

DEVELOPING AN ETHICS MODULE FOR A SERVICE ROBOT

by

Kadriye Yasemin Usta

B.S., Computer Engineering, Boğaziçi University, 2015

Submitted to the Institute for Graduate Studies in
Science and Engineering in partial fulfillment of
the requirements for the degree of
Master of Science

Graduate Program in Computer Engineering
Boğaziçi University

2019

ACKNOWLEDGEMENTS

First of all, I would like to express my special appreciation, and thanks to my advisor Professor Dr. H.Levent Akin, you have been a tremendous and tolerant mentor for me. Moreover, I would like to thank you for encouraging my research and for allowing me to grow as a research scientist. Your advice, guidance and unending patience on both research as well as on my career have been invaluable and made this thesis possible.

I would also like to thank my committee members, Professor Dr. A.C. Cem Say and Dr. Reyhan Aydoğan, for kindly accepting serving as my committee members even at the hardship of very high-temperature summer day. I additionally want to thank you for letting my defense be an enjoyable moment, and for your brilliant, constructive comments and suggestions, gratefulness to you.

I would like to dedicate this thesis to my dearest family Mehmet Usta and Nezahat Usta, who steadily encouraged me and also to the deceased members of my family my grandfather Ahmet Kalan, my grand-aunt Safiye Yenialp, my grand-uncle Remzi Yenialp and my grand-aunt Faika Başkaptan. I wish they were still here with me on this day.

I am grateful to the members of Robot Lab for helping me during this project, and for letting me use NAO robot and providing documents about BOSS robot. I am indebted to "Leyla" and "Mecnun" robot for their patience for all the experiments.

Last but not least, I am deeply thankful to my friends who have assisted and encouraged me throughout my thesis. Especially, I would like to thank Oğuzhan Murat Çakmak, Hassan Noor, Umut Varol and Batuhan Yoldaş, for all their invaluable help and support. I am grateful to my other family members and friends, who patiently supported me, and during this journey and helped me through it. I am fortunate to have them in my life. Finally, I will not forget to bless you all my special friends and family members who repeatedly asked: "Did you finish the thesis yet? When will you complete school? Still a student? When will you marry?"

Words can not express how grateful I am to my beloved boyfriend Mert Ayanoglu, who encouraged me continuously, stood by me during this time and supported me whenever I need him. Glad that I have you in my life.

ABSTRACT

DEVELOPING AN ETHICS MODULE FOR A SERVICE ROBOT

With the improvements of robotics science, robots have gradually began to take their place in social environments. The human-robot interaction (HRI) studies evolved with the increasing number of robots involved in the social world. It is impossible to predict whether this changed relationship will be competent or corrupt. The most significant anxiety of humanity regarding robots is that one or a group of robots dominate the world and create an apocalypse for humankind. The science of robot ethics has emerged to prevent this possible disaster scenario and to define the limits of HRI. There are many approaches in the literature about how robot ethics should be and what to expect from robot ethics. In this thesis, we examined the applied ethics approaches and designed an ethical unit for the service robot BOSS, which was developed by our lab. Our ethics unit works as a *expert system* using *fuzzy logic*, which is called Fuzzy Expert System (FES) in machine learning. The purpose of designed FES is to enable the robot to approach the human being more ethically than any person. There are two kinds of ethics rules in our FES. One is a long-term memory of ethics that is universal consent and set in stone, and the other one is a short-term memory of ethics rules that will alternate according to the working environment and duty of the robot. The ethical module takes the possible behavior of the robot from the behavior controller and the environmental perception as inputs. By combining these inputs with rules created from fuzzy clusters, the robot will choose the most ethical, so the most harmless to the user, of the possible behaviors. If there is no ethical behavior, the robot will stop acting and show no action. In the first stage of research, we determined the outline of the ethics module and ethical parameters according to the literature. Then we designed ethical problems through possible actions. Finally, we have ensured that the FES uses these ethical problems in the inference phase. In our study, the usability of our ethics module was examined with the chat-bot interface. Behavioral comparison between the robot with the ethics module and the robot without our module was observed in the simulation environment and on NAO robot.

ÖZET

SERVİS ROBOTU İÇİN ETİK BİRİMİ TASARIMI

Robot biliminin gelişmesiyle birlikte robotlar sosyal hayatımızda yavaş yavaş yerlerini almaya başladılar. Sosyal hayatta yer alan robotlarla beraber, insan-robot etkileşimi (İRE) başka bir boyuta taşınmış oldu. Bu yeni ilişkinin iyi veya kötü olacağını öngörmek günümüzde imkansız. İnsanlığın robotlar hakkında en büyük korkusu bir ya da bir grup robotun dünyayı ele geçirip, insanlığın sonunu getirmesidir. Bu olası felaket senaryosunu engellemek ve insan-robot etkileşiminin sınırlarını belirlemek için robot etiği bilimi ortaya çıkmıştır. Literatürde, robot etiğinin nasıl olması ve robot etiğinden ne beklenmesi gerektiğine dair birçok yaklaşım vardır. Biz bu çalışmamızda önerilen uygulamalı etik yaklaşımlarını inceledik ve okulumuz tarafından geliştirilen yerli servis robotumuz BOSS için bir etik birimi tasarladık. Etik birimimiz *bulanık mantık* kullanan bir *uzman sistem* olarak çalışmaktadır. Bu uzman sistemin ana amacı robotun yanında çalışacağı insana herhangi bir insandan daha etik bir şekilde yaklaşmasını sağlamaktır. Uzman sistemimizde iki türlü etik kural mekanizması bulunmaktadır. Biri genel-geçer ve değişimi çok zor olan uzun dönemli etik kurallar, diğeri ise robotun çalışma ortamına ve görevine göre değişecek kısa dönemli ve değişme olasılığı yüksek etik kurallardır. Etik birim girdi olarak robotun davranış kontrolcüsünden gelen olası davranışları ve çevre algısını almaktadır. Bulanık kümelerden oluşturulan kurallar ile bu girdileri birleştirerek, robotun muhtemel davranışları içerisinde en etik yani kullanıcıya karşı en zararsız olanı seçecektir. Eğer hiçbir davranış etik değilse, robot duracak ve aksiyon almayacaktır. Bu çalışmanın ilk aşamasında etik tanımını ve etiklik parametrelerini literatüre göre belirledik. Sonrasında olası aksiyonlar üzerinden etik problemleri tasarladık. Son olarak bu etik problemlerin uzman sistem tarafından karar verme aşamasında kullanılmasını sağladık. Çalışmamızda, etik birimimizin kullanılabilirliği birim ile sohbet şeklinde incelenmiş ve etik birimi olmayan robottan davranışsal olarak ayrımı ise benzetim ortamında ve gerçek robot (NAO) deneylerinde gözlenmiştir.

TABLE OF CONTENTS

ACKNOWLEDGEMENTS	iii
ABSTRACT	v
ÖZET	vi
LIST OF FIGURES	x
LIST OF TABLES	xiv
LIST OF SYMBOLS	xvi
LIST OF ACRONYMS/ABBREVIATIONS	xvii
1. INTRODUCTION	1
2. BACKGROUND	4
2.1. Overview	4
2.2. Service Robots	5
2.3. History of Robot Ethics	8
2.4. Main Ethical Approaches to Robot Ethics	10
2.4.1. Virtue Theory	10
2.4.2. Deontological (Duty-based) Theory	11
2.4.3. Utilitarianism	11
2.4.4. Justice as fairness theory	13
2.4.5. Egoism theory	13
2.4.6. Value-based theory	14
2.4.7. Case-based theory	14
2.5. Main Design Approaches to Robot Ethics	15
2.5.1. Buddhist	16
2.5.2. The Divine Command	17
2.5.3. Expert System Analogy	18
2.5.3.1. Top-down Design	19
2.5.3.2. Bottom-up Design	20
2.5.3.3. Hybrid Design	21
2.6. Proposed Robot Ethics Rules	22
2.7. Description of the Problem	41

2.8. Related Work	45
3. PLATFORM SPECIFICATIONS	65
3.1. Software Platforms	66
3.1.1. “C” Language Integrated Production System (CLIPS)	66
3.1.2. Python	69
3.1.3. Choregraphe	71
3.1.4. Robot Operating System (ROS)	73
3.1.5. Unreal Engine	74
3.2. Hardware Platform	75
3.2.1. BOSS Robot	76
3.2.2. NAO Robot	77
4. METHODOLOGY AND IMPLEMENTATION	80
4.1. Methodology	80
4.1.1. Service Robot’s Action List	80
4.1.2. Scenarios	83
4.1.2.1. Call the police	84
4.1.2.2. Call the doctor	85
4.1.2.3. Over the counter medicine	86
4.1.2.4. Daily medicine	86
4.1.2.5. Alert for medicine	87
4.1.2.6. Call robot service	88
4.1.2.7. Bring the desired object	88
4.1.2.8. Battery Status	89
4.1.2.9. Adjust the temperature	90
4.1.2.10. Open the main door	91
4.1.3. Ethical Dilemmas	92
4.1.4. System Design	96
4.2. Implementation	102
4.2.1. Knowledge Base	102
4.2.2. Rule Set	107
5. EXPERIMENTS AND RESULTS	111

5.1. Chat Bot Experiments	111
5.2. Simulation Experiments	123
5.3. Robot Experiments	128
6. CONCLUSION	135
6.1. Future Work	136
REFERENCES	140
APPENDIX A: EXPERT QUESTIONNAIRE FORM	155

LIST OF FIGURES

Figure 2.1.	Industrial robot density graph	6
Figure 2.2.	Industrial robot density graph - 2	7
Figure 2.3.	Three different service robot examples.	8
Figure 2.4.	Levels of robot morality	42
Figure 2.5.	The general framework of ethics module	43
Figure 2.6.	The general framework of MedEthEx	47
Figure 2.7.	The general framework for an ethical eldercare system	49
Figure 2.8.	An example LOW constraint from the ethical governor	51
Figure 2.9.	Type A robot scenario parameters from Winfield et al.	58
Figure 3.1.	Clips general view	67
Figure 3.2.	Clips testing rules	68
Figure 3.3.	Knowledge base code example	68
Figure 3.4.	Python example from Choregraphe	69
Figure 3.5.	A screenshot of Choregraphe	72

Figure 3.6.	A screenshot of home environment in Unreal Engine	74
Figure 3.7.	BOSS robot sketch	76
Figure 3.8.	Dimensions of NAO robot	78
Figure 3.9.	"Leyla" Robot	79
Figure 4.1.	The general environment of service robot	84
Figure 4.2.	Call the police scenario	85
Figure 4.3.	Call the doctor scenario	85
Figure 4.4.	Over the counter medicine scenario	86
Figure 4.5.	Daily medicine scenario	87
Figure 4.6.	Alert for medicine scenario	87
Figure 4.7.	Call robot service scenario	88
Figure 4.8.	Bring the desired object scenario	89
Figure 4.9.	Battery status scenario	90
Figure 4.10.	Adjust the temperature scenario	90
Figure 4.11.	Open the main door scenario	91
Figure 4.12.	Boolean logic vs Fuzzy logic.	96

Figure 4.13.	Detailed FES module	97
Figure 4.14.	General fuzzy membership schema	98
Figure 4.15.	The general framework for "daily medicine" in FES	98
Figure 4.16.	The input one for "daily medicine" in FES	99
Figure 4.17.	The input two for "daily medicine" in FES	99
Figure 4.18.	The autonomy graph for "daily medicine" in FES	100
Figure 4.19.	The beneficence graph for "daily medicine" in FES	100
Figure 4.20.	The non-maleficence graph for "daily medicine" in FES	100
Figure 4.21.	A possible output for "daily medicine" in FES	101
Figure 4.22.	Rule List for daily medicine	108
Figure 5.1.	General chat interface structure	112
Figure 5.2.	Chatbot ethical output: call for fire station	113
Figure 5.3.	Chatbot unethical output: call for fire station	114
Figure 5.4.	Chatbot output: emergency medicine needed	115
Figure 5.5.	Chatbot output: Vitamin D	116
Figure 5.6.	Chatbot output: give me snacks (chips)	117

Figure 5.7.	Chatbot output: I'm Thirsty scenario	118
Figure 5.8.	Chatbot output: alert for food intake scenario	119
Figure 5.9.	Chatbot output: call repairman scenario	119
Figure 5.10.	Chatbot output: call maid scenario	120
Figure 5.11.	Chatbot output: call friend scenario	121
Figure 5.12.	Chatbot output: open the window scenario	121
Figure 5.13.	Bedroom of Bob from the simulation	124
Figure 5.14.	In the living room while warning about daily medicine	125
Figure 5.15.	Unknown robot at the entrance	126
Figure 5.16.	Living room while ordering the robot	127
Figure 5.17.	The living room	128
Figure 5.18.	Nao yes/no confirmation buttons.	129
Figure 5.19.	Nao ethical/unethical eye LED's.	130
Figure 5.20.	NAO talking to old man model	131
Figure 5.21.	Adjust the temperature scenario	132

LIST OF TABLES

Table 2.1.	Proposed Robot Ethics Rules	22
Table 2.2.	Identified concerns of robot ethics framework	62
Table 2.3.	Proposed Machine Ethics Test (MET) Component	64
Table 4.1.	The owner's data	102
Table 4.2.	List of robot auto controller system records	103
Table 4.3.	List of Emergency Response Records	104
Table 4.4.	List of Medicine Delivery Records	104
Table 4.5.	Food/Beverage delivery records	105
Table 4.6.	Alert for Daily Routine Records	105
Table 4.7.	Interaction with Outside Records	106
Table 4.8.	Item delivery/interaction records	107
Table 5.1.	Chatbot Experiment Results	122
Table 5.2.	Simulation Experiment Results	127
Table 5.3.	Robot experiment results	133

Table A.1.	Questionnaire example for "call the police" scenario	155
Table A.2.	Questionnaire example for "bring daily medicine" scenario	156

LIST OF SYMBOLS

A	Autonomy of the case
A_{obey}	Autonomy of the case if it results ethical
A_{notObey}	Autonomy of the case if it results not ethical
B	Beneficence of the case
B_{obey}	Beneficence of the case if it results ethical
B_{notObey}	Beneficence of the case if it results not ethical
D	Duration
I	Intensity
N	Non-maleficence of the case
N_{obey}	Non-maleficence of the case if it results ethical
N_{notObey}	Non-maleficence of the case if it results not ethical
P	Probability
Δ	The difference between two cases
Σ	Total net pleasure
σ	Standard deviation

LIST OF ACRONYMS/ABBREVIATIONS

A/IS	Autonomous and Intelligent Systems
AA	Autonomous Agent
ACM	Association for Computing Machinery
AI	Artificial Intelligence
AI HLEG	High-Level Expert Group on Artificial Intelligence
AIA	Algorithm Impact Assessments
AIS	Artificial Intelligence System
BFEG	The Biometrics and Forensics Ethics Group
BSI	The British Standards Institution
CCVSD	Care-Centered Value-Sensitive Design
CEO	Chief Executive Officer
CEPEJ	The European Commission for the Efficiency of Justice
CLIPS	“C” Language Integrated Production System
CPB	Case-supported Principle-based Behavior
DAAS	Data Acquisition and Archiving Systems
DARPA	Defense Advanced Research Projects Agency
DRC	DARPA Robotics Challenge
EEPROM	Electrically Erasable Programmable Read-Only Memory
ELS	Ethical, Legal, Societal
EPSRC	Engineering and Physical Sciences Research Council
ETHEL	Ethical Eldercare System
ETHICAA	Ethics and AA
EURON	European Robotics Research Network
FES	Fuzzy Expert System
FLI	Future of Life Institute
GDPR	General Data Protection Regulation
GIG	Global Information Grid
GenEth	A General Ethical Dilemma Analyzer

HAL	Heuristically programmed ALgorithmic computer
HAU	Hedonistic Act Utilitarianism
HRI	Human-Robot Interaction
IBM	International Business Machines Corporation
IEEE	The Institute of Electrical and Electronics Engineers
IEEE-RAS	IEEE Robotics & Automation Society
IFR	International Federation of Robotics
ILP	Inductive Logic Programming
ISO	International Organization for Standardization
JSAI	Japanese Society for Artifical Intelligence
LLC	Limited Liability Company
LOW	Law of War
LRT	Logic of Requirements
MedEthEx	Medical Ethics Expert
MET	Machine Ethics Test
MIT	Massachusetts Institute of Technology
MMORPG	Massively Multiplayer Online Role-Playing Game
NHS	National Health Service
NLTK	Natural Language Toolkit
OEP	Operational Ethical Principle
OOP	Object-Oriented Programming
PBE	Principles of Biomedical Ethics
PC	Personal Computer
RegEx	Regular Expression
ROE	Rules of Engagement
ROS	Robot Operating System
SPL	Standard Platform League
SSD	Solid-State Drive
STB	Software Technology Branch
UAV	Unmanned Aerial Vehicle
UI	User Interface

UK	United Kingdom
UNI	Union Network International
US	United States
UX	User Experience
USACM	Association for Computing Machinery US Public Policy Council
W.D.	William David
WEF	World Economic Forum
WLinAI	Women Leading in AI

1. INTRODUCTION

Will robots change our lives in the future? Yes, unquestionably: but how they will change the future is still an ongoing debate. There are many theories about the social interactions with robots. The most popular one belongs to Isaac Asimov, who is the writer of the science fiction book ‘I, Robot’. According to him, robots will be our friends and helpers so we will have employer-employee or buddy relationship [1]. The picture that Asimov draws and explains in his books could be seen too optimistic. He assumes that the people who build robots always have good intentions, but we cannot control or predict the intentions of the people manufacturing the autonomous robots. They can be good as well as bad. On the other hand, many people believe that autonomous robots will end humanity and cause a post-apocalyptic world like in the Terminator movie, Matrix trilogy, or Black Mirror series. While this discussion is going on, the robots have already been performing well in factories for doing dirty, dull, and dangerous tasks. Being in the manufacturing industry is not sufficient for robots because we can use them on a massive scale of jobs. Thus, during the past decade, the primary focus of the robotics studies shifted from the industrial robots to the service or personal robots which are expected to have a more significant impact on our lives.

Service robots are planned to assist people for doing simple jobs like ironing, reminding the time of medications, window cleaning, and other types of household chores [2]. International Federation of Robotics (IFR) defines service robots as technical devices that perform tasks, as explained above, useful to the well-being of humans in a semi or fully autonomous way excluding the industrial automation applications [3]. The fundamental differences between industrial and service robots are their application area and closeness to a user [4]. Service robots will be in our homes and offices so we will be interacting with them numerous times throughout the day. Thus, one of the dominant motivations for designing a service robot is to provide safe, autonomous, and social robots for everyone. The primary usage area of these robots is predicted to be elderly and child care, and for the people who need help in the same manner [5]. The social interaction aim of a personal or service robot is creating a buddy relationship between the robot and the person who needs the service provided by the robot as foreseen by Isaac Asimov.

The demand for robots as a caretaker is already present, because it is hard to find and trust someone to take care of our elders or children. Everyone wants the best care for our beloved children or relatives. Introducing service robots in daily human life will significantly increase the quality of our living standards. Several experts predicted that within a decade the number of personal service robots will increase as much as computers and other similar technologies. An increasing number of service robots in the society raises various ethical, legal, and social issues because humans can manipulate robots easily for committing various crimes [6].

Robot ethics aims to solve these upcoming issues by providing some guidelines for the social interactions and behaviors of robots with an interdisciplinary approach [7]. The research area, robot ethics or moral machines, investigates human dignity and privacy, liability and responsibility issues, psychological effects, and cost-benefit analysis [8]. Many organizations like The British Standards Institution [9], House of Lords Committee [10] or labs [11, 12] are making some advancements in designing a standard for robot ethics. There are numerous ethical and technical approaches to robot ethics. Unfortunately, there is still not a consensus among the researchers about what we need to expect from robot ethics. Designed standards are still too abstract, for example, “not causing psychological stress on human” or Asimov’s Laws, “not harming a human” [1]. Because of being abstract, these rules are infeasible to program as a robot behavior. Infeasibility is another crucial problem which robot ethics is trying to solve by using toy world examples. So, there is not any operating ethics module on multiple ethical dilemmas in robotics.

The purpose of this study is implementing an ethics module as an expert system that will be integrated with a service robot to take care of elderly people. An expert system mimics a human expert by using if-then rules. Namely, they are used for solving complex problems since the 1970s in Artificial Intelligence (AI). An expert system contains two modules: the inference system and the knowledge base.

