	355.51351106204675	1
82	1691	-77884.60688015005	16341.80179881514	1
83	2122	-119640.7616723019	25414.35299333671	1
84	1324	-70624.74418543525	23601.66449352994	1
85	2216	-138310.5053513642	44084.09667239901	1
86	2176	-95513.39386064193	1286.9851816767332	1
87	1864	-88412.25104209002	5814.15763687517	1
88	1786	-92727.62776697357	1498.78091199162	1
89	1486	-72737.9740066259	21488.434672339296	1
90	2496	-119080.5852518373	24854.176572872107	1
91	1948	-101354.26852030956	7127.859841344369	1

Episode	max_timesteps	sum_reward	reward_mean_difference	count_of_failed
92	1958	-110046.45887832022	15820.050199355028	1
93	2006	-101985.01821046043	7758.60953149524	1
94	1059	-75028.5649990723	19197.84367989289	1
95	1431	-88123.49672449603	6102.911954469164	1
96	2109	-96440.1698084757	2213.7611295105016	1
97	2137	-135963.8132262433	41737.404547278114	1
98	1827	-104948.82027461141	10722.41159564622	1
99	1689	-114871.268292174	20644.859613208813	1

Anhang 1/3: Messdaten des Modells aus dem vierten Trainingszenario

Episode	max_timesteps	sum_reward	reward_mean_difference	count_of_failed
0	1455	-139572.38780625432	27287.931729819888	1
1	668	-73773.32807283304	38511.12800360139	1
2	401	-31868.889712731754	80415.56636370267	1
3	2566	-171421.8210663553	59137.36498992087	1
4	568	-84374.66627949841	27909.78979693602	1
5	1369	-49952.79194271706	62331.66413371738	1
6	4053	-155445.12575918084	43160.669682746404	1
7	2124	-105875.03948737303	6409.4165890614095	1
8	2559	-100702.63656066287	11581.819515771567	1
9	3690	-189124.45404525608	76839.99796882164	1
10	10912	-435051.3296926795	322766.8736162451	1
11	4692	-184106.58067775535	71822.12460132092	1
12	486	-60047.127430950866	52237.32864548357	1
13	10080	-66548.61530662145	45735.84076981299	0
14	3312	-122858.94017941065	10574.484102976217	1
15	391	-19717.436219625506	92567.01985680893	1
16	743	-83506.97471680315	28777.48135963129	1
17	1878	-103106.10932391894	9178.3467525155	1
18	1016	-85023.93260886864	27260.5234675658	1
19	1107	-104535.49176005676	7748.964316377678	1
20	1532	-42241.6360096946	70042.82006673983	1
21	3524	-117124.9053152941	4840.44923885967	1
22	203	-58907.64565455435	53376.81042188009	1
23	4397	-121775.39123662916	9490.935160194727	1
24	3493	-151743.9720871852	39459.516010750755	1
25	1468	-94328.55530234546	17955.90077408898	1
26	1667	-103489.34328833577	8795.11278809866	1
27	2227	-203299.81129311383	91015.3552166794	1

Episode	max_timesteps	sum_reward	reward_mean_difference	count_of_failed
28	1663	-80731.42288443842	31553.033191996015	1
29	687	-88827.11715052626	23457.338925908174	1
30	932	-79522.16509663286	32762.290979801575	1
31	1047	-98481.09856043328	13803.357516001153	1
32	3539	-297297.80747770553	185013.3514012711	1
33	6516	-265886.76237301354	153602.3062965791	1
34	1474	-87107.15771247805	25177.29836395639	1
35	2945	-127935.2128341831	15650.756757748662	1
36	289	-58460.85077567921	53823.605300755225	1
37	2454	-92160.30137319927	20124.154703235166	1
38	1126	-85634.7327863673	26649.723290067137	1
39	357	-74902.30855478215	37382.14752165228	1
40	1223	-89556.25027933848	22728.205797095958	1
41	2807	-187277.46219272964	74993.0061162952	1
42	3295	-128034.55609537168	15750.100018937243	1
43	960	-68842.4879640358	43441.96811239864	1
44	370	-63979.30353180085	48305.15254463359	1
45	1259	-97077.13882747927	15207.317248955165	1
46	247	-65317.699554064646	46966.75652236979	1
47	3634	-247874.40333470656	135589.94725827212	1
48	2321	-147716.59825833945	35432.142181905016	1
49	1994	-154343.7654882401	42059.309411805676	1
50	758	-75267.45162457463	37017.0044518598	1
51	803	-96696.74430428469	15587.711772149749	1
52	595	-79127.62418066339	33156.83189577104	1
53	386	-20780.716701622558	91503.73937481188	1
54	928	-74235.79261964842	38048.66345678602	1
55	763	-83532.89926826739	28751.55680816705	1
56	1831	-101537.94911802994	10746.506958404498	1
57	1274	-122700.53145025088	10416.07537381645	1
58	802	-43329.848793566314	68954.60728286812	1
59	10117	-227067.8234075841	114783.36733114967	1
60	1728	-119408.46281307109	7124.006736636657	1
61	849	-99708.72749570434	12575.728580730094	1
62	1246	-82510.6863451565	29773.769731277935	1
63	288	-15740.96284696191	96543.49322947253	1
64	362	-78500.3086232856	33784.147453148835	1
65	541	-66739.89268245646	45544.56339397798	1
66	57	-51426.15702325551	60858.299053178926	1
67	523	-69521.75031680161	42762.70575963282	1