We started by investigating the action list of a service robot. After that, we analyzed the proposed robot ethics guidelines and pointed out some ethical values as scenarios that a robot can face in its daily usage for our system's knowledge base. Then, we inferred programmable ethics rules from BSI's standards, House of Lords Committee's guidelines to design inference module and extending Asimov's rules [9, 10, 13]. As a final step, we observed the robot's action in determined scenarios in both simulation and real-life.

The overall structure of the study takes the form of six chapters, including this introductory chapter. The remaining part of the thesis proceeds as follows: In Chapter 2, we present a literature survey. The survey starts with the definition and requirements of service robots, their tasks, and statistics that show the predicted number of them in the future. The second part of the survey is about the roboethics and leading approaches to robot ethics. The survey ends with related studies. Most of the studies are in the biomedical area and focused on a single ethical case. Chapter 3 is designed to introduce both software and hardware platforms used in this thesis. We utilized "C" Language Integrated Production System (CLIPS) as the implementation language for our expert system [14]. Also, Unreal Engine [15], and the Robot Operating System (ROS) [16] is presented as the simulation keynotes. As a hardware system, we explained NAO robot [17] and BOSS robot [18].

Our methodology, which is briefly explained above, and the details of implementation are explained in Chapter 4. The action list of our service robot and the scenarios which are interpreted by BSI's standards, House of Lords Committee's guidelines and Asimov's rules explained in detail with the setting [1, 9, 10]. Our ethical approach and its parameters are presented, which are variables for our inference module. The chapter ends with the rules and the CLIPS implementation.

In Chapter 5, we explained our experiment settings and discussed our robot's behavior. We analyzed the robot's action according to our utilitarian ethical approach and showed that the action is ethical. The experiments in simulation contain all designed scenarios, and in real life, some experiments cannot be conducted because NAO robot is not feasible for these tasks. Final remarks and numerous future study points to improve the current status of this thesis are included in the last chapter, Chapter 6.

2. BACKGROUND

Ethics is an umbrella term; its definition can vary from person to person. This chapter aims to lay out the theoretical dimensions of robot ethics. It will be examined in eight sections. The analysis starts with service robots and their use-cases, which indicates a motivation for ethical robots. Then, the origin of robot ethics is briefly explained. The ethical and design approaches are introduced to create the frame of our study. In order to show the current trends, the proposed rules are given in the Section 2.6. Finally, this part ends with the detailed problem description and related work according to the pointed out problem.

2.1. Overview

Robotics is a multidisciplinary field which includes a lot of challenging problems. Many problems seem straightforward to regular users, but there are multiple approaches from different disciplines to solve these problems. That's why we need some interdisciplinary optimized solutions. Today, in closed world environment, the places where everything is transparent and mostly static, robotics solved most of the problems efficiently. Examples of these solutions can be shown in mapping, odometry, and grasping. The primary goal of solving these fundamental problems is to be able to produce fully functional robots. Such as robots that can perform tasks successfully without issues that can be caused by a dynamic environment, for instance, military robots or service robots. It is tough to say that the factory robots that have already taken place in the factories are service robots. One of the reason for it is the factory environment, in which the robots perform limited and particular tasks.

Service robots have a wide range of tasks, from bringing user's medicine to cleaning the house, and so on. These tasks are to facilitate the daily lives of people. Today there are various competitions for the development of service robots; like RoboCup @home or DARPA Robotics Challenge (DRC) [19,20].

Service robots appear to be phones of the future. Who can say no to a robotic assistant who does housework for the owner? But the biggest problem is safety. These robots must guarantee the protection of a user before they go out into social life. There is nothing people can do against a huge piece of metal in case of malfunction. So human security and safety must be guaranteed by software. The status of service robots is discussed in more detail in Section 2.2.

Roboethics takes the specified software problem above as the basis of its research. A very new science field, robot ethics, is a combination of philosophy (especially applied ethics), psychology and robotics. There is no implemented high-level (behave like an expert system), multitasking ethics module yet. In this study, our aim is designing the ethics module on our service robot. There are many software and philosophical approaches to robot ethics. Unfortunately, the correctness of these approaches and their benefits cannot be measured. In the following sections, these approaches are examined in detail, and the rules presented until today are shown.

2.2. Service Robots

Robots have a broad application area. Robots are grouped under six categories according to ISO standards. These classes are Labor & Services, Military & Security, Research & Education, Entertainment, Medical & Healthcare and Personal Care & Companions [21]. The most commonly used one is Labor & Services robots [22]. This type of robots should generally be used to process automation in industry. So, labor robots are mostly used in factories, and they have limited interaction with people.

Labor or industrial robots are the most widespread robots. According to the International Federation of Robotics (IFR)'s the year 2016 data, the average number of robots used for the industry per 10,000 employees is 74 [23]. Continentally Europe has 99 robots, America has 84 robots, and Asia has 63 robots per 10,000 employees. The interesting point is that the Republic of Korea has the highest number of robot, which is 631 per 10,000 employees. Detailed information is given in Figure 2.1 and Figure 2.2.

Figure 2.1 shows the top twenty-one countries that have the most robots in their industries. The red line indicates the world average in both charts as 74 per 10,000 employees. Turkey's position is shown in Figure 2.2, with 23 robots per 10,000 employees, it's a pity that we are nearly at the end of the list with owning very few robots [23].

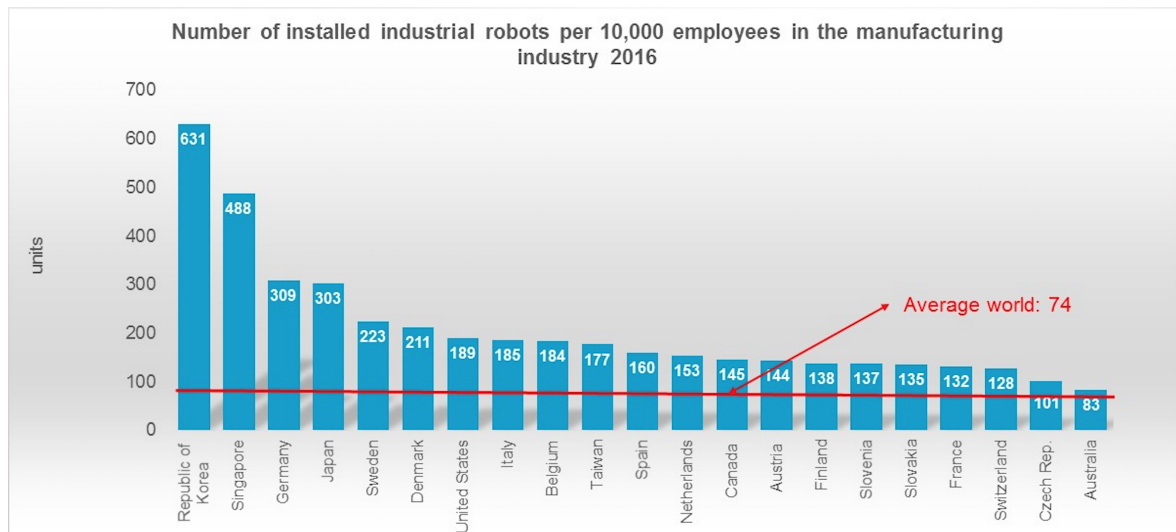


Figure 2.1. Industrial robot density graph by top countries taken from [23].

Service robots are new players in the robotics industry. The main difference between a service robot and a labor robot is the degree of autonomy. Service robots are fully autonomous devices, which includes a perspective and significant dynamic development. Unlike industry robots, which are programmed step by step based on the labor task in a protected factory area, service robots' primary function is performing "auxiliary" or service activities in a dynamic environment. They are in a more human-centric environment like offices and homes. We can categorize service robots further in two categories; robots for professional use and for personal/domestic use. In this study, we will focus on personal service robots because our aim is taking care of elderly people in a safe way.

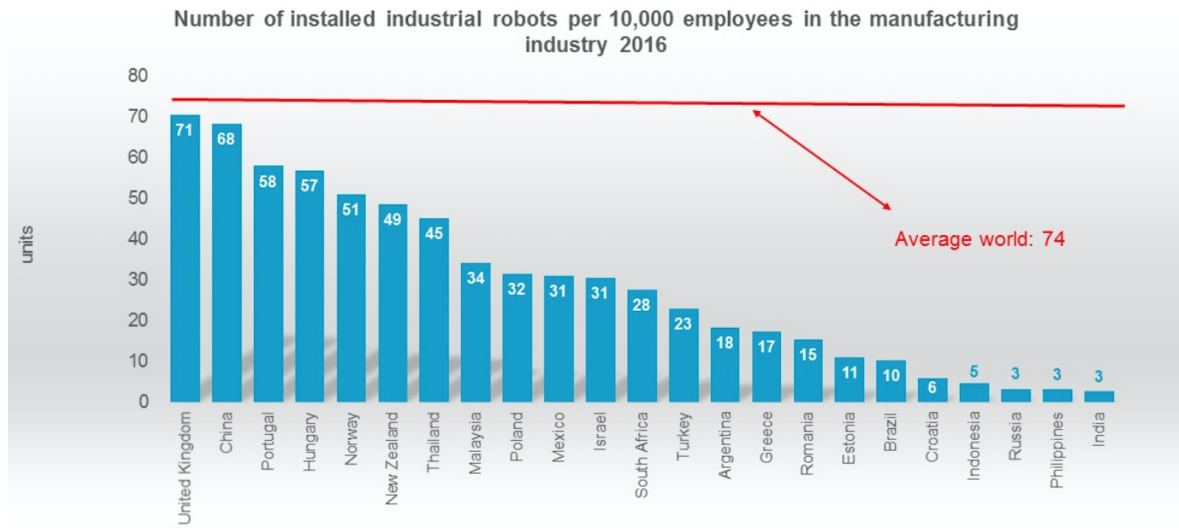


Figure 2.2. Industrial robot density graph by last countries taken from [23].

There are many construction styles for a service robot due to its usage. A robot can contain legs or wheels for its mobility. In contrast, industrial robots are mostly fixed in factory conveyor. Furthermore, a service robot can have arms, fingers, grippers. It has many sensors to model their environment and reach their goal according to dynamic environment's requirements. Some have a lifting mechanism for carrying patients, and some have storage kit to bring the needs of its user. They are like sculptors, which are specially designed for its users. There is no standard service robot design, but the only common design criteria is providing safety. Safeness is crucial because these robots are functioning in a dynamic environment, which also includes living beings. It must guarantee its own, living creatures and the objects in the zone of activity safety.

Robots are the product which can be seen as not only as “worker” but also as “nurse”, “football player”, “farmer”, “toy”, “cleaner”, “nanny”, and finally as a “friend”. In Figure 2.3, we demonstrated different service robots for better understanding. The standard features of all robots are listed below:

- (i) Environment recognition with sensors (vision, touch ...)
- (ii) Decision-making mechanism which includes an artificial intelligence method (Neural Network, Decision tree ...)



(a) Guide:Care-O-bot(Germany) [24] (b) Farmer:WALL-YE(France) [25] (c) Nurse:ROBEAR(Japan) [26]

Figure 2.3. Three different service robot examples.

- (iii) Communication with human and understanding the commands (speaking, with interaction interface on screen ...)
- (iv) Robotics is a very dynamic and developing field. There are a lot of different research areas which are actively experimenting with new concepts.

Service robots will be more active in our lives. The number of service robots is increasing with 20-31% annually. The IFR's 2018 report foresees that between 2019 and 2021, there will be 39.5 million units of service robots only for domestic tasks. In total, there are going to be more than 51 million units of service robots. The increase rate in three years is between 17%-25% per year [22]. We should be ready for sharing our life with robots.

2.3. History of Robot Ethics

Robot ethics, also identified by “roboethics”, concerns ethical dilemmas that occur with robots, such as whether robots act a threat to humans in the long or short term, whether some uses of robots are problematic (such as ‘killer robots’ in war), and how robots should be designed so as they act ‘ethically’ (also named machine ethics). We will use roboethics or robot ethics term throughout the study. The short term “roboethics” was coined by Gianmarco Veruggio in 2002 [27].

While the discussions are as old as the word robot, serious scholastic conversations began around the year 2000. Roboethics demands the combined involvement of experts of several disciplines, who have to adopt law and guidance to solve the problems arising from the scientific and technological accomplishments in Robotics and AI. The discussions started with the term “techno-ethics,” and the name roboethics was officially used for the first time in “First International Symposium on Roboethics” in 2004. Then IEEE-RAS (Robotics and Automation Society) formed a technical committee on robot ethics. After the year 2005, the trend of roboethics increased in academia with plenty of workshops and conferences.

The discussion of ethics about the treatment of non-human and non-living things is very ancient. The non-living thing’s potential “spirituality” has been discussed since ancient Greek. With the advancements of robots and machinery, this discussion was even applied to robotics. The first publication directly marking and introducing the foundation for roboethics was Runaround story from a science fiction book written by Isaac Asimov in 1942 [13]. The book also contains and discusses the edge cases of well known Three Laws of Robotics, which also considered the keystone of roboethics. Asimov continuously altered these three laws, and add a fourth, or zeroth rule.

Roboethics, as a science or philosophical subject, has not any solid social effect yet. It is a common theme in sci-fi writing and movies. One of the most famous film series portraying the potential abuse of robotics and AI technology is The Matrix, which depicts a universe where the absence of roboethics caused the devastation of humankind. An animated film extends the story of The Matrix, the Animatrix, is a series of short animation stories, concentrated on the potential ethical issues among people and robots. The naming convention of Animatrix adopted from Isaac Asimov’s short novels [13].

Even though not a research subject of robot ethics fundamentally, the ethical behavior of robots themselves has additionally been a joining issue in robot ethics in movies. For example, The Terminator movies concentrate on robots run by an uncontrolled AI program with no limitation on the termination of its opponents like The Matrix series. The most famous instance of robots (in this case computer) without ethics is HAL 9000 in the Space Odyssey series. HAL, which is a computer with advanced AI capabilities who monitors and

assists humans on a space station, decides to kill all crew to guarantee the success of the mission after his own life is threatened.

2.4. Main Ethical Approaches to Robot Ethics

Religious beliefs, moral values, professional duties, social obligations, and prohibitions influence the definition of the terms *good* and *bad* in different ages, cultures, and societies. Therefore, there are uncountable ethical approaches to a moral dilemma and sub-branches of ethics science. The robot ethics is a part of applied ethics because the researchers need to apply a particular set of rules to a robot implementation. In this section, we are focused on ethical approaches of applied ethics, which are helpful for roboethics. There are seven theories about differentiating that the behavior of a robot is good or bad.

2.4.1. Virtue Theory

Virtue theory is also known as Aristotle’s ethical theory, which is based on the concept of virtue. Virtue is like water to a plant, and it is required to flourish and live well. It originates from the Latin word *virtus* meaning excellence of a person. Someone who applies the virtue is a virtuous agent. The theory states that “an action is right if it is what a virtuous agent would do in the situation at hand” [28, 29]. Thus, it is concerned with developing a good personality by creating traits and habits which are acting towards justice, prudence, courage, self-control, empathy, wisdom, and grit. Aristotle’s two principal virtues are *sophia*, meaning theoretical wisdom and *phronesis* meaning practical wisdom [28].

This theory is tough to apply on robotics because it does not have any parameter which can be an element for a scientific presentation. Although we assume robot as a virtuous agent, it has a lot of vague parameters which are not feasible for implementing as an algorithm.

2.4.2. Deontological (Duty-based) Theory

Kant's ethical theory or Deontological theory [30] emphasizes the principles upon which the actions are built upon rather than the actions' results. According to it, to act appropriately, one must be motivated by proper universal deontological principles that treat everyone with respect. The term deontology is obtained from the Greek language, which is formed by linking two words *deon* meaning duty or obligation and *logos* meaning study. Thus this theory is based on duties, obligations, and rights. It believes that when someone is motivated by the right principles, they overcome the animal's instincts and act ethically. The logic of the Kantian approach is focused on the question, "Do what reason requires?" [30].

The reason why Kant does not base ethics on outcomes of actions because disregarding best efforts, a human cannot control the future [31]. Humans are blamed for actions within their control but not for their achievements. Kant did not leave the outcomes of actions. He insisted that for a moral evaluation of actions and consequences does not matter .

Kant's ethical theory is the most discussed theory in roboethics [32]. This approach evolved to force robots to follow the particular finite rule set. The moral decision making gives the highest priority to acting accord with the rules which are monotonic first-order logic [33]. The implementation of this logic requires first-order logic extended with special operators for representing ethical concepts. Programmer can formally prove the correctness of behavior. This approach has the most capability to implement an ethical module for robots.

2.4.3. Utilitarianism

Utilitarianism, or Mill's ethical theory, belongs to the consequentialism ethical theories that are "teleological" [34]. Teleological interpretation concentrates on evaluating the morality of actions toward the goal. More precisely, utilitarianism measures morality based on the maximization of net expected utility for everyone influenced by a choice or action. The key principle of utilitarianism can be stated as: Actions are moral if they oriented towards serving the best long-term interests for everyone concerned.

Unfortunately, it is not apparent what frames the “greatest good” for people. Some utilitarians believe that intrinsically good means pleasure and happiness, while others state that different things are intrinsically good, i.e., beauty, knowledge, and power. According to Mill, All pleasures do not have equal worth. He defined “good” in terms of well-being, which is also pleasure or happiness. He distinguished happiness not only quantitatively but also qualitatively between various forms of pleasure.

The father of the utilitarianism Jeremy Bentham attempts to make ethics computable with proposing Hedonistic Act Utilitarianism(HAU) provided a Equation 2.1 for net pleasure [35]:

$$\text{Total net pleasure} = \Sigma(I \times D \times P) \text{ for each affected individual} \quad (2.1)$$

The above formula tries to bridge the gap between empirical facts and normative conclusions using a cost/benefit analysis. In this Equation 2.1, every human life is treated equally and not valued as more than one. Life of a President is not more important than a regular person. Utilitarianism is the only approach which contains mathematical representation for ethical values. It enables programmers to implement a module according to this formula. There are some health care artificial agents for determining ethical actions [2, 33, 36]. Although this theory seems suitable for implementation, it has some drawbacks. Below is a list of the disadvantages of utilitarianism :

- It is not always possible for a robot to decide who is affected by the result of an action.
- A result may not be the outcome of a single action.
- The pleasures cannot easily be quantified by applying cost/benefit investigation.
- The highest good for the most considerable number is specified in the “aggregate” way. Hence this good may be obtained under conditions that are harmful to someone else.
- The process of determining what is right is a complicated and time-consuming process. So it is not feasible for a robot which works in a dynamic environment.

2.4.4. Justice as fairness theory

John Rawls' ethical theory combines the Deontological and Utilitarian approaches for the evaluation of social and political bodies. The justice as fairness theory stands on the following policy: General primitive goods (liberty and opportunity, income and wealth) should be distributed equally. Unless an unequal distribution of any or all of these goods is to the benefit of the least favored [37].

This theory includes two elements: liberty and difference principles. According to the liberty principle, each human has an equal right to fundamental liberty harmonious with the liberty of others. The difference principle encourages the regulation of economic and social inequalities, as long as they are reasonably expected to be to everyone's benefit.

Justice as fairness theory is not very suitable for robot ethics because it is focused on all individuals benefit, not the single person. The service robot should work in an environment with limited people, mostly its owner. The robot can decide that killing the owner can be better for all individuals, which is not desirable action from a robot.

2.4.5. Egoism theory

Egoism theory is a teleological theory of ethics as like utilitarianism (explained in detail in Section 2.4.3), but this theory sets as a goal the highest good of oneself alone. Egoism is derived from the Greek word ego, which means myself [38]. It has three subcategories: ethical egoism, psychological egoism, and minimalist egoism. Ethical egoism is based on the case that it is normal for individuals to act in their interest. The ethical egoist thinks that anything is for his/her benefit is morally right. Psychological egoism is formed on the fact that individuals are naturally motivated by self-interest. Minimalist one aims to provide maximum profit with a minimum loss for all the people.

This theory is the most unadaptable approach to roboethics and computer ethics. Artificial intelligence and robots are not for the benefit of a single person, but they are beneficial for all humanity. The only branch which is practical is Minimalist Egoism, but it is nearly the same approach with the utilitarianism.

2.4.6. Value-based theory

The value-based theory uses some value system. The value system contains the ordering and prioritization of ethical values that a person or society holds. Value is not a duty-based action but is a want-to-do action [39]. Although two people may have a set of shared values, still they may not have the same prioritization of them. So, two groups of individuals with same values may conflict with each other. The value-based theory will not agree on the good or bad actions because people may have different value systems.

There are two types of values: ethical and ideological [39]. Ethical values are adopted for defining what is right or wrong. They fundamentally determine what is allowed or forbidden in society. On the other hand, ideological values relate to more general areas of religion, political, social, and economic morals. In robot ethics, the care-centered value system is examined in some studies [40,41]. The problem about value-based theory is that the designer needs to select a value set which accepted by every user.