Episode	max_timesteps	sum_reward	reward_mean_difference	count_of_failed
68	471	-63306.9524595423	48977.503616892136	1
69	1451	-57554.22148378302	54730.23459265142	1
70	677	-82659.93849905756	29624.51757737687	1
71	888	-96043.9975492387	16240.458527195733	1
72	1402	-105520.38171801696	6764.074358417478	1
73	184	-56333.05585696981	55951.400219464624	1
74	1608	-85394.50950372484	26889.946572709596	1
75	2691	-100272.7277743975	12011.728302036936	1
76	5126	-241840.43650942043	129555.98043298599	1
77	371	-29223.30313552408	83061.15294091035	1
78	168	-11817.715457175556	100466.74061925888	1
79	565	-73200.09634163367	39084.35973480076	1
80	2239	-134441.5242029784	22157.068126543978	1
81	2165	-121493.00241384219	9208.546337407752	1
82	202	-54308.91006303916	57975.54601339527	1
83	1288	-87881.2787390615	24403.177337372937	1
84	5135	-173571.8733679027	61287.41729146827	1
85	257	-66944.76429320117	45339.69178323327	1
86	61548	-601478.1387180209	489193.6826415864	1
87	823	-84664.817011253	27619.63906518143	1
88	272	-59390.91893604745	52893.53714038699	1
89	2467	-192004.36119959655	79719.90512316211	1
90	1043	-74929.26090746139	37355.19516897305	1
91	2149	-105531.93105444917	6752.5250219852605	1
92	5500	-476829.080297996	364544.6242215616	1
93	907	-62796.06332880723	49488.39274762721	1
94	916	-96615.1085959957	15669.347480438737	1
95	984	-79178.95030215883	33105.5057742756	1
96	45	-52242.449477733804	60042.00659870063	1
97	424	-76972.4110787145	35312.04499771993	1
98	2467	-50969.970664024964	61314.48541240947	1
99	2836	-118707.5851452389	6423.129068804468	1