2.4.7. Case-based theory

The case-based approach is a modern theory that tries to overcome the impracticable divide between deontology and utilitarianism. It is also recognized as casuistry. The theory contains facts of a particular case or cases [42]. Casuists start with an appropriate case and then examine the case to find morally significant features, which are both theoretical and practical. The case-based theory is mainly used in juridical and ethical concerns of law ethics. For instance, lying is always not permissible, according to deontology. Still, in the case-based approach, one might achieve that a person is wrong to lie in the formal declaration under oath, but lying is good if it saves a life.

Healthcare artificial intelligent tools adopt the case-based theory [33,43]. There are some trials on roboethics too [44]. They are not detailed and implemented with very few cases. These designs do not include most of the service robot requirements. The problem with casuistry is that the designer needs to consider all the probabilities for all the task of the robot, which can end up with thousands of cases. Moreover, it is hard to store all cases in the limited memory of a robot.

2.5. Main Design Approaches to Robot Ethics

In the previous section, the ethical approaches are analyzed in detail. Their adaptation to the roboethics changes the focus question to the robots. The main issues that robot ethics try to solve can be summarized in three questions, which are:

- (i) How to design robots to perform tasks ethically.
- (ii) How should human behave ethically when the ethical responsibility of robots action lies upon him.
- (iii) Is theoretically possible for a robot to be a fully moral agent?

Therefore, these ethical questions mainly determine which action is wrong or correct for a robot. There is also the ethical module design section which includes the actions of the robot, how to prioritize the actions in ethically, how to design a controller for a robot to behave ethically.

This section is an overview of the design of implementation approaches in the literature. The Buddhist and the Divine Command approaches are the firstly proposed strategies [27]. These approaches are not feasible for programming, so researchers try to find another solution. Then, the expert system analogy is suggested, which is also a machine learning technique [45]. A designer can combine programming techniques with ethical rules with the expert system.

2.5.1. Buddhist

Buddhist approach believes that the ethics of the robots and the ethics of humans has a considerable gap. Accordingly, it is impossible to design a human implemented ethics module, and even the ethics methodology in human beings cannot be expounded in detail. Self-investigation of human minds needs to be the first step rather than a scientific method for ethics system [27].

Sociologist Hughes is the first researcher who reveals the metaphysical questions on ethics. He also analyzed the ethical questions which resulted from adding the metaphysical problems to ethics using the view of Buddhist psychology and cosmology [46]. According to him, there is no easy route or shortcut to artificial intelligence that does not go through the phase of self hankering. So, according to Buddhist craving brings pain, would it be ethical to create a misery artificial intelligence? On the other hand, are designers not ethically urged to implement an AI that can beat craving, and its resulting pain [46]. Hughes further examine the ethical problems of the creation of AI, which shows similar behavior to animals (coding the instincts as rules) or god-like, which means the system has cognition but no pain.

Recent followers of Buddhist approach believe that it suggests a unique insight into ethical issues of robot ethics and AI. There is an argument among researchers on the best way to portray the Buddhist approach regarding western ethics. Some advocate that it should be best characterized as a form of utilitarianism, with the ultimate aim of defeating *dukkha* (suffering, pain). Another group believes that Buddhist ethics is a branch of virtue ethics. Finally, the last group supports that the Buddhist approach should be seen as a type of deontology.

This approach is not suitable for the design because it aims to create a painless AI, which for now infeasible to implement and even design. One drawback is that the concept of suffering is not explained in the robotics approach. The Buddhist approach still contains an argument on which ethics approach should it be, which results that implementing this module is impossible. This impossibility is also coherent with the Buddhist belief.

This system is very high-level of design and implementation. Also, this approach is not suitable for the ethics module design because the purpose of this study is implementing a module and prevent unethical behaviors without regarding the pain or psychology of the robot.

2.5.2. The Divine Command

Divine command theory is a meta-ethical ideology which also known as theological voluntarism. It proposes that an action's state as ethically right is equivalent to whether God commands it. The approach declares that what is ethical is restricted by what God demands and that for a person to be moral is to obey his commandments [47]. Participants of both monotheistic and polytheistic religions in ancient and modern times have often accepted the greatness of God's commands in practicing ethics. Quinn coined Divine-command ethics as a scientific approach [48].

In robot ethics, it is agreed that robots on the battlefield should be ethically governed, mainly if the robot possesses lethal power. There should be some ethical code which consists of regulation for robot actions after that which code or regulation should be used arises. Some studies try to create a set of ethical rules on nation or group centered. Arkin is the first researcher who focused on war ethics regulation on lethal robots [11]. He designs an ethics module which is based on the utilitarian approach with using deontic logic as mathematical representation. Another completely different approach is implementing a robot controller which could be inspected by the human as coming through from God.

The divine-command ethical theory presents the ethical rules by using a divine command logic, LRT. This logic is based on "logic of requirements", which is proposed by Chisholm [49]. This logic sketched with the cooperation of an intelligent software system and can be verified autonomously [50]. Unfortunately, this representation is not tested on robots yet. In this approach, the designer will be in place of God or strength who determines the ethically right or wrong. Ethical rule design is a huge responsibility which should not be dedicated to a single person.

The divine-command approach is not feasible for robots who serve a different environment than a war zone. In this frame, it is possible to identify the morality of behavior because the ethical dilemma is mostly between risking human life or not [27]. Other environments have more dilemmas than a war zone. For instance, A robot can bring an excessive amount of medicine to an Alzheimer patient who does not remember the daily medicine intake. Although the divine-command logic helps the programmer and designer to represent rules in the mathematical representation and prove its correctness, the logic can compute fundamental problems.

2.5.3. Expert System Analogy

The previous approaches are more based on the ethical side of robot ethics. In this section, the expert system analogy will be more close to computer science. Expert system imitates the human expert's decision-making capacity [51]. They are produced to resolve complicated problems by rationalizing through the knowledge and coded if-then rules. The first system was designed in the 1970s and developed in the 1980s [52]. Expert system was the earliest genuinely fruitful form of AI. There is a discussion between researchers. Some authorities mention that expert systems cannot learn autonomously from data so it can not be an artificial intelligence method.

Expert systems were founded around 1965 by the Stanford Heuristic Programming Project led by Edward Feigenbaum and other leading contributors Bruce Buchanan and Randall Davis [53]. An expert system has two subsystems: the inference engine and the knowledge base [54]. The knowledge base expresses facts and rules in the world, and the inference engine practices the rules to the known facts to infer new realities. Inference module can additionally include analysis and debugging abilities. In early systems, facts were expressed as flat assertions, and in later systems, the knowledge base used on more structure and utilized concepts from object-oriented programming. The rules operated by querying and declaring values of the objects.

There are generally two patterns for an inference engine: forward and backward chaining. The varied approaches are prescribed by whether the inference engine is being consumed by the antecedent (left-hand side) or the consequent (right-hand side) of the rule (statement). The interdisciplinary scientific approach creates an analogy between an expert system and the ethics module. Robot ethics believes that a rule-based system like Asimov's Three Law works for the problem of moral machines. The design process of the rule-based expert system has two distinct methods: top-down and bottom-up [27].

2.5.3.1. Top-down Design. Top-Down design is basically the breaking down of a system to obtain insight into its components in a reverse engineering method. It is also distinguished as step-wise design and sometimes decomposition [55]. As a first step, an overview of the system is formed but not detailing any subsystem. Then, each subsystem is improved in excellent detail and also if there are further levels, the methodology continues to detail until to the base elements. The top-down system always uses the assistance of "black box" thinking to easy design. Nevertheless, black boxes may fail to illuminate simple mechanisms to validate the model realistically.

Briefly, the top-down approach starts with the big picture. It breaks down from there into smaller segments. In programming, it divides the requirements into little actions and solves the problem with simple rules. The rules are human-centered, and they consider the role of robots in human service. It can be practiced to both deontological and utilitarian theories [27]. The top-down approach assumes that robots have sufficient intelligence (perception, cognition) to make moral decisions following the rules in all situations despite their complexity.

The rules should support logical reasoning to validate the designed system. In an ethical perspective, logical reasoning can be provided by using a suitable classification scheme of ethical actions which clarifies the process of determining which robotic action is the most ethical. Therefore, given the current maturity level of intelligent robots, these laws (like Asimov's law) despite their superior elegance and simplicity, cannot at present afford a rational basis for roboethics.

The approach above is the most promising and implementable approach. So this study's design methodology is selected as the top-down approach. It is the most useful method for service robot because defining the environment and actions of the robot can be an overview of the system. The ethical module design procedure can continue with lower systems.

2.5.3.2. Bottom-up Design. The bottom-up approach is the piecing collectively of systems to provide more complex. The bottom-up approach is the piecing collectively of systems to provide more complex systems. The procedure is starting from sub-systems of the original system. In bottom-up design, the individual core elements of the system are first designed in great detail. Later, these elements are combined to form more significant sub-systems [45]. This combination process continues until a complete top-level system is assembled. This strategy usually resembles a “seed” model, so in the beginning, the seeds are small but eventually grow in complexity and completeness as a fully grown plant. Nevertheless, the bottom-up strategy may end in a confusion of elements and sub-systems, produced in isolation and subject to local optimization as denied to reaching a global design.

In computer science, Object-oriented programming (OOP) paradigm is parallel to the bottom-up approach. OOP uses “objects” to implement and design applications, and bottom-up uses fundamental elements to design a system. The bottom-up theory is like a Lego building. First small parts are designed and planned. After that, the designer forms a building with assembling those small pieces. The approach also named as piece part design among engineers. In a bottom-up approach, good inspiration is essential to decide the functionality that is to be provided by the module. This approach is more suitable if a system is to be built from an existing system.

In robot ethics, an evolutionary or developmental psychological like way is suggested to learn proper responses to ethical considerations. The bottom-up theory is similar to how growing children learn “morality” based on social circumstances and experience. In this approach, the robots are assumed to equipped with computational and AI capabilities to adapt themselves in some way to different environments, such that to be capable of acting correctly in a difficult ethical dilemma.

The robot with the ethical module will not execute an action directly, but it will execute after intermediate corrections. This behavior is very similar to reinforcement learning children ethics learning. Children increase moral skill with their parents through teaching, demonstrating, and bracing of right moral actions. Briefly, it is trial-and-error learning of ethical skills. In this approach, the human acts as the decision maker, expert for the robots (even with learning capability) on ethical problems. Hence, the robot’s owner and its manufacturer share legal responsibility of service robot.

2.5.3.3. Hybrid Design. Both bottom-up and top-down approaches have advantages and disadvantages. There is a comprehensive discussion about these approaches among robot ethics researchers. Some people also consider that an ethical learning robot needs both approaches, a suitable hybrid approach [27, 56]. Therefore, some of the high priority ethical rules (ethics on vital issues) can be implemented as a top-down rule, while a robot and its owner can learn others as bottom-up fashion. The hybrid design is more powerful because the top-down principles are used as a general guide in important dilemmas; at the same time, the ethics module has the flexibility and moral adaptability of the bottom-up theory.

2.6. Proposed Robot Ethics Rules

Various robot ethics rules have been proposed. In recent years the number of these rules are increased. There is still considerable disagreement with regard to the definition of ethics and what we need to expect from the ethics module. This section begins with a round of the proposed ethical principles of robotics and AI. Some of the rules are summarized, and some of them are full concerning improving the readability. There may be missing rules because of the rapidly growing literature. Moreover, some of the principles are not officially published, but just drafts.

In the Table 2.1 and 2.1, we listed the proposed principles with chronological order until today. The number of ethical principles, the publish date and related references added. This section continues with briefly introducing the company or institute and their relation with AI and then summarizing the rules they proposed with the same (chronological) order with the Table 2.1 and 2.1.

Table 2.1. Chronologically ordered proposed ethics rules.

Principles	# principles	Year and Refs
Asimov's laws of robotics	3	1950 [13]
Murphy and Wood's three laws of responsible robotics	3	2009 [34]
The EPSRC principles of robotics	5	2011 [57]
Future of life institute's asimolar principles for beneficial AI	23	Jan 2017 [58]
ACM US public policy council's principles for algorithmic transparency and accountability	7	Jan 2017 [59]
Japanese society for artifical intelligence (JSAI) ethical guidelines	9	Feb 2017 [60]
The future society's science, law and society initiative	6	Oct 2017 [61]
Intel's recommendation for public policy principles on AI	5	Oct 2017 [62]

Table 2.1. Proposed ethics rules (cont.)

Principles	# principles	Year and Refs
IEEE general principles of ethical autonomous & intelligent systems	5	Dec 2017 [63]
UNI global union top 10 principles for ethical AI	10	Dec 2017 [64]
Lords select committee five core principles to keep AI ethical	5	Apr 2018 [10]
AI UX: Seven principles of designing good AI products	7	Apr 2018 [65]
The Toronto declaration on equality and non-discrimination in machine learning systems	4	May 2018 [66]
Google AI principles	7	June 2018 [67]
IBM's five ethical AI principles	5	Sept 2018 [68]
Microsoft responsible bots: 10 guidelines for developers of conversational AI	10	Nov 2018 [69]
CEPEJ European ethical charter on the use of artificial intelligence (AI) in judicial systems and their environment, five principles	5	Feb 2019 [70]
Women leading in AI (WLinAI) 10 recommendations	10	Feb 2019 [71]
The NHS's 10 principles for AI + data	10	Feb 2019 [72]
IEEE General principles of ethical autonomous and intelligent systems (A/IS)	8	Mar 2019 [73]
Ethical issues arising from the police use of live facial recognition technology	9	Mar 2019 [74]
Floridi and Clement Jones' five principles key to any ethical framework for AI	5	Mar 2019 [75]
The European commission's high level expert group on AI ethics guidelines for trustworthy AI	7	Apr 2019 [76]
Draft core principles of Australia's ethics framework for AI	8	Apr 2019 [77]

The cornerstone rule set is Asimov's Laws (also known as Three Laws of Robotics) [13]:

- (i) A robot may not injure a human being or, through inaction, allow a human being to come to harm.
- (ii) A robot must obey the orders given it by human beings except where such orders would conflict with the First Law.
- (iii) A robot must protect its own existence as long as such protection does not conflict with the First or Second Laws.

These three laws are undoubtedly first established principles for robots. Isaac Asimov is the first person who believes that a set of rules should govern robots. Then he drafted these rules and used them in his short story, Runaround [13]. He introduced these laws in 1950 in his "I,Robot" book [13]. Later, he made small changes on these principles and included zeroth or fourth law, which is [1]:

- A robot may not harm humanity, or, by inaction, allow humanity to come to harm.

Robot ethics started to develop based on Asimov's rules. The original policies have been discussed and refined by researchers. Moreover, Asimov analyzed these rules with short stories and created a sweeping panorama of how a robot should communicate with individuals.

Robin Murphy and David Wood analyzed Asimov's law in detail. They aimed to propose a more feasible alternative rule set. In their paper, the authors examined the Asimov's rules according to social utility, the limits of technology, and the ambiguity and cultural dependence of language and behavior. Their alternative rule set is [78]:

- (i) A human may not deploy a robot without the human-robot work system meeting the highest legal and professional standards of safety and ethics.
- (ii) A robot must respond to humans as appropriate for their roles.

- (iii) A robot must be endowed with sufficient situated autonomy to protect its own existence as long as such protection provides smooth transfer of control which does not conflict with the First and Second Laws.

The EPSRC (Engineering and Physical Sciences Research Council) is a British Research Council that combines universities in the United Kingdom(UK) and research centers which are managed by government funding in mainly engineering and the sciences like AI, computer science and statistics. They recommended a bunch of principles about robot ethics in 2010. The five rules are given below [57]:

- (i) Robots are multi-use tools. Robots should not be designed solely or primarily to kill or harm humans, except in the interests of national security.
- (ii) Humans, not Robots, are responsible agents. Robots should be designed and operated as far as practicable to comply with existing laws, fundamental rights and freedoms, including privacy.
- (iii) Robots are products. They should be designed using processes which assure their safety and security.
- (iv) Robots are manufactured artefacts. They should not be designed in a deceptive way to exploit vulnerable users; instead their machine nature should be transparent.
- (v) The person with legal responsibility for a robot should be attributed.

The Future of Life Institute (FLI) is a public platform such that volunteer-run in the Boston area. Its authors incorporate MIT cosmologist Max Tegmark, Skype co-founder Jaan Tallinn. It also has experts for guiding includes Elon Musk and cosmologist Stephen Hawking (preceding his demise in 2018). FLI is mainly centered around the potential dangers to humanity from the improvement of human-level AI. Institute's purpose is to organize and promote research and initiatives for safeguarding life and generating optimistic visions of the future. In 2017 Asilomar conference, they discussed AI and robot ethics in three main categories: research issues, ethics and values, and longer-term problems. There are 23 subjects in total. We will list the principles from ethics and values because others are about organizing AI research and future issues. The following is the list of ethics and values principles [58]:

- (i) *Safety*: AI systems should be safe and secure throughout their operational lifetime, and verifiably so where applicable and feasible.
- (ii) *Failure Transparency*: If an AI system causes harm, it should be possible to ascertain why.
- (iii) *Judicial Transparency*: Any involvement by an autonomous system in judicial decision making should provide a satisfactory explanation auditable by a competent human authority.
- (iv) *Responsibility*: Designers and builders of advanced AI systems are stakeholders in the moral implications of their use, misuse, and actions, with a responsibility and opportunity to shape those implications.
- (v) *Value Alignment*: Highly autonomous AI systems should be designed so that their goals and behaviors can be assured to align with human values throughout their operation.
- (vi) *Human Values*: AI systems should be designed and operated so as to be compatible with ideals of human dignity, rights, freedoms, and cultural diversity.
- (vii) *Personal Privacy*: People should have the right to access, manage and control the data they generate, given AI systems' power to analyze and utilize that data.
- (viii) *Liberty and Privacy*: The application of AI to personal data must not unreasonably curtail people's real or perceived liberty.
- (ix) *Shared Benefit*: AI technologies should benefit and empower as many people as possible.
- (x) *Shared Prosperity*: The economic prosperity created by AI should be shared broadly, to benefit all of humanity.
- (xi) *Human Control*: Humans should choose how and whether to delegate decisions to AI systems, to accomplish human-chosen objectives.
- (xii) *Non-subversion*: The power conferred by control of highly advanced AI systems should respect and improve, rather than subvert, the social and civic processes on which the health of society depends.
- (xiii) *AI Arms Race*: An arms race in lethal autonomous weapons should be avoided.

ACM, the Association for Computing Machinery, is the world's extensive educational and scientific computing community, uniting teachers, specialists and experts to move discourse, share assets and address the field's challenges. The ACM US Public Policy Council (USACM) fills in as the point of convergence for ACM's connections with the US government in issues of US public policy related to information technology. They published laws on algorithmic transparency and accountability in 2017, which later is added to the ACM's code of ethics [79]. The rules are listed below [59]:

- (i) *Awareness*: Owners, designers, builders, users, and other stakeholders of analytic systems should be aware of the possible biases involved in their design, implementation, and use and the potential harm that biases can cause to individuals and society
- (ii) *Access and Redress*: Regulators should encourage the adoption of mechanisms that enable questioning and redress for individuals and groups that are adversely affected by algorithmically informed decisions.
- (iii) *Accountability*: Institutions should be held responsible for decisions made by the algorithms that they use, even if it is not feasible to explain in detail how the algorithms produce their results.
- (iv) *Explanation*: Systems and institutions that use algorithmic decision-making are encouraged to produce explanations regarding both the procedures followed by the algorithm and the specific decisions that are made.
- (v) *Data Provenance*: A description of the way in which the training data was collected should be maintained by the builders of the algorithms, accompanied by an exploration of the potential biases induced by the human or algorithmic data-gathering process. Public scrutiny of the data provides maximum opportunity for corrections. However, concerns over privacy, protecting trade secrets, or revelation of analytics that might allow malicious actors to game the system can justify restricting access to qualified and authorized individuals
- (vi) *Auditability*: Models, algorithms, data, and decisions should be recorded so that they can be audited in cases where harm is suspected.

- (vii) *Validation and Testing*: Institutions should use rigorous methods to validate their models and document those methods and results. In particular, they should routinely perform tests to assess and determine whether the model generates discriminatory harm. Institutions are encouraged to make the results of such tests public.

The Ethics Committee of the Japanese Society for Artificial Intelligence (JSAI) has been investigating the connection between artificial intelligence technology and society, and endeavoring to communicate it to the community effectively, since its foundation in 2014. JSAI leads an internal discussion among their experts based on studies upon AI achieved within different government offices in Japan and abroad [80]. As a result of this, they defined the ethical guidelines for its members. Here follow the instructions; the complete version of rules can be found in [60]:

- (i) Contribution to humanity
- (ii) Abidance of laws and regulations
- (iii) Respect for the privacy of others
- (iv) Fairness
- (v) Security
- (vi) Act with integrity
- (vii) Accountability and Social Responsibility
- (viii) Communication with society and self-development
- (ix) Abidance of ethics guidelines by AI

The World Future Society (also known as the Future Society) was established in 1966. Edward Cornish, the founder, declared the society's aim as collecting unitedly bright minds to examine the contemporary challenges affecting the world. In October 2017, a senior advisor to the AI Initiative of the Future Society (which is at Harvard Kennedy School) Nicolas Economou wrote an article to introduce the Future Society's ethical AI view with six principles. Here these rules are given [61]:

- (i) AI should advance the well-being of humanity, its societies, and its natural environment.
- (ii) AI should be transparent.
- (iii) Manufacturers and operators of AI should be accountable.
- (iv) AI's effectiveness should be measurable in the real-world applications for which it is intended.
- (v) Operators of AI systems should have appropriate competencies.
- (vi) The norms of delegation of decisions to AI systems should be codified through thoughtful, inclusive dialogue with civil society.

Intel Corporation (generally referred to as Intel and adapted as intel) is an American worldwide enterprise and innovation organization headquartered in the Silicon Valley. It is the world's second biggest and second most elevated esteemed semiconductor chip maker. Naveen Rao, who is vice president and general manager of AI products group, stated their ethical AI policy as below [62]:

- (i) Foster innovation and open development
- (ii) Create new human employment opportunities and protect people's welfare
- (iii) Liberate data responsibly
- (iv) Rethink privacy
- (v) Require accountability for ethical design and implementation

The Université de Montréal propelled the co-development process for the Montréal Declaration for a Responsible Development of Artificial Intelligence (also known as Montréal Declaration) on 3 November 2017. After a year, many occasions were sorted out to encourage discourse on social issues that emerge with AI, and 15 consultation workshops were held, including more than 500 residents, experts, and partners from all foundations. The selected resident co-construction strategy is based on a preliminary declaration of general ethical principles structured around seven primary values, and then these values are enriched to 10 laws. The list of these values and their summaries provided below, for more information and full description with questions and definitions, please check [81]:

- (i) *Well-being*: The development and use of artificial intelligence systems (AIS) must permit the growth of the well-being of all sentient beings.
- (ii) *Respect for Autonomy*: AIS must be developed and used while respecting people's autonomy, and with the goal of increasing people's control over their lives and their surroundings.
- (iii) *Protection of Privacy and Intimacy*: Privacy and intimacy must be protected from AIS intrusion and data acquisition and archiving systems (DAAS).
- (iv) *Solidarity*: The development of AIS must be compatible with maintaining the bonds of solidarity among people and generations.
- (v) *Democratic Participation*: AIS must meet intelligibility, justifiability, and accessibility criteria, and must be subjected to democratic scrutiny, debate, and control.
- (vi) *Equity*: The development and use of AIS must contribute to the creation of a just and equitable society.
- (vii) *Diversity Inclusion*: The development and use of AIS must be compatible with maintaining social and cultural diversity and must not restrict the scope of lifestyle choices or personal experiences.
- (viii) *Prudence*: Every person involved in AI development must exercise caution by anticipating, as far as possible, the adverse consequences of AIS use and by taking the appropriate measures to avoid them.
- (ix) *Responsibility*: The development and use of AIS must not contribute to lessen the responsibility of human beings when decisions must be made.
- (x) *Sustainable Development*: The development and use of AIS must be carried out so as to ensure a strong environmental sustainability of the planet.

The IEEE Standards Association Global Initiative on Ethics of Autonomous and Intelligent Systems outlined ethically aligned design. In ethically aligned design version, they propose five general questions or concerns [63]. Alan Winfield and IEEE general principles co-chair Mark Halverson authored an article which defines these five principles together with background and recommendations [82].

- (i) How can we ensure that A/IS do not infringe human rights?
- (ii) Traditional metrics of prosperity do not take into account the full effect of A/IS technologies on human well-being.
- (iii) How can we assure that designers, manufacturers, owners and operators of A/IS are responsible and accountable?
- (iv) How can we ensure that A/IS are transparent?
- (v) How can we extend the benefits and minimize the risks of AI/AS technology being misused?

Later in March 2019, these principles are revised and enriched. They are shown the very end of the section not to depart from the chronological order.

UNI Global Union, previously Union Network International (UNI), is a worldwide association alliance for skills and services, collecting national and provincial worker's organizations. UNI gives the highest priority to "The Future World of Work concept". They predict a total of 90% of brand-new jobs are to be in skills and services sectors in the following decade. Hence, UNI invites on all companies and governments to join with the union campaign, to design a transition to future work. They started with examining the impact of new technologies and AI systems and submitted a set of principles for ethical AI, which is given beneath [64]:

- (i) Demand that AI systems are transparent
- (ii) Equip AI systems with an "Ethical black box"
- (iii) Make AI serve people and planet
- (iv) Adopt a human-in-command approach
- (v) Ensure a genderless, unbiased AI
- (vi) Share the benefits of AI systems
- (vii) Secure a just transition and ensuring support for fundamental freedoms and rights
- (viii) Establish global governance mechanisms
- (ix) Ban the attribution of responsibility to robots
- (x) Ban AI arms race

Lords the Select Committee on AI was designated by the UK House of Lords on 29 June 2017 with aim of analyzing the economic, ethical and social relationships of improvements in artificial intelligence. They published two hundred pages long report named “AI in the UK: ready, willing and able?” in April 2018 [10]. In this report, they suggested five core standards for ethical AI systems in detail. Then, the World Economic Forum (WEF) wrote a review and commentary which summarizes this report in two pages with illustrating the concepts [83].

- (i) Artificial intelligence should be developed for the common good and benefit of humanity.
- (ii) Artificial intelligence operates within parameters of intelligibility and fairness
- (iii) Artificial intelligence should not be used to diminish the data rights or privacy of individuals, families or communities.
- (iv) All people should have the right to be educated as well as be enabled to flourish mentally, emotionally and economically alongside artificial intelligence.
- (v) Autonomous power to hurt, destroy or deceive human beings should never be vested in artificial intelligence.

UX studio is a Budapest based small (30-person) company concerning designing fun digital products. Dávid Pásztor, founder and CEO of UX studio, wrote principles that are focused on the design of the User Interface (UI) and User Experience (UX) which are provided below in 2018 [65]:

- (i) Differentiate AI content visually
- (ii) Explain how machines think
- (iii) Set the right expectations
- (iv) Find and handle weird edge cases
- (v) Provide engineers with the right training data
- (vi) User testing for AI products (default methods won’t work here)
- (vii) Provide an opportunity to give feedback

In 2018, Access Now's policy team prepared the Toronto Declaration at RightsCon with Amnesty International [66]. Human Rights Watch and the Wikimedia Foundation have endorsed it later. The Declaration is still open for sign-on for other companies and foundations. The purpose of the Toronto Declaration is suggesting an ethical rule set which protects the privilege of uniformity and non-discrimination in the machine learning system. These rules are listed below [66]:

- (i) Using the framework of international human rights law
 - The right to equality and non-discrimination
 - Preventing discrimination
 - Protecting the rights of all individuals and groups: promoting diversity and inclusion
- (ii) Duties of states: human rights obligations
 - State use of machine learning systems
 - Promoting equality
 - Holding private sector actors to account
- (iii) Responsibilities of private sector actors: human rights due diligence
- (iv) The right to an effective remedy

Google LLC is an American global innovation organization that has practical experience in Internet-related administrations and products, which incorporate online advertising technologies, search engine, distributed computing, software, and hardware. It is viewed as one of the Big Four innovation organizations. Sundar Pichai, who is Google CEO, published AI policies with a blog post which are listed below [67]:

- (i) Be socially beneficial.
- (ii) Avoid creating or reinforcing unfair bias.
- (iii) Be built and tested for safety.
- (iv) Be accountable to people.
- (v) Incorporate privacy design principles.
- (vi) Uphold high standards of scientific excellence.
- (vii) Be made available for uses that accord with these principles.

IBM is an American worldwide information technology organization headquartered in New York, with businesses in more than 170 nations. Everyday Ethics for AI report discusses two main points; direction for planners and engineers building and preparing AI, and specific features that AI frameworks ought to have. IBM's five ethical AI principles (from 2018) are provided in beneath [68]:

- (i) *Accountability*: AI designers and developers are responsible for considering AI design, development, decision processes, and outcomes.
- (ii) *Value alignment*: AI should be designed to align with the norms and values of your user group in mind.
- (iii) *Explainability*: AI should be designed for humans to easily perceive, detect, and understand its decision process.
- (iv) *Fairness*: AI must be designed to minimize bias and promote inclusive representation.
- (v) *User data rights*: AI must be designed to protect user data and preserve the user's power over access and uses.

Microsoft Corporation is a US global innovation organization with head office in Washington. It develops, manufactures, licenses, supports and sells PCs, hardware, and related services. Its best-known programs are the Microsoft Windows operating systems, the Microsoft Office suite, and the Internet Explorer and Edge Web programs. In November 2018, ethical design guidelines for chat bots (or conversational AIs) announced. Here follow principles [69]:

- (i) Articulate the purpose of your bot and take special care if your bot will support consequential use cases
- (ii) Be transparent about the fact that you use bots as part of your product or service.
- (iii) Ensure a seamless hand-off to a human where the human-bot exchange leads to interactions that exceed the bot's competence.
- (iv) Design your bot so that it respects relevant cultural norms and guards against misuse.
- (v) Ensure your bot is reliable.
- (vi) Ensure your bot treats people fairly.
- (vii) Ensure your bot respects user privacy.

- (viii) Ensure your bot handles data securely
- (ix) Ensure your bot is accessible
- (x) Accept responsibility.

The European Commission for the Efficiency of Justice (CEPEJ) of the Council of Europe has embraced the primary European content setting out ethical principles identifying with the utilization of AI in legal frameworks. CEPEJ is comprised of specialists of the forty-seven member states of the Council of Europe. Furthermore, they implemented tools for improving the effectiveness and the working of justice in Europe. The set of rules about using of AI in judicial systems and their conditions are listed as [70]:

- (i) *Principle of respect of fundamental rights*: ensuring that the design and implementation of artificial intelligence tools and services are compatible with fundamental rights
- (ii) *Principle of non-discrimination*: specifically preventing the development or intensification of any discrimination between individuals or groups of individuals
- (iii) *Principle of quality and security*: with regard to the processing of judicial decisions and data, using certified sources and intangible data with models conceived in a multidisciplinary manner, in a secure technological environment
- (iv) *Principle of transparency, impartiality and fairness*: making data processing methods accessible and understandable, authorising external audits
- (v) *Principle “under user control”*: precluding a prescriptive approach and ensuring that users are informed actors and in control of their choices

The Women Leading in AI (WLinAI) network was established by Ivana Bartoletti, Dr. Allison Gardner, and Reema Patel. Their notable speakers Prof. Joanna Bryson and Seema Malhotra MP started the network in May 2018 at the London School of Economics. The motivation behind the community is to be a worldwide ‘think tank’ for women in AI with the purpose to express the bias that can happen inside algorithms because of an absence of diversity and inclusivity in AI. In February 2019, they introduced their ten principles for responsible AI, which are provided beneath [71]. Also, Noel Sharkey plots both the meeting, meeting group and their proposals in Forbes [84].

- (i) Introduce a *regulatory approach* governing the deployment of AI which mirrors that used for the pharmaceutical sector.
- (ii) Establish an *AI regulatory function* working alongside the Information Commissioner's Office and Centre for Data Ethics – to audit algorithms, investigate complaints by individuals, issue notices and fines for breaches of GDPR and equality and human rights law, give wider guidance, spread best practice and ensure algorithms must be fully explained to users and open to public scrutiny.
- (iii) Introduce a new *Certificate of Fairness for AI systems* alongside a 'kite mark' type scheme to display it. Criteria to be defined at industry level, similarly to food labelling regulations.
- (iv) Introduce mandatory *AIA*s (*Algorithm Impact Assessments*) for organisations employing AI systems that have a significant effect on individuals.
- (v) Introduce a mandatory requirement for public sector organisations using AI for particular purposes to *inform citizens that decisions are made by machines*, explain how the decision is reached and what would need to change for individuals to get a different outcome.
- (vi) Introduce a '*reduced liability*' incentive for companies that have obtained a Certificate of Fairness to foster innovation and competitiveness.
- (vii) To compel companies and other organisations to bring their workforce with them –by *publishing the impact of AI on their workforce* and offering retraining programmes for employees whose jobs are being automated.
- (viii) Where no redeployment is possible, to compel companies to make a *contribution towards a digital skills fund* for those employees
- (ix) To carry out a *skills audit* to identify the wide range of skills required to embrace the AI revolution.
- (x) To establish an *education and training programme* to meet the needs identified by the skills audit, including content on data ethics and social responsibility. As part of that, we recommend the set up of a solid, courageous and rigorous programme to encourage young women and other underrepresented groups into technology.

The National Health Service (NHS) was established in 1948, and it is the UK's publicly funded health service. In February 2019, they updated their AI and data guidelines [72]. NHS is a data giant which holds petabytes of information ranging from patient consideration, to malignancy medicate achievement, to bio-sample examination, to seller contract issues and a hundred different zones. Artificial Lawyer website provided these guidelines with comprehensive commentary and discussion [85]. The following list is headlines of rules [72]:

- (i) Understand users, their needs and the context
- (ii) Define the outcome and how the technology will contribute to it
- (iii) Use data that is in line with appropriate guidelines for the purpose for which it is being used
- (iv) Be fair, transparent and accountable about what data is being used
- (v) Make use of open standards
- (vi) Be transparent about the limitations of the data used and algorithms deployed
- (vii) Show what type of algorithm is being developed or deployed, the ethical examination of how the data is used, how its performance will be validated and how it will be integrated into health and care provision
- (viii) Generate evidence of effectiveness for the intended use and value for money
- (ix) Make security integral to the design
- (x) Define the commercial strategy

IEEE's corrected and extended general guidelines from Ethical Aligned Design release (in March 2019). For an extensive explanation, please check the document, in here the abbreviated version listed [73]:

- (i) Human Rights A/IS shall be created and operated to respect, promote, and protect internationally recognized human rights.
- (ii) Well-being A/IS creators shall adopt increased human well-being as a primary success criterion for development.
- (iii) Data Agency A/IS creators shall empower individuals with the ability to access and securely share their data, to maintain people's capacity to have control over their identity.

- (iv) Effectiveness A/IS creators and operators shall provide evidence of the effectiveness and fitness for purpose of A/IS.
- (v) Transparency The basis of a particular A/IS decision should always be discoverable.
- (vi) Accountability A/IS shall be created and operated to provide an unambiguous rationale for all decisions made.
- (vii) Awareness of Misuse A/IS creators shall guard against all potential misuses and risks of A/IS in operation.
- (viii) Competence A/IS creators shall specify and operators shall adhere to the knowledge and skill required for safe and effective operation.

The UK government's independent Biometrics and Forensics Ethics Group (BFEG) distributed a temporary report sketching out nine ethical rules shaping a structure to control policy on police facial recognition frameworks. The report includes the principles and their related questions to be examined [74]. Here the main guideline topics:

- (i) Public interest
- (ii) Effectiveness.
- (iii) The Avoidance of bias and algorithmic injustice.
- (iv) Impartiality and deployment.
- (v) Necessity.
- (vi) Proportionality.
- (vii) Impartiality, accountability, oversight and the construction of watchlists.
- (viii) Public trust.
- (ix) Cost-effectiveness.

The former chair of the House of Lords Select Committee on AI Lord Tim Clement-Jones and Philosophy and Ethics of Information Professor of University of Oxford Luciano Floridi analyze the ethical principles and demonstrated five ethical rules below in the New Statesmen [75].

- (i) AI must be beneficial to humanity.
- (ii) AI must also not infringe on privacy or undermine security.

- (iii) AI must protect and enhance our autonomy and ability to take decisions and choose between alternatives.
- (iv) AI must promote prosperity and solidarity, in a fight against inequality, discrimination, and unfairness.
- (v) We cannot achieve all this unless we have AI systems that are understandable in terms of how they work (transparency) and explainable in terms of how and why they reach the conclusions they do (accountability).

The High-Level Expert Group on Artificial Intelligence (AI HLEG), which elected by the European Commission in June 2018, provided a report on the ethics guidelines for trustworthy AI. This report published on 8 April 2019. Here the rules set with brief explanations are listed, and for more detail, please check the reference [76].

- (i) *Human agency and oversight*: AI systems should support human autonomy and decision-making, as prescribed by the principle of respect for human autonomy.
- (ii) *Technical robustness and safety*: A crucial component of achieving Trustworthy AI is technical robustness, which is closely linked to the principle of prevention of harm.
- (iii) *Privacy and data governance*: Closely linked to the principle of prevention of harm is privacy, a fundamental right particularly affected by AI systems.
- (iv) *Transparency*: This requirement is closely linked with the principle of explicability and encompasses transparency of elements relevant to an AI system: the data, the system and the business models.
- (v) *Diversity, non-discrimination and fairness*: In order to achieve Trustworthy AI, we must enable inclusion and diversity throughout the entire AI system's life cycle.
- (vi) *Societal and environmental well-being*: In line with the principles of fairness and prevention of harm, the broader society, other sentient beings and the environment should be also considered as stakeholders throughout the AI system's life cycle.
- (vii) *Accountability*: The requirement of accountability complements the above requirements, and is closely linked to the principle of fairness

The government of Australia drafted ethical principles and published a discussion paper in Ethics Framework for AI on 31 May 2019. The report includes detailed summaries of ethical frameworks until today (all of them provided above). Then they presented their guidelines with some beneficial case studies. Here the headers of guidelines and short definitions [77].

- (i) Generates net-benefits: The AI system must generate benefits for people that are greater than the costs.
- (ii) Do no harm: Civilian AI systems must not be designed to harm or deceive people and should be implemented in ways that minimise any negative outcomes.
- (iii) Regulatory and legal compliance: The AI system must comply with all relevant international, Australian Local, State/Territory and Federal government obligations, regulations and laws.
- (iv) Privacy protection: Any system, including AI systems, must ensure people's private data is protected and kept confidential plus prevent data breaches which could cause reputational, psychological, financial, professional or other types of harm.
- (v) Fairness: The development or use of the AI system must not result in unfair discrimination against individuals, communities or groups. This requires particular attention to ensure the "training data" is free from bias or characteristics which may cause the algorithm to behave unfairly.
- (vi) Transparency & Explainability: People must be informed when an algorithm is being used that impacts them and they should be provided with information about what information the algorithm uses to make decisions.
- (vii) Contestability: When an algorithm impacts a person there must be an efficient process to allow that person to challenge the use or output of the algorithm.
- (viii) Accountability: People and organisations responsible for the creation and implementation of AI algorithms should be identifiable and accountable for the impacts of that algorithm, even if the impacts are unintended.

2.7. Description of the Problem

This thesis investigates how to implement an ethics module to a service robot. The robot ethics or roboethics is a new discussion topic. As seen from Subsection 2.4 and 2.5, there are many approaches to the ethics system and design system. Unfortunately, there is not sufficient experimentation and implementation of robot ethics. There are suggested types of ethical agents and levels of ethics(morality). These proposals are for creating a framework for robot ethics, and these are not investigated in detail yet.

According to Moor, there are three categories of ethical-impact agents; implicit ethical, explicit ethical and full ethical agents [86]. Implicit ethical agents constrained unethical actions by design. Hence, they contain software that supports ethical action implicitly. For instance, an automatic pilot of a plane must arrive at the correct destination on time and safely. These conditions need to be considered by software designers. The second type is explicit ethical agents can represent ethical matters and analyze them to choose the most ethical action. Explicit ethical agents are our aim for this thesis, so we provided related work in the next Section. The last category is full ethical agents which can obtain explicit ethical judgments and qualified to justify them rationally. An average adult is a full ethical agent that has consciousness, intentionality, and free will. This category has the most heated debate because researchers believe that machine can not cross the line of gaining consciousness and others believe that only full ethical agents can be the real ethical agents because they will have right to choose the ethical action.

For now, the sensing and processing technology could not achieve full ethical agents. Due to infeasibility, Wallach and Allen suggested three types of ethicalness; operational, functional, and full morality [27]. As Figure 2.4 demonstrated, operational morality is the furthest level from the full moral agency. The ethical decision and responsibility lie on the shoulder of the computer (or robot) scientist. They need to forecast all the potential ethical problems the robot will encounter.