Literaturverzeichnis

- ACM Digital Library (2/28/2023): ACM Digital Library. URL: <https://dl.acm.org/>.
- Alghonaim, R./Johns, E. (2021): Benchmarking Domain Randomisation for Visual Sim-to-Real Transfer. In: *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, S. 12802–12808. ISBN: 978-1-7281-9077-8. DOI: 10.1109/ICRA48506.2021.9561134.
- Arulkumaran, K./Deisenroth, M. P./Brundage, M./Bharath, A. A. (2017): Deep Reinforcement Learning: A Brief Survey. In: *IEEE Signal Processing Magazine* 34.6, S. 26–38. DOI: 10.1109/MSP.2017.2743240.
- Attia, A./Dayan, S. (2018): Global overview of Imitation Learning. arXiv: 1801.06503 [stat.ML].
- Ayala, A./Cruz, F./Campos, D./Rubio, R./Fernandes, B./Dazeley, R. (2020): A Comparison of Humanoid Robot Simulators: A Quantitative Approach. In: *2020 Joint IEEE 10th International Conference on Development and Learning and Epigenetic Robotics (ICDL-EpiRob)*, S. 1–6. DOI: 10.1109/ICDL-EpiRob48136.2020.9278116.
- Balakrishna, A./Thananjeyan, B./Lee, J./Li, F./Zahed, A./Gonzalez, J. E./Goldberg, K. (2020): On-Policy Robot Imitation Learning from a Converging Supervisor. In: *Proceedings of the Conference on Robot Learning*. Hrsg. von Leslie Pack Kaelbling/Danica Kragic/Komei Sugiura. Bd. 100. Proceedings of Machine Learning Research. PMLR, S. 24–41. URL: <https://proceedings.mlr.press/v100/balakrishna20a.html>.
- Bellman, R. (1966): Dynamic Programming. In: *Science* 153.3731, S. 34–37. DOI: 10.1126/science.153.3731.34. eprint: <https://www.science.org/doi/pdf/10.1126/science.153.3731.34>. URL: <https://www.science.org/doi/abs/10.1126/science.153.3731.34>.
- Bharadhwaj, H./Wang, Z./Bengio, Y./Paull, L. (2019): A Data-Efficient Framework for Training and Sim-to-Real Transfer of Navigation Policies. In: *2019 International Conference on Robotics and Automation (ICRA)*. IEEE. DOI: 10.1109/icra.2019.8794310.
- Brockman, G./Cheung, V./Pettersson, L./Schneider, J./Schulman, J./Tang, J./Zaremba, W. (2016): OpenAI Gym. In: *CoRR* abs/1606.01540. arXiv: 1606.01540. URL: <http://arxiv.org/abs/1606.01540>.
- Canese, L./Cardarilli, G. C./Di Nunzio, L./Fazzolari, R./Giardino, D./Re, M./Spanò, S. (2021): Multi-Agent Reinforcement Learning: A Review of Challenges and Applications. In: *Applied Sciences* 11.11, S. 4948. DOI: 10.3390/app11114948. URL: <https://www.mdpi.com/2076-3417/11/11/4948>.
- Chen, X./Hu, J./Jin, C./Li, L./Wang, L. (2021): Understanding Domain Randomization for Sim-to-real Transfer. In: *CoRR* abs/2110.03239. arXiv: 2110.03239. URL: <https://arxiv.org/abs/2110.03239>.
- Collins, J./Ketter, W. (2022): Power TAC: Software architecture for a competitive simulation of sustainable smart energy markets. In: *SoftwareX* 20, S. 101217. ISSN: 2352-7110. DOI: <https://doi.org/10.1016/j.softx.2022.101217>. URL: <https://www.sciencedirect.com/science/article/pii/S2352711022001352>.

- Cutler, M./Walsh, T. J./How, J. P. (2014):** Reinforcement learning with multi-fidelity simulators. In: *2014 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE. DOI: 10.1109/icra.2014.6907423.
- Deshpande, A. M./Kumar, R./Minai, A. A./Kumar, M. (2020):** Developmental Reinforcement Learning of Control Policy of a Quadcopter UAV with Thrust Vectoring Rotors. URL: <https://arxiv.org/pdf/2007.07793>.
- Deshpande, A. M./Minai, A. A./Kumar, M. (2021):** Robust Deep Reinforcement Learning for Quadcopter Control. In: *IFAC-PapersOnLine* 54.20, S. 90–95. ISSN: 24058963. DOI: 10.1016/j.ifacol.2021.11.158.
- Fang, B./Jia, S./Guo, D./Xu, M./Wen, S./Sun, F. (2019):** Survey of imitation learning for robotic manipulation. In: *International Journal of Intelligent Robotics and Applications* 3, S. 362–369.
- Foronda, C. L. (2021):** What Is Virtual Simulation? In: *Clinical Simulation in Nursing* 52, S. 8. ISSN: 18761399. DOI: 10.1016/j.ecns.2020.12.004.
- Furrer, F./Burri, M./Achtelik, M./Siegwart, R. (2016):** RotorS—A Modular Gazebo MAV Simulator Framework. In: *Robot Operating System (ROS): The Complete Reference (Volume 1)*. Hrsg. von Anis Koubaa. Cham: Springer International Publishing, S. 595–625. ISBN: 978-3-319-26054-9. DOI: 10.1007/978-3-319-26054-9_23. URL: https://doi.org/10.1007/978-3-319-26054-9_23.
- Gao, Y./Vedula, S. S./Reiley, C. E./Ahmidi, N./Varadarajan, B./Lin, H. C./Tao, L./Zappella, L./Béjar, B./Yuh, D. D. u. a. (2014):** Jhu-isi gesture and skill assessment working set (jigsaws): A surgical activity dataset for human motion modeling. In: *MICCAI workshop: M2cai*. Bd. 3. 3.
- GitHub (4/4/2023):** Farama-Foundation/Gymnasium: A standard API for single-agent reinforcement learning environments, with popular reference environments and related utilities (formerly Gym). URL: <https://github.com/Farama-Foundation/Gymnasium>.
- Google Scholar (2/28/2023). URL: <https://scholar.google.de/>.
- Haarnoja, T./Zhou, A./Abbeel, P./Levine, S. (2018):** Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor. In: *CoRR* abs/1801.01290. arXiv: 1801.01290. URL: <http://arxiv.org/abs/1801.01290>.
- Hentati, A. I./Krichen, L./Fourati, M./Fourati, L. C. (2018):** Simulation Tools, Environments and Frameworks for UAV Systems Performance Analysis. In: *2018 14th International Wireless Communications & Mobile Computing Conference (IWCMC)*. IEEE. DOI: 10.1109/iwcmc.2018.8450505.
- Holzweißig, K. (2022):** Wissenschaftliches Arbeiten. In: 6.6.
- Hsu, K.-C./Ren, A. Z./Nguyen, D. P./Majumdar, A./Fisac, J. F. (2023):** Sim-to-Lab-to-Real: Safe reinforcement learning with shielding and generalization guarantees. In: *Artificial Intelligence* 314, S. 103811. ISSN: 0004-3702. DOI: <https://doi.org/10.1016/j.artint.2022.103811>. URL: <https://www.sciencedirect.com/science/article/pii/S0004370222001515>.