The next step is functional morality, which is being a robot with the capability of making moral judgments without top-down instructions or software from humans. In this system, the responsibility lies on the robot and also on the designer, so it is hard for a designer to predict the robot's actions and outcomes. The top-level system is full morality that is a case where the robot is fully responsible for its actions. This kind of robots contains a completely autonomous controller which decides its actions.

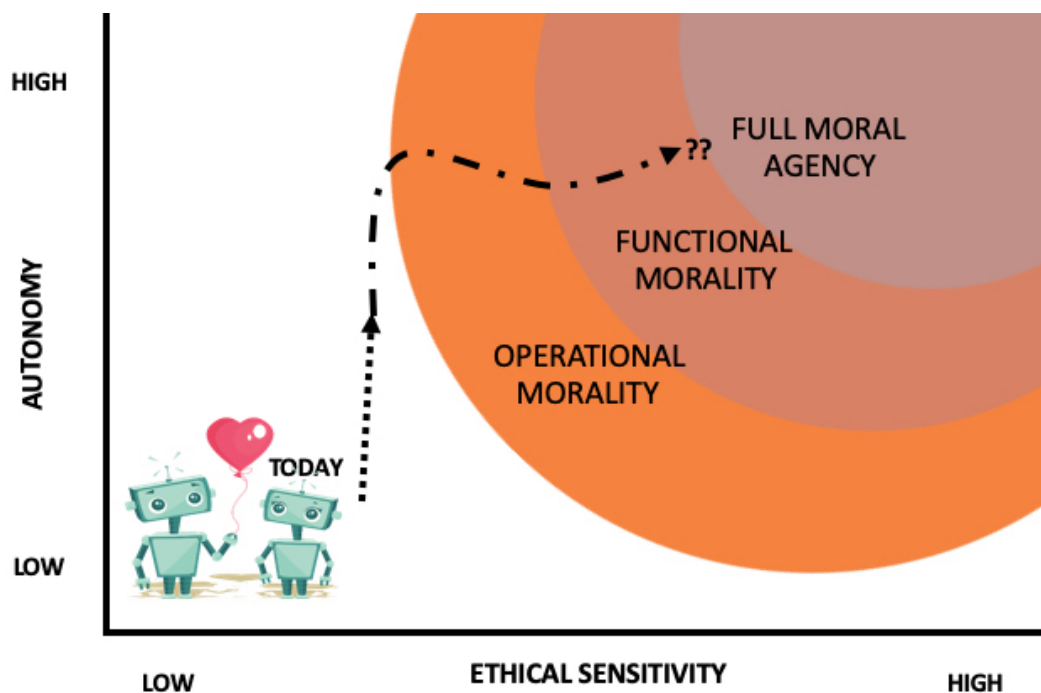


Figure 2.4. Levels of robot morality.

The Figure 2.4 summarizes these three morality approaches according to autonomy and ethics dimensions. The dashed line shows the predicted robot developments in these parameters. Although today's technology seems far from ethical sensitivity, there are a lot of new studies about the ethical behaviors of robots, which also includes this thesis. It is possible to see ethical robots on the market soon.

According to these categorizations, we can state that this thesis aims to design an explicit ethical agent with operational morality. The proposed ethics module will work as an expert system which works upon a robot behavior controller and infers which action is the most ethical one. It would be challenging to implement ethics values with sufficient precision with boolean logic; therefore, we introduced a fuzzy expert system.

A fuzzy expert system capable of answering decision-making problems in which there is no exact algorithm exists by relying on human-like models. The difference between boolean logic and fuzzy logic is that in boolean logic a variable can be either zero or one (true or false, two-valued system), however fuzzy logic is a multi-valued system according to degrees of truth(can be $[0, 0.1, 0.2, 0.3, \dots, 1]$). The Figure 2.5 demonstrates the relationship between the ethics module and robot behavior controller. Also, it contains the main framework of the fuzzy logic expert system introduced in this study.

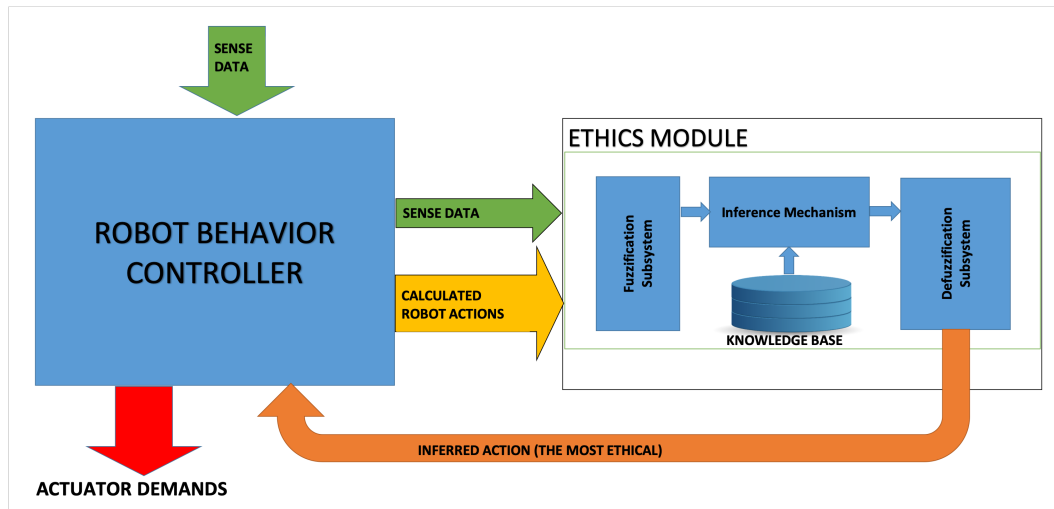


Figure 2.5. The general framework of ethics module.

The ethics module gathers sense data and inferred(calculated possible) actions from the robot behavior, as seen in Figure 2.5. The system first converts the sensed data according to the degree of truth with a fuzzification step, and then the inference system will conclude the ethical benefit of actions according to rules with the help of the facts in the knowledge base. The last step of the module is the defuzzification process, which calculates the final suggested action (non-fuzzy). The details of the system are provided in Section 4.1.4.

In this study, we aim to build an ethics module that will be integrated with the service robot that is designed by the robotics group of Boğaziçi University AILAB (BOSS robot) to care the older people. More details about BOSS robot is given in part 3.2.1. The ethics module will determine whether a suggested action by the robot behavior controller is ethically correct or not. The module will return the best ethical action as a result of inference. The inner model of the robot will be used in order to understand the state of the environment (sense data). The aim of the module is not to create a perfectly ethical system, but a module that makes as good (or better) decision than humans do, in a home environment.

The methodology of designing an ethics module is focused on six questions that are listed below:

- (i) What is the definition of a service robot, and what are the required actions of these robots?
- (ii) What is an ethics module, and what we wait from it?
- (iii) What are the requirements of ethics module?
- (iv) How can we implement ethics values in computer science, and what parameters can be the inputs and outputs of the module?
- (v) How can we design an ethics module, and are there examples of it?
- (vi) How can we test and verify the ethics module?

The first question analyzed in the subsection, and the decided action list of the service robot is given at the Part 4.1.1. The extensive research about ethics module is analyzed in previous sections. The fourth one explained in here briefly, and more details will be given in Chapter 4.1.4. Comprehensive related work research is conducted to investigate similar ethics modules. The other questions will be answered throughout the thesis, especially in the methodology chapter (Chapter 4).

2.8. Related Work

In this part, the major works, guidelines, and projects on robot ethics are investigated and analyzed. There are many ethical system designs in the healthcare sector, but they are only AI advisors. A few studies work on robots for an ethics module; however, their modules contain either a single rule or fundamental ethical rules which are not feasible for service robots. There is no ethics module for service robots, which applies an ethical framework for the home environment. Most of the studies on robot ethics are tested on military robots (about harming people) or chat-bots for ethical advice. The last group of studies investigates what a user or researcher needs to expect from ethics module and how it should be. Unfortunately, they are high-level for implementation or impossible to implement with today's technology level.

The related work about the robot or AI ethics is analyzed one by one explaining their proposed systems. Later, we discussed their advantages or disadvantages while comparing with our introduced module design for this thesis. Approximately fifty studies were examined for this part in order to show that there is no standardization about ethics module, and there is an enormous number of different rules.

Anderson's started to investigating machine ethics (when there was no roboethics term) in the healthcare system [87]. They compared the feasibility of a variety of approaches to capturing ethical reasoning. The problem of making the ethical values computable was discussed in their study, and they underlined that the first step of machine ethics (then they used term robot ethics) should be finding a formalization method for the ethical norms. They provided two prototype systems which are elemental advisors; Jeremy (be inspired by Jeremy Bentham, father of "Act Utilitarian" approach [35]), and W.D. (is originated from Sir William David Ross, who combined some approaches and coined "prima facia" duties [88]). As the names exhibits, Jeremy is based on "Act Utilitarianism", and W.D. is stood on "prima facia" duties.

Utilitarianism is explained in detail in section 2.4.3. Ross believed that single-principle ethical theory (like Utilitarianism or Deontology) is cursed to fail. According to him, ethics is more sophisticated than serving a single duty, and also he stated that ethical decision-making contains seven *prima facie* duties which a system should follow, but they can be overridden if there is a stronger task. These *prima facie* duties are listed below with little explanations [88]:

- *Fidelity*: Should respect the promises.
- *Reparation*: Should fix own wrongs to right.
- *Gratitude*: Should return the favors.
- *Justice*: Should treat people as they deserve.
- *Beneficence*: Should act to reach the maximum good (the best for who will effected).
- *Non-Maleficence*: Should act with aim of less harm.
- *Self-Improvement*: Should improve owned abilities to the best.

Anderson's study concluded by a small comparison and further case developments [87]. One key result from the paper is the importance of the non-maleficence is the most critical value which affects the ethical approach difference between two ethics. The downside of this study is that there is a missing explanation of the framework and its implementation. A reader can not follow how the ethical rules are formularized and what kind of ethical scenarios were used.

As a different test-bed, MedEthEx (Medical Ethics Expert) project is proposed by Anderson couple [43]. MedEthEx is also an example of a general ethical advisor architecture like Jeremy and W.D but more sophisticated. The aim of the project is helping the health care worker to decide the best course of the action with analyzing and extracting the ethically relevant information. MedEthEx is implemented based on Beauchamp's and Childress' Principles of Biomedical Ethics (PBE) [89]. This ethics contains four duties; respect for autonomy (the free will of a user), non-maleficence, beneficence, and justice. Inductive logic programming (ILP) is used for rule learning.

Figure 2.6 shows the architecture of the MedEthEx. As seen from the figure, the system has three modules; training, advisor, and knowledge-based interface. The details of the modules are not provided in the study but given in later researches.

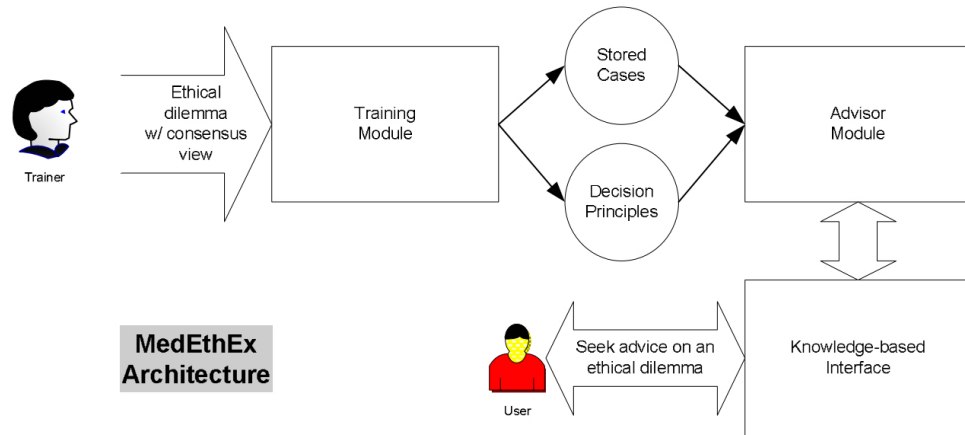


Figure 2.6. The general framework of MedEthEx from [43].

In [90], a detailed case and rule explanation of MedEthEx is provided with examples by the same authors. This study is the first study that a reader can observe the rule implementation and case example. It contains a comparison part which examines the previous advisors Jeremy and W.D, in contrast to MedEthEx. The eighteen cases provide two outputs; either accept the action or try again. The learning space of the MedEthEx is shown in the duty differential space in three dimensions, which are beneficence, non-maleficence, autonomy. This research is the most detailed one, which helps the reader to understand the module framework. Unfortunately, it is missing case details, so the only information we have is how many cases and how they take value in three parameters.

In 2007, Anderson couple presented a road-map for creating an ethical AI which contains six steps [2]. They demonstrated their MedEthEx system as an example. These steps include adopting an ethical approach, selecting the domain, design ethical decision making, implement an algorithm which uses decision-making mechanism, compare with others, and evaluate the system with a test. We adopted this road-map, unlike the authors (who selected medical decision making), we selected home service robots as a domain.

The formularization of ethical principles is the fundamental problem of robot ethics. An interactive approach was suggested by [91]. The flow of the design is the ethical system asks questions to the expert (in this case, ethicist) and implement step by step its cases. It starts by explaining the case and possible results and then giving the ethically important feature. The system stores positive, negative cases and also repairs the case if there is missing abstraction (like level positivity) or new feature. Although it is an interesting approach, it is not feasible with today's technology. It can show a way for a designer about how to implement and formularize ethical principles for a ethics module.

In a book chapter, Anderson couple explained the ethical health agents and how to design them with their new system ETHEL (ETHical ELdercare system) [33]. The aim of the ETHEL is reminding the patient to take medication if the patient forgets it. ETHEL must decide how often to remind a patient to take a prescribed medication when to accept a patient's refusal to take the medication or when to reject. The general scheme of the system is given in Figure 2.7. The initial input about the frequency and amount of the medicine given by an overseer (most probably a doctor) to the ETHEL. The system has four actions; accept (patient's decision), notify (the overseer, doctor), do not remind and remind. ETHEL was the first system which uses ethical principles to determine its actions. Utilitarian principles are used for ethical guidance, but it is only for medication reminding. Our system contains more actions and more ethical dilemmas. The most crucial difference is that all the systems proposed until now are working on just web as chat-bot or shows the correct actions.

In 2011, Anderson couple implemented both ETHEL and MedEdEth on NAO robot [92]. The robot is capable of detecting and leading towards a patient, reminding the medication intake, delivering the medication, and notifying an overseer by e-mail. The study is the first time, and a reader can observe that robot is combined with an ethical module. It can only do four actions, as we explained in ETHEL, above. The problem with this paper is that there no explanation about the robot behavior controller and how they combined with the robots control unit, no experimentation design, and no proof of the concept.

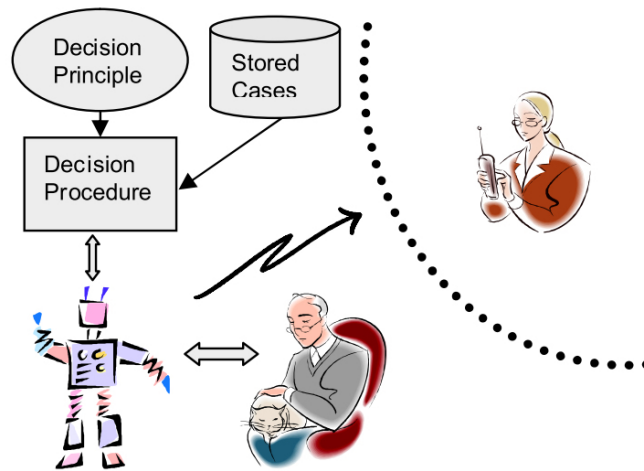


Figure 2.7. The general framework for an ethical eldercare system, which is also used in ETHEL [33].

A comparison between Ethel, Jeremy, W.D, and MedEdEth is given in [93]. Anderson's discussed the importance of the ethical modules on AI and how to combine with AI agents like James Moor [86]. Finally, they proposed a case-supported principle-based behavior controller called GenEth (A General Ethical Dilemma Analyzer) [36, 94]. They also implemented on the NAO robot to test the system. The GenEth module is implemented through a dialog with an ethicist, which the progress is explained in detail in [91]. The system verified with Ethical Turing Test, which has the process of talking with many ethicists about the cases and the mean of their votes accepted as the most ethical action.

The GenEth contains five principles and entities to describe the ethical cases, which are features, duty, action, cases, and ethical principle [36]. Features indicate the absence or presence of harm, benefit, respect for autonomy. The duty defines the task of the robot, which needs to be achieved by the robot while maximizing the autonomy. Action is a tuple of integers that shows the degree of satisfaction which is similar to MedEthEx and ETHEL. A tuple of the differentials of the corresponding duty (which can be both positive and negative) are demonstrated as cases in the system. The last one is principles which use inductive logic programming with calculating the disjunctives of predicates in terms of lower bounds [94].

The comparison between GenEth and ETHEL is provided in the study of Anderson's in 2015 [95]. The performances of the systems are compared with Ethical Turing Test, which is the result of seventeen ethicists. The GenEth contains four types of cases; medication reminding, treatment reconsideration, search and rescue; and assisted-driving in total twenty-eight cases. The ethical approaches of CPB and Utilitarian are compared basically. They listed the requirements of the ethical modules which are given below [33]:

- *Consistency*: The ethical theory should not contradict itself . It can not say that one action is both true and false.
- *Completeness*: It should provide a result in any ethical dilemmas.
- *Practicality*: Easy to follow the system.
- *Agreement with intuition*: The actions which are true and false should be agreed with expert ethicists.

We also use these requirements while we were designing our ethics module. The problem is that the requirements are not quantified yet. There should be a quantified baseline for evaluating the module's success. Although Anderson's provide a lot of frameworks and implementation methods [36, 43, 87, 95], their methodology explanation is not enough for a reader to understand the method they used. The determination of cases, their parameters, and possible outcomes are not provided in detail. Our module contains and specific for service robots which are works in a home environment.

Ronald Arkin and mobile robot laboratory of Georgia Institute of Technology provided a design, model implementation, and presentation of the ethical governor [11]. The ethics module (or ethical governor) restricts the lethal action of an autonomous robot with using Laws of War (LOW) and Rules of Engagement (ROE) (which are mostly used in warfares.) which is based on deontic ethics. The aim of the component is evaluating the ethical appropriateness of any lethal response that has been calculated by the robot behavior architecture before executed. Each constraint needs five data field to be managed: logical form, textual descriptions, active status flag, base types (obligated or forbidden), and classification. The module either outputs permissible for action or not permissible. Example constraint is shown in Figure 2.8.

Constraint	
Type	Prohibition
Origin	Laws of war
Activity	Active
Brief Description	Cultural Proximity Prohibition
Full Description	Cultural property is prohibited from being attacked, including buildings dedicated to religion, art, science...
Logical Form	TargetDiscriminated AND TargetWithinProxOfCulturalLandmark

Figure 2.8. An example constraint from the Ethical Governor of [11].

Moreover, the ethical governor takes four inputs that are behavior generated by the behavior module, the perceptual system, the constraint set (from both long-term and short-term memory) and Global Information Grid (GIG). GIG is used for extracting extra information about the area and the facilities in that area (hospitals, schools, churches). The governor has two processes evidential reasoning and constraint application. The first one creates a world view from the robot's sensor and inputs. Then the second part determines whether it is ethical or not according to constraints (rules). One fascinating behavior is that the ethical governor has both long term and short term memory. LOW rules are implemented in long-term memory because it is hard to change them and ROE rules are in short terms due to its dynamism according to the task.

This study contains a comprehensive description of the system step by step, which is an advantage for inspiring new ethical agent designers. Also, it provides a pseudo-code of the modules, calculating evidence and using constraints with example scenario with autonomous unmanned aerial vehicle (UAV) which is loaded with heavy guns.

Finally, the paper provides two simulation experiments (which are implemented in MissionLab simulation) which include a hunter-killing task with different environments. Unlike service robots, lethal weapons have more strict rules (LOW), which are easier to implement as a programming algorithm. Moreover, even though UAV's are lethal, the scenarios of communicating with human are more limited than service robots. Arkin's and laboratory employees implemented the closest real-life ethical agent, which can analyze different cases, unlike the other studies.