- Huang, S. H./Papernot, N./Goodfellow, I. J./Duan, Y./Abbeel, P. (2017):** Adversarial Attacks on Neural Network Policies. In: *CoRR* abs/1702.02284. arXiv: 1702.02284. URL: <http://arxiv.org/abs/1702.02284>.
- Hussein, A./Gaber, M. M./Elyan, E./Jayne, C. (2017):** Imitation Learning: A Survey of Learning Methods. In: *ACM Comput. Surv.* 50.2. ISSN: 0360-0300. DOI: 10.1145/3054912. URL: <https://doi.org/10.1145/3054912>.
- IEEE Xplore (2/28/2023). URL: <https://ieeexplore.ieee.org/Xplore/home.jsp>.
- Ivaldi, S./Padois, V./Nori, F. (2014):** Tools for dynamics simulation of robots: a survey based on user feedback. URL: <https://arxiv.org/pdf/1402.7050>.
- Koch, W./Mancuso, R./West, R./Bestavros, A. (2018):** Reinforcement Learning for UAV Attitude Control. In: *CoRR* abs/1804.04154. arXiv: 1804.04154. URL: <http://arxiv.org/abs/1804.04154>.
- Körber, M./Lange, J./Rediske, S./Steinmann, S./Glück, R. (2021):** Comparing Popular Simulation Environments in the Scope of Robotics and Reinforcement Learning. URL: <https://arxiv.org/pdf/2103.04616>.
- Li, Y. (2019):** Reinforcement Learning Applications. DOI: 10.48550/ARXIV.1908.06973. URL: <https://arxiv.org/abs/1908.06973>.
- Liu, Z./Guo, Z./Cen, Z./Zhang, H./Tan, J./Li, B./Zhao, D. (2023):** On the Robustness of Safe Reinforcement Learning under Observational Perturbations. arXiv: 2205.14691 [cs.LG].
- Maria, A. (1997):** Introduction to modeling and simulation. In: *Proceedings of the 29th conference on Winter simulation - WSC '97*. New York, New York, USA: ACM Press. DOI: 10.1145/268437.268440.
- Mnih, V./Badia, A. P./Mirza, M./Graves, A./Lillicrap, T. P./Harley, T./Silver, D./Kavukcuoglu, K. (2016):** Asynchronous Methods for Deep Reinforcement Learning. In: *CoRR* abs/1602.01783. arXiv: 1602.01783. URL: <http://arxiv.org/abs/1602.01783>.
- Mnih, V./Kavukcuoglu, K./Silver, D./Graves, A./Antonoglou, I./Wierstra, D./Riedmiller, M. A. (2013):** Playing Atari with Deep Reinforcement Learning. In: *CoRR* abs/1312.5602. arXiv: 1312.5602. URL: <http://arxiv.org/abs/1312.5602>.
- Molchanov, A./Chen, T./Hönig, W./Preiss, J. A./Ayanian, N./Sukhatme, G. S. (2019):** Sim-to-(Multi)-Real: Transfer of Low-Level Robust Control Policies to Multiple Quadrotors. In: *CoRR* abs/1903.04628. arXiv: 1903.04628. URL: <http://arxiv.org/abs/1903.04628>.
- Moos, J./Hansel, K./Abdulsamad, H./Stark, S./Clever, D./Peters, J. (2022):** Robust Reinforcement Learning: A Review of Foundations and Recent Advances. In: *Machine Learning and Knowledge Extraction* 4.1, S. 276–315. ISSN: 2504-4990. DOI: 10.3390/make4010013. URL: <https://www.mdpi.com/2504-4990/4/1/13>.
- Ningombam, D. D. (2022):** Deep Reinforcement Learning Algorithms for Machine-to-Machine Communications: A Review. In: *2022 13th International Conference on Computing Communication and Networking Technologies (ICCCNT)*. IEEE. DOI: 10.1109/icccnt54827.2022.9984457.