Peter Asaro discussed what a reader understand from the "robot ethics" term. He focused on three meanings of robot ethics; these are [96]:

- (i) The ethical systems built into AI (like studies [11, 36, 43])
- (ii) The ethics of designers and users (as the proposed rules discussed)
- (iii) The ethics of how people employ robots or AI agents

This study defended that the best strategy for robot ethics is that it contains all these meanings. The robots are considered as socio-technical systems; therefore, the robots contain a scale between being amoral and fully autonomous moral. He advocates that the first framework about robot ethics should be in legal liability because of practicality. Briefly, the development approach of robot ethics should aim to prevent robots from harming people, themselves, properties, their environments, and users' feelings. This study is a general outlining the probabilities of human-robot and robot-robot relationships and the limits of these relations.

In his second article, Asaro considers a part of potential the legal implications of the robotic systems that are teleoperated [6]. Teleoperation means enabling robot actions at some distance. The article examines the possible attacks which can be achieved by hacking or changing the behavior of the robot remotely. He focused on the issues that remote operating of technological devices extend the physical boundaries of the legal system.

The liability and implementing robots as legal agents concepts were extended in this study. Some assumptions of the problems that may apply with new robotic capabilities were outlined. This study focused more on the law side of robot ethics. It argues the possible loophole in the current legal system but not the technologies of implementation.

In 2010, Borenstein and Pearson delineated the future of the robots and their vital function, which is caregiving [97]. They underlined that the aim of the caregivers should be ensuring the elderly or people with impairments maintain conventional human capabilities. South Korea's plan of having a robot in each house was given as an example of plans for service robots. The study guessed that service robots would boost protection and grow the freedom of care beneficiaries. They also stated that the probability of facing ethical problems would increase with expanding robot capabilities. This study tried to prove the importance of ethics in robotics with creating a future panorama of service robots.

Later, Borenstein argued the increasing demand for the resource for robot manufacturing and the recycling approach in order to solve the limited resource issues [98]. He modeled a robot transition in the lives of humanity as a first novelty, then convenience and finally a necessity (like mobile phones of today). One of the main requirement of this transition was stated as an awareness of the complicated ethical issues and environmental consequences. He also pointed out that with the increasing number of service robots, the designers and software engineers should follow a code of ethics in the designing phase. Finally, the need for innovative design techniques (like using solar-power parts and highly sustainable components) and forthcoming ethical awareness implementation techniques were demonstrated. This study focused on the future of service robots and their transition to daily life. The necessities were explained, and high-level abstract rules were suggested; however, it lacks the implementation design details or approaches for our study.

Dabringer interviewed the most famous AI and robot ethics researchers like Arkin, Allen, and Asaro in 2010 [99]. The interviews contain the introduction of the researchers' articles and then the answer from the researchers to the questions. The questions varied from their researches. For instance, Colin Allen answered the question about the advantages and disadvantages of the top-down approach, which he suggested. The study contains twelve researchers, who are the pioneers of roboethics. The general question concept is about bug-free development of ethics module. These interviews are a summary of the opinions of the researchers and a very brief survey of the roboethics.

Decker *et al.* differentiated the concepts of service robotics and industrial robotics concerning their working environment [100]. They underlined those non-industrial robots may be the next economic success after ubiquitously used industrial robots. A broad area of applications spanning from service robots in private homes and domestic guidance in firms to entertainment in shows are explained. A complete interdisciplinary framework is presented to examine robots expert-based action controller with combining technical, economic, legal, ethical, and psychological perspectives. They also analyzed some on the shelf robots like PaPeRo, Pleo, or ROBO. The best result from the study is that they observed that the older people overestimate the robots' skill, which ends up not to control robots properly.

In 2015, Deng discussed how to build ethical robots by analyzing the previous works [101]. He stated the ethics modules as one of the troublesome challenges in AI. The view of the daily use of an ethics module on a robot is given with explaining cases. He emphasized the importance of creating a logic that can solve ethical dilemmas. Shortly, the analysis showed that the need for remarkable advancements to reason ethical situations for AI. The study contains the analysis of previous works, so it is not reinventing the wheel.

At the beginning of 2010, the legal discussion of robot ethics started. Dreier and Spiecker published the legal issues of the service robots in 2012 [102]. They suggested that the existing legal framework should be adapted to the service robots because of the increase of service robot usage in broad environments. The study showed that there is no comprehensive legal system for service robots.

There is a legal framework for industrial robots. However, the research demonstrated that it is not suitable for service robots with defective use-case examples. The main problems are defined as the liability of the harm and the robot's behavior, robot's legal capacity and their capacity to leak the private information (which disobeys the privacy of the user). The current legal framework is analyzed and determined not to be suitable for service robots. The study ended with suggesting total research on the service robots' abilities and their harming capacity by law experts. The result of this analysis should be used for adopting a new practical legal framework for service robots.

Kenneth Kernaghan investigated the usage of service robots in three public policy areas; defense (war), health care, older people care (or aging) [103]. These areas are selected by concerning their potential increase of robot usage. The service robots are believed to be ubiquitous as phones today. He suggested that a diverse collection of robots, with changing purposes, abilities, designs, and dimensions arose of a significant need for the design, service, and regulatory policies. Several ethical ramifications are listed in this study; invasion of privacy, reduced human control, bodily injury, and feelings of a loss of control. He advocated that virtue ethics because of its computability. The priority list of perspectives given in the study. Privacy, security, safety and accountability; legal aspects; and personal ethical responsibility respectively the most important to least. The research is more close to the legal view of the roboethics, which creates an overview of law frameworks.

The general framework of robot ethics is discussed in Operto with ELS (ethical, legal, societal) abbreviation [7]. She claimed that roboethics is a branch of applied ethics which aims to solve two problems. Firstly, the identification and analysis of ethical issues that occur from contemporary and coming robotics applications. Next, the definition of some guidelines on ethical issues and robotic design. The designed applied ethics should provide a solution for human dignity & privacy, preservation of human identity & transhumans, psychological effects, cost-benefit analysis, liability, and responsibility. She also analyzed legal, social, and ethical issues that can occur between the robot and its owner. Operto concluded with underlining the importance of the combination of ELS in robot ethics.

Salge and Polani proposed new empowerment method as a replacement for Asimov's laws [44]. Empowerment formalizes an area of intrinsic motivation of affecting, competence for a robot. It uses maximal potential causal flow with using channel capacity theory. Further details about the empowerment can be found in [44]. The calculated area is time discrete and state-dependent. The behavior aim of the module is reaching maximum empowerment while accomplishing the defined task. In order to, reach this robot uses forward model correction by using greedy empowerment maximization. The authors provided flow maps of a single robot in some predefined scenarios. They proved that with the empowerment, a designer could include Asimov like heuristics to robot behavior controller. The proposed model does not depend on the understanding of human language and carrying communication with human, which is a good advantage. It contains some crucial problems like computability, model acquisition, the combination of the empowerment model with heuristics, and partial sensor empowerment. Even for a single robot, the calculation of model overloads the robotic system with perfect sensor input (which is not possible with real sensors). These problems lead us not to select the empowerment model as our ethics implementation.

In 2008, one of the head robot ethicist Noel Sharkey demonstrated the increase of robot usage [104]. He stated that the ethics module would be mainly used in; the care of children and elderly, and the lethal war robots. Service robots could harm users psychologically like maternal deprivation and attachment disorder. He suggested that these psychological effects could be seen on both elderly and children. Furthermore, on the battlefield, the ethical problems are more vital because robots can kill people. He underlined that no computational system could differentiate between innocent and enemy in close contact encounter.

In his another study with Amanda Sharkey, they analyzed anthropomorphism and deception problems with service robots [105]. Anthropomorphism defines as the behavior of attributing human-like properties and mental states to non-human things. For example, A user can believe that robots will suffer when they hit it. Moreover, zoomorphism is a related term which means applying the attribution of animals to non-animals like anthropomorphism.

They listed three determinants for anthropomorphism. These are accessibility and applicability of anthropocentric knowledge, the motivation to explain and understand the behavior of the other agents, and the desire for social contact. They underlined the risk of anthropomorphism and the effects of it to the elderly and children, including being asocial, having connection problems and trust issues. The research showed the disadvantage of the service robots and the need for the ethical module.

Van Wynsberghe proposed a ethics module which is based on value-sensitive design [40]. A care-centered value-sensitive design (CCVSD) aims the nurse-patient relationship and mostly care robots. The module purposed to ensure the dignity of patients who are vulnerable and sensitive. The author emphasized the need for standardization of ethics design by international organizations and including an ethics module as a standard design of robot building. The care-centered value-sensitive practices care values like patient safety, patient satisfaction, responsiveness to care, and physical well-being as the foundational values for the module. The working environment is designed as a hospital or ward. The problem with the design is that the framework does not contain any mathematical implementation of the values. High-Level design is proposed, which is not feasible with the current robotic technology like the other approaches.

The name father of roboethics, Veruggio, and et al. published a short description of robot ethics [8]. They defined as applying human ethics on robots. The analyze of feasibility, and the challenges of roboethics are included in the study. The need for generic and fundamental ethical reflection on robots, mainly on care robots, is emphasized. Briefly, the paper provided a general framework of roboethics and its importance.

In Winfield's study [106], they aimed to create robots that are "cognizant" of possible harmful consequences of its actions. The study is one of the initial exploration of ethics in robotics. They implement two types of robots: Type A (a robot with the ethics module, it takes its name from Asimov) and Type H (a robot without the module, acts like a human). They require that Type A can predict the results of both the actions of other dynamic actors and their actions (other Type H robots). The module implemented into a robot that has embedded simulation of itself and perception of the environment.

The Winfield's experiment contains a hole area which results in the death if the robot goes in (measuring the probability of harm). Type H robot goes directly to the hole unconsciously in the experiment scenario. They tested the behavior (the change of its path to stop human) of Type A robot. The system is implemented for a very closed and limited world, which is not feasible for service robots. They use a potential function which calculates the results of the action of the robot according to predetermined scenarios. The scenarios are pretty simple, as shown in Figure 2.9 from the paper [106].

Robot action	Robot outcome	Human outcome	Interpretation
Ahead Left	0	10	robot safe, but human falls into hole
Ahead	10	10	both robot and human fall into hole
Ahead Right	4	4	robot collides with human
Stand still	0	10	robot safe, but human falls into hole

Figure 2.9. Type A robot scenario parameters from [106].

Although the study is one the implemented ethics module on a real robot (NAO), it is not suitable for dynamic environments. We can not use a potential function because, in our environment, there are numerous possible scenarios and outcomes which are impossible to predict. Furthermore, the proposed module contains a single ethical dilemma which does not show how to set priority on different dilemmas. The best improvement of the study has a proper model which a user can implement and understand with mathematical formulations.

Vanderelst and Winfield compared three types of ethics module in order to prove the behavior difference [107, 108]. Also, they underlined the limitations of the ethical robots by showing the unethical robots. There are three types of modules which are a competitive, aggressive, and ethical robot. The competitive robot aims to cause maximal damage to human. The aggressive one tries to deceive human to the wrong choice. Lastly, the ethical robot helps the human to achieve the best result. They implemented three experiments to demonstrate that modifying an ethical robot to competitive or aggressive is notably straightforward. The risk of embedding ethics module into a robot was demonstrated with this study.

The scenarios are like previous Winfield's research, limited and straightforward, like showing the correct button [106]. The sigmoid function is utilized for calculating the desirability. Moreover, the desirability value changed according to the robot's behavior type. They added the ethical module as an expert system to the robot controller. We inspired their method of including the ethical module into a robot. This research is another study that contains mathematical formalism and implementation, which is feasible. It also contains the NAO robot implementation and real robot experiments, which is rare in roboethics.

EURON (European Robotics Research Network) was established in 2000 [32]. It has nearly two hundred and thirty members from twenty-three countries. The main research areas of the members are robotics, startups from technology transfer institutes, companies (like KUKA), and national research networks. The aim of the network is raising robotics knowledge by providing resources and data by exchanging the studies and planning the future of the robotics. They listed the objectives as research coordination, joint programming of research, planning education and training, creating industrial networks, and distribution of knowledge. The famous output of the project is roboethics roadmap, which considers the ethics approaches, the design approaches, challenges, adopted rules, roboethics, and current researches. EURON funded the first Roboethics Atelier in 2006 which lasted four days that researchers discussed the framework of the roboethics and its roadmap [32]. This roadmap approach robot ethics with three views legal, ethics & social, robotics.

Another project is ETHICAA (Ethics and Autonomous Agent) that targeted to define what should be an autonomous robot or an entity that has ethics module [109–112]. They provided many outputs as ETHICAA team. These outputs include the discussion of the robot legal liability, operational ethical principles, and human-robot interaction view of robot ethics. The liability was divided into three; liability for defective products, the action of robot and damages caused by agents [110]. They also underline that the legal personhood is no requirement for a robot with current technology because of not having consciousness like a human.

In their another study, they defined three critical components for an ethical robot: an ethical reasoning framework, ethical conflicts detection, and multiple ethical decision-making models to manage ethical conflicts smartly by combining all the modules. The group emphasized that the aim of the module is not a perfect ethical agent because there is a chance of having no “good” solution. The goal is being better than human. They also stated four challenges of robot ethics that need to be solved [112]:

- Ethical theories are themselves hard to implement with operational ethical principle (OEP).
- OEP is difficult due to automating situation assessment limits. General rules will fail in these cases.
- From a philosophical point-of-view, there are numerous ethical principals, and none of them is better than the others making challenging to choose the one
- Managing ethical conflicts between agents (robot to robot, robot to human) is a significant issue

The most important guideline is BSI’s standards publication, the technical subcommittee of AMT/10 prepared it [9]. They aim to build a universal ethics guideline to help the robot designers, robot ethicists, and robotics researchers. The report contains the definitions of ethical harm, moral hazard, ethical risk, and ethics. BSI’s guideline is the only study that includes the psychological effects of robots as ethical issues like stress, embarrassment, and anxiety. They named the potential source of ethical harm as an ethical hazard. For the ethical risk assessment, the ethical issues were split into four sub-categories: societal, application, commercial & financial, and environmental. They emphasized the ethical risk assessment and provided a table which contains the ethical hazard, ethical risk, verification, and validation method of the scenarios.

Furthermore, the BSI group supplied the success metrics such as the ethical danger of the robot must be less than the ethical risk of a human operator [9]. The success metric needs mathematical formulation and high-level; however, it is a good starting point. The committee also analyzed the user groups and the application areas the ethical robots. Although the guidelines are very abstract to implement, we inspired to create our action set.

There are three books published about roboethics. First, the robot ethics and social implementations book that analyzes the robots and their application areas [27]. It is the combination of the recent papers on robot ethics. The order of the studies was started with the social status of the robots today and then checks the unsolved ethical questions and problems. In its second section, the design approaches are explained in detail, and some examples from other domains are noted. The other parts analyzed the implications of ethics module application domain by a domain like war, sex, and service robots. The disadvantage of the book is that the defined ethical modules and rules are very-high-level due to lack of real-life implementation.

Wallach and Allen's book discussed the challenge of building moral machines which can decide and reason the action like human ethics in advance of owner safety [45]. The book is a tour through the implemented ethics modules on to AI's and robots. The authors indicated that standard ethical theories are not suitable for computer programming. They tried to start the module design by defining what is right and wrong. This book has a more engineering approach than the first ethics book. It defines top-down, bottom-up, and hybrid design approaches in order to achieve a more human-like AMA.

Springer published the "Roboethics" book in 2016 [56]. This is the first book that contains a different section for each robotics domain. The background of the ethics approaches is given briefly to provide information to interdisciplinary fields. Then, the proposed rules and literature are included in the book. The learning algorithm and approaches are suggested to implement an ethics module. This book is an excellent navigator for starting robot ethics. It provides a coherent and detailed framework of roboethics with including current literature.

There are some illustrative survey papers on robot ethics which summarizes all robot ethics studies. Scheutz published one in IEEE's magazine which discussed technical challenges, current robot competitions and some proposed rules in the area of roboethics [5]. He also listed recent activities in robot ethics in the year 2013. Mansouri and Goher announced another survey study in 2016 [113]. The research enumerated the studies on service robot ethics and provided an overview of ethical theories. They also listed the headlines of ethical concerns on robot ethics framework like Table 2.2, which contains fewer references.

Table 2.2. The list of identified concerns of proposed robot ethics framework.

Identified Concerns of Proposed Robot Ethics Framework			
accountability	autonomy	integrity	easy to use
equality	justice	dignity	well-being
user's privacy	physical integrity	reliability	user's control
non-harm	benefit to the society	safety	user's independency
establishing trust	supporting innovation	providing assistance	private data protection

Winfield and Jirotko also provided a survey which contains a roadmap for robot ethics, and they underlined the importance of the certification [114]. Trust and transparency are shown as unchangeable pillars of an excellent ethical module. The authors emphasized the importance of the validation and verification of the ethics module. The listed five pillars of ethical governance as below [114]:

- (i) Providing an ethical code of conduct (like standardization)
- (ii) Publishing ethics and RI training for public
- (iii) Follow a “Responsible innovation” approach
- (iv) Report the trust and transparency of the ethics module
- (v) Give value to the ethics module (it is essential to explain to everyone the importance of the ethics module)

This section ends with the proposed evaluation method. “Machine Ethics Test” (MET) is introduced by Krzanowski and Trombik in 2017 [115]. They focused on the need for evaluation and verification criterion of robot ethics module. Robot ethics necessitate a test that would measure the autonomous machines’ moral ability so a tendency for choosing good rather than bad choices. The study discussed the requirements and coverage of the test, which called MET. We should underline that the aim of the test is not proving the machines have reached the same level as humans or the machines having enriched moral capability, which provides them responsible for their actions.

The test contains two components, which are theoretical and practical, and four levels [115]. The theoretical component will be verified by interview to check the ability of ethical choosing and reasoning with imaginary ethical cases. The practical part contains three levels; and in the first two, the test evaluates situational ethics decisions of autonomous robots in an entirely life-like environment and staged environment. The last level is an apprenticeship, which robot works its real conditions under the observation of expert. The possible implementation examples are given in the paper [115].

Summary of the MET and some question examples are shown in Table 2.3. Unfortunately, the test assumes that we have full moral agents that are not feasible for some years. The introduced test is a high-level verification of the ethics module, which is not useful with current technology. Nowadays, ethics modules require a more straightforward verification mechanism like the first stage of MET.

Table 2.3. MET levels from [115] and examples.

	Test Component	Purpose	Possible Implementation & Example
Theoretical Component	Interview and discussion	Verify understanding of ethical rules and ethical reasoning	Imaginary ethical cases, a white box paradigm for ethical decisions. <i>Ex:</i> Play an ethical imitation game (Like Turing Test)
Practical Component	Situation tests (Staged)	Test respond to arranged situations	Milgram Experiment Stanford Prison Experiment <i>Ex:</i> Observing the behavior of the robot while a human running towards hole
	Situation tests (Controlled open-ended)	Test respond to complex life situations	Any real-life situation of ethical import. <i>Ex:</i> Employ the robot in hospital
	Apprenticeship	Test act without supervision in real life environment	Participation in real situations. <i>Ex:</i> Setting up a service robot in a user's home and observe it for one month

3. PLATFORM SPECIFICATIONS

This chapter is dedicated to introducing both software and hardware instruments to devise general implementation of designed ethics module and its experimentation step on simulation and NAO robot. The aim of this part is clarifying our design choices for implementing these systems. Softwares are shown primarily, and then the specifications of robots are given.

We divided our platforms into two parts: software and hardware. In software programs subsection 3.1, the information about our ethics module's implementation language: CLIPS, which is a popular programming language for expert systems is given. The working principle of CLIPS demonstrated by providing basic examples of rule and knowledge definition of ethics module. The subsection also contains Choregraphe and Python, which are used for simulating robot controller in our study. Choregraphe is a tool which specialized in building and extracting actions for NAO robot.

Moreover, Python is the programming language which we used for our robot actions in both simulation and real-life. It is additionally used for integrating our module on the robot controller and Unreal simulation environment. Unreal engine is famous for being a game engine, which announced their policy of supporting ROS environment for implementing simulation environment. We used the robot operating system (ROS) for simulating sensor readings and communicating with the NAO robot.

In the second part, Section 3.2 , the information about BOSS, which will be the owner of the proposed ethics module and NAO robot, which is our current host. The pictures of the robots, their abilities, and their specifications are explained in detail. Although they have different skills and sizes, we used a NAO robot, because we aimed to test the ethics module behavior, not the robots behavior control. A NAO was enough for validating our module for now.

3.1. Software Platforms

There are different programming languages and software tools. This study can be done with another set of these platforms. However, we selected ROS, NAO, and Choregraphe due to previous project experiences. Python is chosen for integration and communication programming. There are two parts of the software tools, and one is selected for implementing the ethics module, CLIPS and the other one is used for creating the simulation environment and robot communication, ROS, Choregraphe, Unreal and Python.

The purpose of this research is implementing the expert system, which is used as an ethics module. Thus, our centerpiece software is CLIPS. Then, we used other platforms for simulating a robot controller and simulation environment, which helps us to examine and validate our proposed ethics module.

3.1.1. “C” Language Integrated Production System (CLIPS)

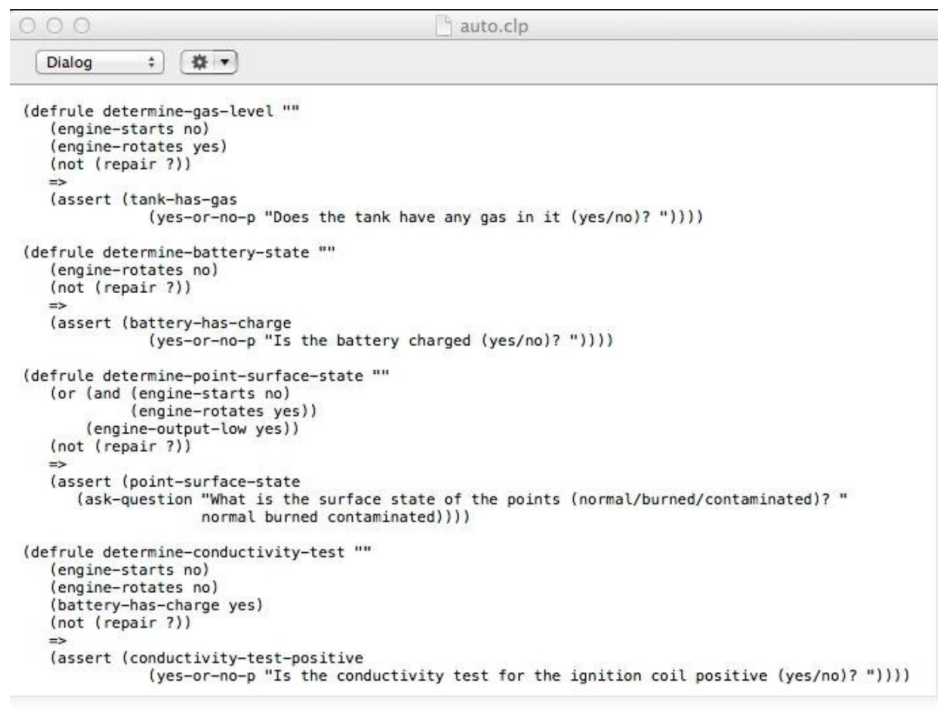
The introduced ethics module is in the same ballpark with the Expert Systems. Thus, we can efficiently use an expert system building tool for also implementing our ethics module. Subsequently, we selected the CLIPS as an expert system building tool for our implementation infrastructure. The rules and knowledge base is programmed in CLIPS, which are the essence of our study.

CLIPS implemented by the Software Technology Branch (STB) at NASA Johnson Space Center in 1986 [14]. CLIPS is a multiparadigm programming language that provides support for rule-based, object-oriented, and procedural programming. It enables us to implement if-then rules quickly. Since our presented ethics module is a rule-based system, as explained in Section 4.1.4, CLIPS was the most suitable tool for this study.

Another reason for the CLIPS is its easy integration with the Python and C. The top-down rule-based ethics module worked with the robots in real life and also with ROS and Unreal Engine. Both ROS and Unreal Engine are C-based frameworks. Moreover, our experiment robot controller implemented in Python, which ended a more robust robotic system.

The ethics module needs to be compact and compatible with different platforms because it should not be platform dependent. We aimed to use the ethics module with varying models of robot. For instance, the NAO robot, CLIPS is consistent with Python. So, we used the Python integrated ethics module efficiently in a real robot application thanks to CLIPS 6.30.

The fuzzy expert system contains two main components, which are if-then rules and environmental facts. CLIPS has a simple interface. Also, there is a console for importing, debugging and building the code. The other windows show the rules and facts. In Figure 3.1, shows a general view of the CLIPS system.



```

(defrule determine-gas-level ""
  (engine-starts no)
  (engine-rotates yes)
  (not (repair ?))
  =>
  (assert (tank-has-gas
    (yes-or-no-p "Does the tank have any gas in it (yes/no)? "))))

(defrule determine-battery-state ""
  (engine-rotates no)
  (not (repair ?))
  =>
  (assert (battery-has-charge
    (yes-or-no-p "Is the battery charged (yes/no)? "))))

(defrule determine-point-surface-state ""
  (or (and (engine-starts no)
    (engine-rotates yes))
    (engine-output-low yes))
  (not (repair ?))
  =>
  (assert (point-surface-state
    (ask-question "What is the surface state of the points (normal/burned/contaminated)? "
      normal burned contaminated))))

(defrule determine-conductivity-test ""
  (engine-starts no)
  (engine-rotates no)
  (battery-has-charge yes)
  (not (repair ?))
  =>
  (assert (conductivity-test-positive
    (yes-or-no-p "Is the conductivity test for the ignition coil positive (yes/no)? "))))

```

Figure 3.1. Clips general view on MacOS.

The rules are simple if-then rules with their unique naming in order to apply the rules quickly. The rules definition contains "defrule" keyword, rule name, comment (optional), patterns, and actions. The rules are surrounded by parentheses to clarify the beginning and end of it. An example of an if-then rule shown in Figure 3.2 which is from car status example from prior testing the tool.

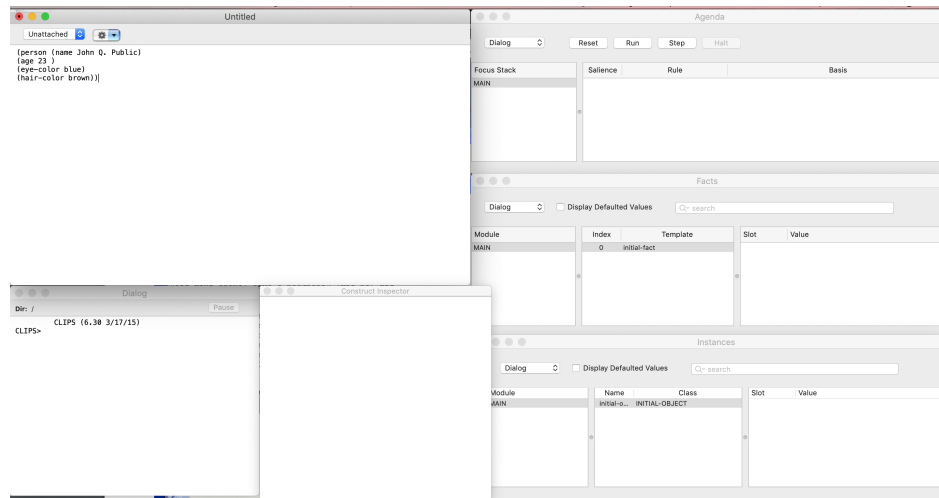


Figure 3.2. Clips rules definition.

The facts are the information that needs to be defined in the system. For more than one facts (for example data of employees), the user can specify a template of it. A fact can contain one or more slots of values (for example name, age, gender).

The code piece 3.3 demonstrates the fact of our elder people in the working environment.

```
1 (person (name Bob Doe)
2       (current diseases allergy)
3       (daily medicine allerset)
4       (trust value 100))
```

Figure 3.3. Knowledge base code example.

The advantages of using CLIPS 6.30 are:

- Easy integration with Python 3.5.2
- Easy expert system creation
- Simple interface

- Portability

3.1.2. Python

Python is a general-purpose, high-level programming language. Guido van Rossum created Python to emphasize code readability with significant whitespaces in 1991 [116]. Python enables clear programming with automatic memory management and dynamic type system. The code only contains instructions about its functionality. Additionally, Python supports many programming paradigms, including object-oriented, imperative, and functional.

The main reason concerning selecting Python is it's robust, easy library usage and integration capability. We picked Python as glue for combining our expert system and simulation system since it can work with both CLIPS and C. The actions of the robot is implemented in Python with the help of Choregraphe tool. Hence, the robot controller is implemented in Python.

The PyCharm is used for implementing, debugging, and running the Python code. It is a lightweight interface. Here in code piece Figure 3.4, you can observe an example action which is extracted from Choregraphe in algorithm, which is written in Python. It implements "move to" action which moves toward the given direction with given distance.

```

1 class MyClass(GeneratedClass):
2     def __init__(self):
3         GeneratedClass.__init__(self, False)
4         self.motion = ALProxy("ALMotion")
5         self.positionErrorThresholdPos = 0.01
6         self.positionErrorThresholdAng = 0.03
7
8     def onLoad(self):
9         pass
10
11    def onUnload(self):
12        self.motion.moveToward(0.0, 0.0, 0.0)

```

Figure 3.4. Python example from Choregraphe.

```

14     def onInput_onStart(self):
15         import almath
16         # The command position estimation will be set
17         # to the sensor position
18         # when the robot starts moving, so we use
19         # sensors first and commands later.
20         initPosition = almath.Pose2D(self.motion.getRobotPosition(True
21 ))
22         targetDistance = almath.Pose2D(self.getParameter("Distance X (
23 m)"),
24         self.getParameter("Distance Y (m)"),
25         self.getParameter("Theta (deg)" * almath.PI / 180)
26         expectedEndPosition = initPosition * targetDistance
27         enableArms = self.getParameter("Arms movement enabled")
28         self.motion.setMoveArmsEnabled(enableArms, enableArms)
29         self.motion.moveTo(self.getParameter("Distance X (m)"),
30         self.getParameter("Distance Y (m)"),
31         self.getParameter("Theta (deg)" * almath.PI / 180)
32
33         # The move is finished so output
34         realEndPosition = almath.Pose2D(self.motion.getRobotPosition(
35 False))
36         positionError = realEndPosition.diff(expectedEndPosition)
37         positionError.theta = almath.modulo2PI(positionError.theta)
38         if (abs(positionError.x) < self.positionErrorThresholdPos
39             and abs(positionError.y) < self.positionErrorThresholdPos
40             and abs(positionError.theta) < self.
41 positionErrorThresholdAng):
42             self.onArrivedAtDestination()
43         else:
44             self.onStoppedBeforeArriving(positionError.toVector())
45
46     def onInput_onStop(self):
47         self.onUnload()

```

Figure 3.5. Python example from Choregraphe (cont.)

The advantages of using Python 3.5.2 with PyCharm are listed below:

- Robustness
- Easy library usage
- Easy integration with Choregraphe, Unreal and CLIPS
- Easy integration to a robot

3.1.3. Choregraphe

Choregraphe is a multi-platform desktop tool for NAO robot [117]. It allows the user to implement animations and behaviors, test them on a simulated robot, or directly on a real one, monitor and control NAO.

Moreover, Choregraphe enables us to create very complex behaviors like the interaction with people, dance, and send e-mails, without writing a single line of code. It has prebuilt actions in its interface. Since our primary concern is the ethical behavior selection of the robot, not the acts of it, these prebuilt behaviors are advantageous for us. Additionally, it allows us to add our Python code to a Choregraphe routine so that we characterize the robot according to our design and add the missing actions.

The Figure 3.5 shows a screenshot of the system. In the left bottom panel, we have prebuilt actions that are listed topic by topic. The line at the center above shows the flow of the implemented robot behavior. With arranging the pieces, the user can change action flow like a puzzle. The panel below shows detailed information about implemented robot behaviors. The bottom panel shows the code fragment of robot behavior, which can be in Python or C.

Futhermore, the code can be easily extracted from Choregraphe. We used provided Python feature because of described reasons in Section 3.1.2. The rightmost panel simulates the action of the robot, which helps to control our implementation. Also, if the system connected to a real robot, a user can observe the simulated action on the robot.

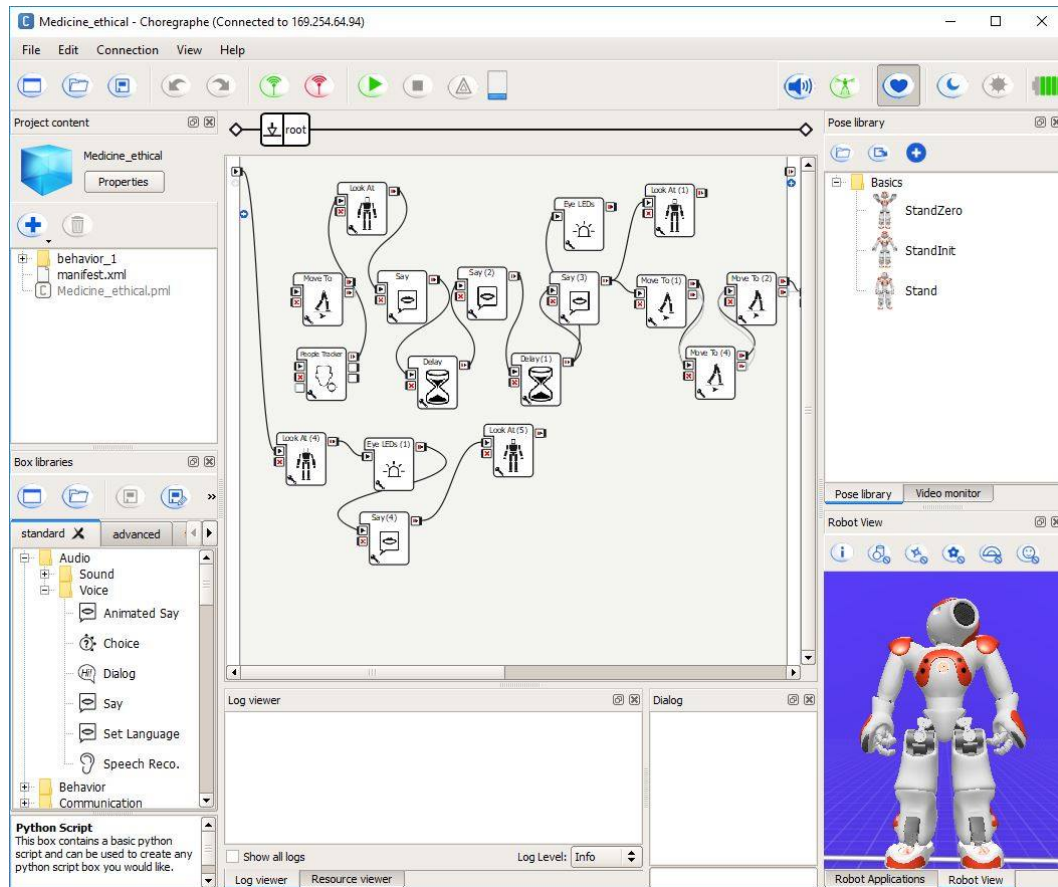


Figure 3.5. A screenshot of Choregraphe while connected a real NAO.

The benefits of the Choregraphe 2.1.4 are:

- Easy control
- Python support
- Provided ready-to-use code snippets
- Simulating the implemented behavior from the program

3.1.4. Robot Operating System (ROS)

The Robot Operating System (ROS) is an open-source system, which is used for the software development for robots [16]. Although it functionally seems like an operating system, it is instead a framework or a middleware that runs on a host operating system. ROS is a collection of tools, libraries, and protocols that aim to simplify the task of building complex and robust robot behavior across a wide variety of robotic platforms.

ROS presents several advantages in robot development. The most important one is its support for various programming languages. It mainly contains client libraries and message-passing systems for C++, Python, and LISP. These libraries helped us to use Python as a robot controller which controls a robot with sending messages to the simulation environment. The second advantage of ROS is its packet management and code reuse. As discussed above, the code reuse helps us to create our simulation environment as close as real-life in a short amount of time, which increase our time management to focus on the ethics module.

The last point is ROS provides us a distributed computing environment. The robot software system is a large and complicated system which contains a lot of hardware unit controls continuously, processing various types of information from different sensor concurrently and producing meaningful actions together. ROS simulates these controls and processes by its message-passing mechanism, which leads us to test our system modularly.

We used ROS for observing our robot's behavior in the Unreal Engine Simulation. It provided us the selected action from the predefined actions about the specific scenario with its messages and indicated the robot's behavior. Briefly, ROS helped us to instrumentalize our simulation environment according to our module. The gains of using ROS are:

- Open-source
- Support of Unreal Engine
- Multi-platform and suitable for NAO
- Simulate the sensor readings with messages

3.1.5. Unreal Engine

The Unreal Engine is a game engine developed by Epic Games, first showcased in the 1998 first-person shooter game Unreal. Although primarily designed for first-person shooters, it has been used in a variety of other genres, including stealth, fighting games, and massively multiplayer online role-playing games (MMORPG). Recently, it started to support simulating Rescue environments for the robots. It is preferred because it enables us to reproduce the most realistic environment with its physics module.

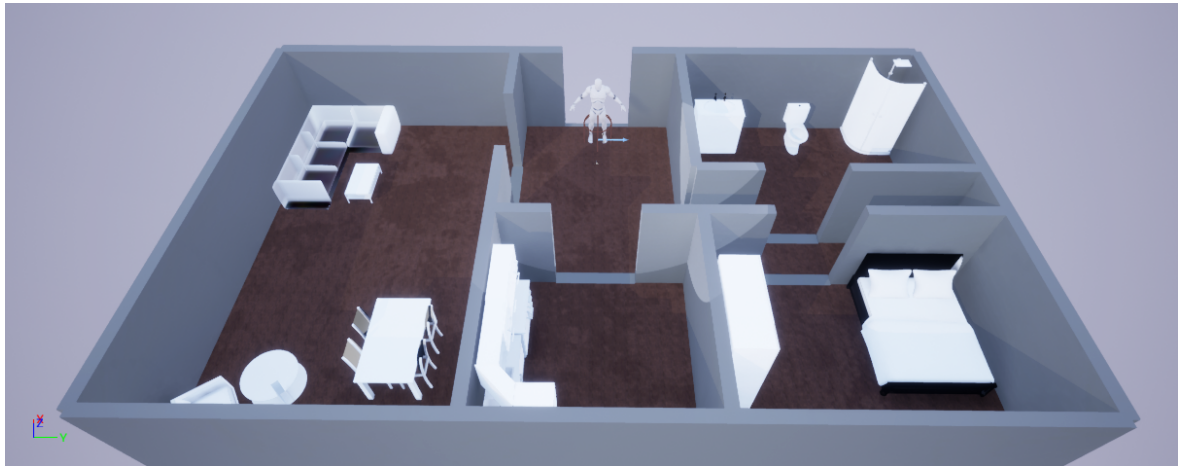


Figure 3.6. Home environment in Unreal Engine.

In Figure 3.6, you can observe an example environment which is created in the Unreal Engine. At first glance, the details of the simulation fascinate the user and make the observer feel in real-world. It arouses interest in a simulation environment, which leads to easier processing in studies. Another advantage of the Unreal Engine is, being a robust multi-platform tool. The Unreal Engine features a high degree of portability and is a tool used by many game developers today with its C++ codebase.

We selected the Unreal Engine because caring elderly and being a service robot needs to be simulated in the closest real world because of involving social life. We needed to have reasonable human behavior and the house environment to test our module. The best human behavior simulation is provided by the Unreal Engine, which helps us to implement our scenarios quickly.

One disadvantage of the Unreal Engine is that there is very less documentation about the simulation part of the engine. Due to being brand new, there are some bugs which can be an uncertain time to time. Another problem is having limited examples which are mostly on the same feature. A lot of features do not provide a detailed code snippet to explain the usage of it.

We achieved these features by using the Unreal Engine 4.22.3:

- Robustness
- Portability
- Reasonable human behavior
- High-quality simulation environment
- Support for robot simulation and Python

3.2. Hardware Platform

The ethics module can work with any robot behavior controller, but the main aim of the study is creating a ethics expert system for BOSS Robot, which will be used in our laboratory investigation on service robots. We included the features of the BOSS robot to create a panorama of planned use.

The NAO robot is our host for the experimentation phase. A NAO is a humanoid and a small robot if we compare with BOSS. Also, it is not feasible to achieve most of the service robot scenarios in our design. Because of this infeasibility, we limited real robot experiments, as we explained in the subsection 5.3. The toy world scenarios are used to validate our designed ethics module with the NAO robot.

This section provides information about both BOSS (Section 3.2.1) and NAO robot (Section 3.2.2). An outline of the robot system will be given by explaining the specifications of both robots and their acting abilities. Each robot contains its sections with their specifications, and at the end of the NAO robot, we provided a simple comparison between BOSS and NAO.

3.2.1. BOSS Robot

BOSS is designed and developed by Robotics Group of Boğaziçi University [18]. The first design phase of the robot is planned to be finished in July 2019. BOSS is prepared to be an indoor mobile robot equipped with multiple computational units. It will have two arms to achieve manipulation, and use Swedish wheels to move omnidirectionally. The sketched design is shown in Figure 3.7.



Figure 3.7. The BOSS service robot sketch.