- Pan, A./Lee, Y./Zhang, H./Chen, Y./Shi, Y. (2021):** Improving Robustness of Reinforcement Learning for Power System Control with Adversarial Training. In: *arXiv e-prints*, arXiv:2110.08956, arXiv:2110.08956. DOI: 10.48550/arXiv.2110.08956. arXiv: 2110.08956 [eess.SY].
- Panerati, J./Zheng, H./Zhou, S./Xu, J./Prorok, A./Schoellig, A. P. (2021):** Learning to Fly – a Gym Environment with PyBullet Physics for Reinforcement Learning of Multi-agent Quadcopter Control. URL: <https://arxiv.org/pdf/2103.02142>.
- Pinto, L./Davidson, J./Sukthankar, R./Gupta, A. (2017):** Robust Adversarial Reinforcement Learning. In: *CoRR* abs/1703.02702. arXiv: 1703.02702. URL: <http://arxiv.org/abs/1703.02702>.
- Pullum, L. L. (2022):** Review of Metrics to Measure the Stability, Robustness and Resilience of Reinforcement Learning. arXiv: 2203.12048 [cs.LG].
- Raffin, A./Hill, A./Gleave, A./Kanervisto, A./Ernestus, M./Dormann, N. (2021):** Stable-Baselines3: Reliable Reinforcement Learning Implementations. In: *J. Mach. Learn. Res.* 22.1. ISSN: 1532-4435.
- Recker, J. (2021):** Scientific research in information systems: A beginner’s guide. Second Edition. Progress in IS. Cham: Springer International Publishing. ISBN: 9783030854362. URL: <https://ebookcentral.proquest.com/lib/kxp/detail.action?docID=6789173>.
- Reda, D./Tao, T./van de Panne, M. (2020):** Learning to Locomote: Understanding How Environment Design Matters for Deep Reinforcement Learning. In: *Motion, Interaction and Games*. Hrsg. von Daniele Reda/Tianxin Tao/Michiel van de Panne. New York, NY, USA: ACM, S. 1–10. DOI: 10.1145/3424636.3426907.
- Sadeghi, F./Levine, S. (2016):** CAD2RL: Real Single-Image Flight without a Single Real Image. In: *CoRR* abs/1611.04201. arXiv: 1611.04201. URL: <http://arxiv.org/abs/1611.04201>.
- Schott, L./Hajri, H./Lamprier, S. (2022):** Improving Robustness of Deep Reinforcement Learning Agents: Environment Attack based on the Critic Network. In: *2022 International Joint Conference on Neural Networks (IJCNN)*, S. 1–8. DOI: 10.1109/IJCNN55064.2022.9892901.
- Schuderer, A./Bromuri, S./van Eekelen, M. (2021):** Sim-Env: Decoupling OpenAI Gym Environments from Simulation Models. In: *International Conference on Practical Applications of Agents and Multi-Agent Systems*. Springer, Cham, S. 390–393. DOI: 10.1007/978-3-030-85739-4_{\text{underscore}}39. URL: https://link.springer.com/chapter/10.1007/978-3-030-85739-4_39.
- Schulman, J./Levine, S./Moritz, P./Jordan, M. I./Abbeel, P. (2015):** Trust Region Policy Optimization. In: *CoRR* abs/1502.05477. arXiv: 1502.05477. URL: <http://arxiv.org/abs/1502.05477>.
- Schulman, J./Wolski, F./Dhariwal, P./Radford, A./Klimov, O. (2017):** Proximal Policy Optimization Algorithms. In: *CoRR* abs/1707.06347. arXiv: 1707.06347. URL: <http://arxiv.org/abs/1707.06347>.

- Shah, S./Dey, D./Lovett, C./Kapoor, A. (2017):** AirSim: High-Fidelity Visual and Physical Simulation for Autonomous Vehicles. In: *CoRR* abs/1705.05065. arXiv: 1705.05065. URL: <http://arxiv.org/abs/1705.05065>.
- Silano, G./Iannelli, L. (2019):** CrazyS: A Software-in-the-Loop Simulation Platform for the Crazyflie 2.0 Nano-Quadcopter. In: *Robot Operating System (ROS): The Complete Reference (Volume 4)*. Hrsg. von Anis Koubaa. Cham: Springer International Publishing, S. 81–115. ISBN: 978-3-030-20190-6. DOI: 10.1007/978-3-030-20190-6_4. URL: https://doi.org/10.1007/978-3-030-20190-6_4.
- Slaoui, R. B./Clements, W. R./Foerster, J. N./Toth, S. (2019):** Robust Domain Randomization for Reinforcement Learning. In: *CoRR* abs/1910.10537. arXiv: 1910.10537. URL: <http://arxiv.org/abs/1910.10537>.
- Sutton, R. S./Barto, A. G. (2018):** Reinforcement Learning, second edition: An Introduction. MIT Press. ISBN: 9780262352703.
- Tobin, J./Fong, R./Ray, A./Schneider, J./Zaremba, W./Abbeel, P. (2017):** Domain randomization for transferring deep neural networks from simulation to the real world. In: *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, S. 23–30. DOI: 10.1109/IROS.2017.8202133.
- Todorov, E./Erez, T./Tassa, Y. (2012):** MuJoCo: A physics engine for model-based control. In: *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, S. 5026–5033. DOI: 10.1109/IROS.2012.6386109.
- Wang, Z./Hong, T. (2020):** Reinforcement learning for building controls: The opportunities and challenges. In: *Applied Energy* 269, S. 115036. ISSN: 0306-2619. DOI: 10.1016/j.apenergy.2020.115036.
- Webster, J./Watson, R. T. (2002):** Analyzing the Past to Prepare for the Future: Writing a Literature Review. In: *MIS Q.* 26.2, S. xiii–xxiii. ISSN: 0276-7783.
- Wong, A./Bäck, T./Kononova, A. V./Plaat, A. (2022):** Deep multiagent reinforcement learning: challenges and directions. In: *Artificial Intelligence Review*. ISSN: 0269-2821. DOI: 10.1007/s10462-022-10299-x.
- Yan Duan/Xi Chen/Rein Houthooft/John Schulman/Pieter Abbeel (2016):** Benchmarking Deep Reinforcement Learning for Continuous Control. In: *International Conference on Machine Learning*, S. 1329–1338. ISSN: 1938-7228. URL: <https://proceedings.mlr.press/v48/duan16.html>.
- Zhai, P./Hou, T./Ji, X./Dong, Z./Zhang, L. (2022):** Robust Adaptive Ensemble Adversary Reinforcement Learning. In: *IEEE Robotics and Automation Letters* 7.4, S. 12562–12568. DOI: 10.1109/LRA.2022.3220531.
- Zhang, A./Wu, Y./Pineau, J. (2018):** Natural Environment Benchmarks for Reinforcement Learning. DOI: 10.48550/ARXIV.1811.06032. URL: <https://arxiv.org/abs/1811.06032>.
- Zhao, W./Queralta, J. P./Westerlund, T. (2020):** Sim-to-Real Transfer in Deep Reinforcement Learning for Robotics: a Survey. In: *2020 IEEE Symposium Series on Computational Intelligence (SSCI)*. IEEE. DOI: 10.1109/ssci47803.2020.9308468.

Erklärung

Ich versichere hiermit, dass ich meine Bachelorarbeit mit dem Thema: *Experiment zur Verbesserung der Robustheit von Reinforcement Learning Modellen in Drohnensimulationen anhand trainiertem Gegenspieler* selbstständig verfasst und keine anderen als die angegebenen Quellen und Hilfsmittel benutzt habe. Ich versichere zudem, dass die eingereichte elektronische Fassung mit der gedruckten Fassung übereinstimmt.

(Ort, Datum)

(Unterschrift)