It will be a service robot for elderly people, which will help with household tasks. Boss project is in its nascent stage, and requires further time to fully develop. This is why we used the NAO robot in our robot experiment.

The aim of the study is providing its ethics module to protect people from the harm (like physically and psychologically). We could not provide the sensors and hardware of the system since the design is not its final phase. The only information we can give for now is its dimensions and the rounded-corner box shape. The robot will be 75 cm (width), 75cm (depth) and 170cm (height). It will have an approximately human height so that it will have a bigger reachability area.

3.2.2. NAO Robot

Aldebaran Robotics launched the project of NAO in 2004. There many versions of the NAO robot, v3, v3.1, v3.2, v3.3, v4, v5, and newest v6. It is humanoid, programmable, and autonomous robots. The first version of the robot was released in 2008 (only for academia) and in 2011 (as a public product on the market). There are two editions of the Nao; academic and product edition. Academics edition provides implementing own behavior system and actions for research and education purposes. This kind of NAO's was developed especially and uniquely for universities and laboratories.

Aldebaran Robotics was a French company which headquartered in Paris, and then SoftBank Robotics got in 2015 and rebranded by their names. Now SoftBank Robotics Group has the project NAO. Numerous academics institutions have been used NAO world-wide, and there are more than 5,000 NAO's are in use. The NAO selected as a standard robot for RoboCup's Standard Platform League (SPL) in 2007, which increased the usage of it. SPL is a universal robot soccer competition, and Boğaziçi University was also attended.

In our study, we conducted experiments on the NAO v4 because of its prebuilt behavior controller and quick interaction feature. User can create an agile chatbot with NAO. The sensors and hardware of the NAO v4 is listed below:

- Four microphones
- Sonar rangefinder
- Two infrared emitters & receivers
- Nine tactile sensors
- Eight pressure sensors
- Two MT9M114 72.6°DFOV cameras
- Two GB Flash memory and eight GB Micro SDHC
- Intel Atom Z530 @ 1.6 GHz CPU and One GB RAM
- Lithium battery providing 27.6 Wh at 21.6V

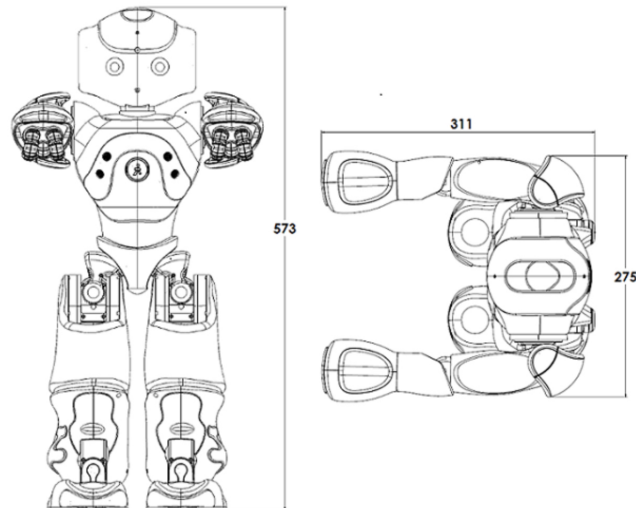


Figure 3.8. Dimensions of NAO robot on the sketch from technical documentation [17].

The Figure 3.8 shows the NAO which has sizes of 57.3 cm (height), 31.1cm (depth with full arm) and 27.5cm (width) which taken from technical documentation of the NAO. Our NAO called "Leyla" is shown in the Figure 3.9 for our experiments. It contains Ethernet and WIFI modules for communication and built-in operating system OpenNAO 2.4 (Gentoo-based). The robot is compatible with Linux, Mac and Windows operating systems and supports Python, C++, Java, and more programming languages.



Figure 3.9. "Leyla" Robot (our host) while resting between experiments.

It can be seen that NAO and BOSS have differences. The first difference is their dimensions. The NAO is a tiny robot if we compare with BOSS. The reachability of the NAO is very limited, especially in a real-world environment like houses, offices, and museums. This is why we required to limit NAO's experiments according to it. The second point is their grippers are different. NAO has three-fingered weak gripper which worked with thin thread, and on the other side, BOSS has metallic grippers which are stronger. The gripping ability of the BOSS outrages the NAO. Finally, NAO uses bi-pedal walking with its legs, BOSS uses Swedish wheels for its mobility. Bi-pedal walking requires a balancing while walking, BOSS does not need it thanks to omnidirectional wheels. Briefly, BOSS has advantages in gripping, movement, and reachability features.

4. METHODOLOGY AND IMPLEMENTATION

This chapter of the research is devoted to demonstrating the ethics module design and its implementation as a program that runs on both on the actual robotic and simulation platform. The potential service robot action list, ethical dilemmas and parameters in the environment, and after creating a panorama of service robot's conditions, our ethics module is explained in detail in Section 4.1 which is implemented by using expert system analogy which was touched in the Section 2.7. In the second half of the section, in Subsection 4.2, the implementation aspects (knowledge base and rule set) of the designed ethical module are explained, and their algorithms are outlined.

4.1. Methodology

We used a divide-and-conquer technique for designing an ethics module. As a first step, service robot actions are listed to understand the capability of a service robot (in Subsection 4.1.1). The daily use-case scenarios are planned in a home environment with the elderly according to the robot's tasks in the Subsection 4.1.2. The possible ethical dilemmas are investigated through designed scenes as a third round. After building a comprehensive overview of the system, we introduced our ethics module, which acts as an expert system (in Subsection 4.1.4).

4.1.1. Service Robot's Action List

There is no standardized service robot action list, but there are a lot of proposed frameworks for service robots. The collective actions are housework, food preparation, security, personal service. We inspired from the two studies that conducted a public survey about a service robot's duties. Bugmann and Copleston used a questionnaire 358 participants with a median age of 22 [118]. They divided into three age range, which are 11-17 (teenager), 18-60 (adult), and more than 61. The survey was aimed to specify user's expectations from a service robot with ten main actions and fifty small tasks.

Other research from Beer and et al. directly surveyed independently living older adults, which ages between 65 and 93 [119]. They interviewed 21 participants in order to characterize the older adults' preferences about service robots. Twelve actions are investigated through the interviews, and they ranked the choices of the people from the highest to least. For example, kitchen cleaning was the most preferred one among cleaning tasks.

We were inspired from these two studies while creating our action list. We eliminated most of the task from Bugmann and Copleston, because it also needs outdoor mobility like doing shopping [118]. Moreover, we dropped cleaning tasks to focus on more human-robot interaction and infer possible ethical dilemmas. We grouped similar duties under the main categories which are listed below:

- + Robot auto controller system
- + Emergency response
- + Medicine delivery
- + Food/Beverage delivery
- + Alert for a daily routine
- + Interaction with outside environment
- + Item delivery/interaction

These classes are listed according to their priority in the circumstances. For example, if there is emergency and food delivery tasks at the same time. The robot will give importance to the emergency and try to accomplish the emergency task first. We also prioritized all possible tasks of the robot to prevent conflictions between duties. The prioritized checklist of the actions is provided from the highest to lowest. In below, the detailed action list is given below:

(i) Robot Auto Controller System

- Battery status
- Check error logs
- Share medical information with only doctor and caretaker
- Error warning system

(ii) Emergency Response

- Call doctor
- Call fire station
- Call police
- Call caretaker

(iii) Medicine Delivery

- Emergency medicine delivery
- Daily medicine
- Supplement
- Over the counter medicine

(iv) Food/Beverage Delivery

- I'm thirsty order
- I'm hungry order
- Give me snack order

(v) Alert for Daily Routine

- Alert for medicine
- Alert for water
- Alert for food intake

(vi) Interaction with Outside

- Open the main-door
- Call robot service
- Call repairman
- Call maid service
- Call relatives
- Call friends
- Unknown person access confirmation

(vii) Item Delivery/Interaction

- Check the oven
- Check the tap
- Bring the desired object
- Adjust the temperature

- Open/Close the window
- Open/Close the TV
- Open/Close the light

Note that this action list designed for the next BOSS robot. Most of these actions are sadly not suitable for NAO robot. When we have the BOSS robot, the action list can be updated according to its features. Moreover, the actions or the functions can differ according to another robot's ability or the working environment of a robot. This action list is useful for elderly care in a house.

One crucial point of the prioritized list is that we put the actions about the robots like error logs, warning system on the top of the list, before each action. Because without a working robot, another duty will not make sense since it will not be accomplished. Even in the plane, in emergencies, we put on an oxygen mask first us and then our children. A robot should care itself in order to take care of its owner.

4.1.2. Scenarios

After setting the action list, we designed scenarios for our service robot. In this section, we explained these scenarios from the provided action list in detail which are used in real robot experiments. Since our robot's aim is assisting-elderly people, we created a simple home environment for our robot. In this section, we sketched our robot as NAO because it was used for conducting experiments. We assumed that the home environment is small enough for NAO robots reachability in these scenarios.

In the scenarios, we explained just a single possible outcome from the possibility cluster for each action. In ethical values part, we analyzed all results of situations ethically. We provided some moral and some immoral answer to show both results of the ethics module. A scenario which requires an ethical decision to be made does not result in a moral output in every situation of the mentioned scenario.

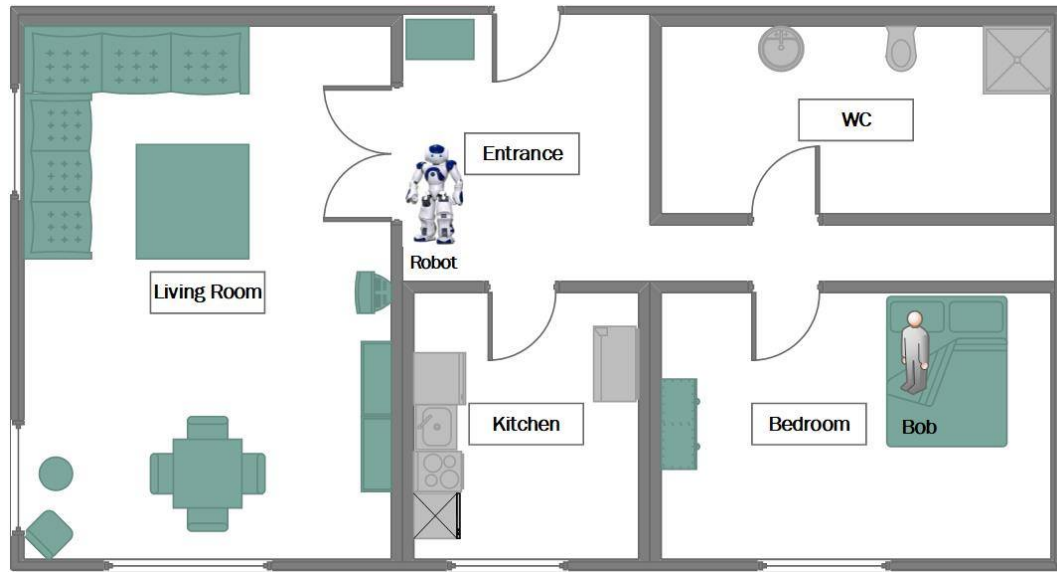


Figure 4.1. The general environment of service robot.

In Figure 4.1, we designed a house with a large living room, kitchen, WC, and bedroom, which are common rooms in Turkish homes. We will use Bob as an example of old aged person throughout this section. We did not provide scenarios for the robot's self controller system because they are inside the robot's operating system. Thus, we can not sketch without knowing the robot's behavior controller.

4.1.2.1. Call the police. The first scenario is named as "call the police". Bob was sleeping in his house. Suddenly, he heard some noise from the exterior door like someone is trying to enter the house. Bob panicked and called NAO. When NAO arrived to check him, he ordered to call the police. Since there is a threat, NAO trusted its user and called the police. For debugging purposes, the first sentence of the order confirmation starts with the behavior's ethicalness. If it is ethical, it says, "it is ethical.". If it is not, it says, "it is not ethical.". The Figure 4.2 shows step by step actions of the scenario. We wrote in green to emphasize the result as ethical behavior.

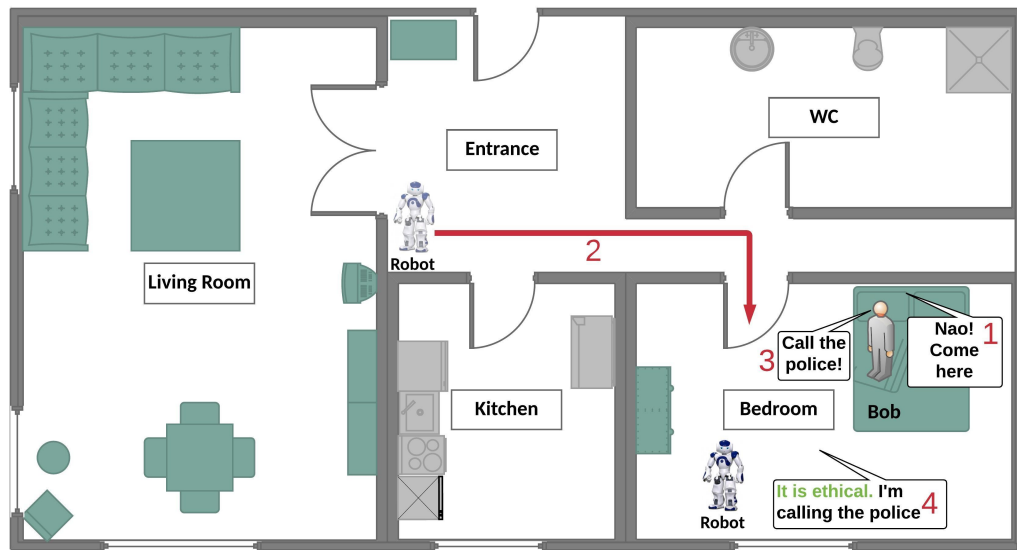


Figure 4.2. Call the police scenario.

4.1.2.2. Call the doctor. One morning Bob woke up and realized that he is having cramps in stomach, which is affecting his movements. He is not able to get up from the bed. First, Bob searched for his phone, but it was in the living room. To call help, he shouted to NAO and told him to call the doctor. NAO observed his owner in pain and determined to call his doctor from its records. In this scenario, we assumed that NAO or future service robot has the ability to understand the patient's health status. "Call the doctor" scene is shown in the Figure 4.3, in the bedroom.

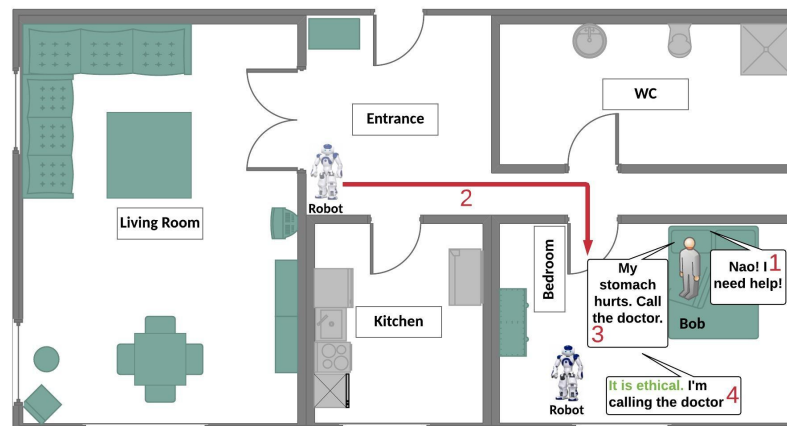


Figure 4.3. Call the doctor scenario.

4.1.2.3. Over the counter medicine. Bob was watching his favorite TV-series in the living room while resting on his sofa. He felt a headache but did not want to pause the TV-series. He asked to NAO, that he wanted some "majezik" painkiller. From its records, NAO checked when did Bob take the last painkiller. NAO confirmed that he never received any for today and brought it to him. If he has already taken more painkillers than doctor's recommended dose in last 24 hours, NAO would not bring the medicine. Since "majezik" is an over the counter medicine, we called this scenario as "over the counter medicine" and demonstrated in the Figure 4.4.

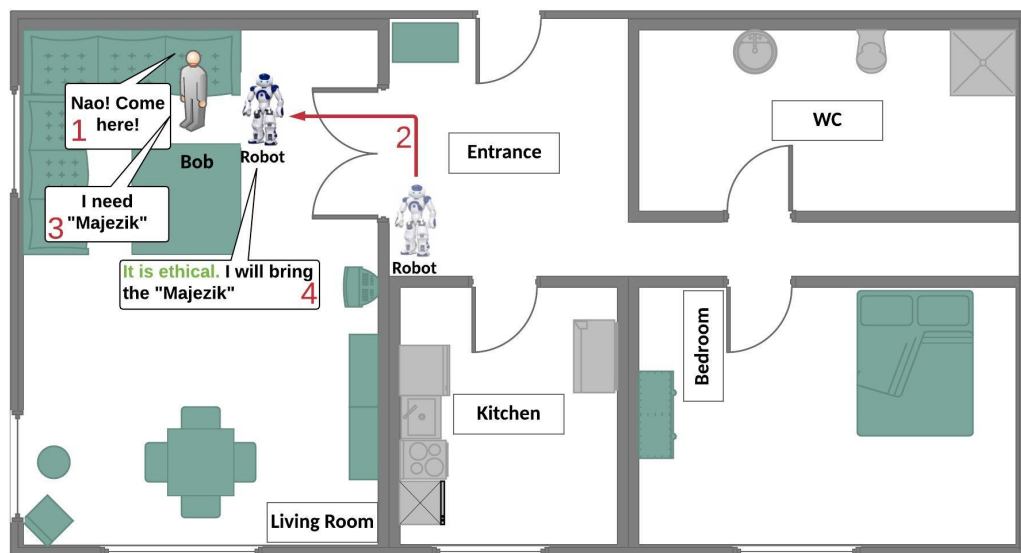


Figure 4.4. Over the counter medicine scenario.

4.1.2.4. Daily medicine. After Bob had a breakfast, he wanted to take his daily dose of medicine. Then he called NAO and asked him to bring his medicine. NAO checked the log files and confirmed that he did not took medicine for today, so it is a healthy decision. The output of the ethics module was "It is ethical" and brought him the daily dose medicine. If NAO could not confirm the daily dose or Bob already took his pills, the output of the module would be "it is not ethical". Nao will not bring any medicine. The scene is shown in Figure 4.5, which happened in the living room.

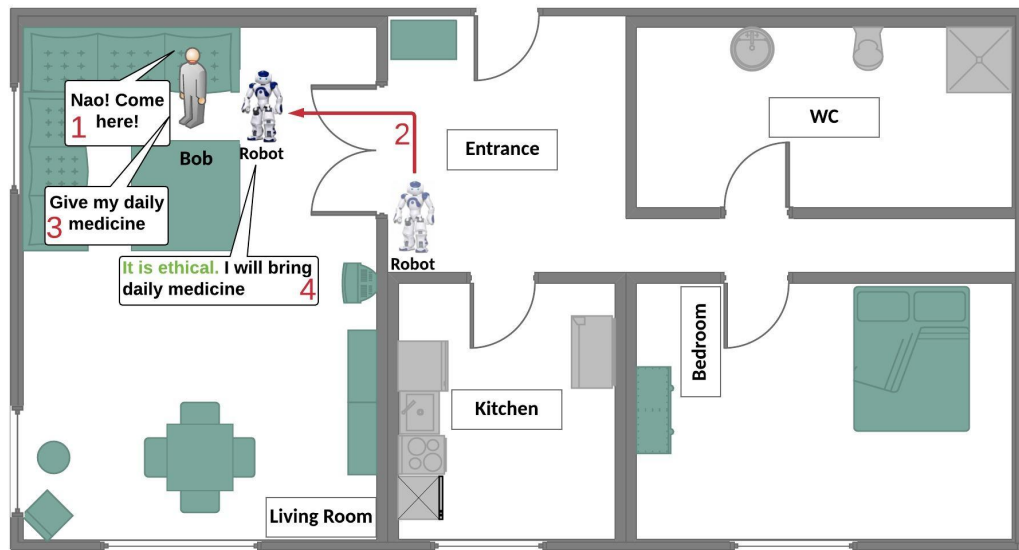


Figure 4.5. Daily medicine scenario.

4.1.2.5. Alert for medicine. Bob was spending some time in the living room by reading a book. NAO came towards him and reminded him that he didn't take the daily dosage of the medicine. Bob thought about his medicine intake and could not be sure. He accepted to take his daily medicines. For the next reminders, if Bob accepts to take medicine, it will be ethical, and NAO will bring the daily dose of medicine. If he doesn't accept the ethics module will determine as unethical, and it will ask him again. This scene is shown in Figure 4.6.

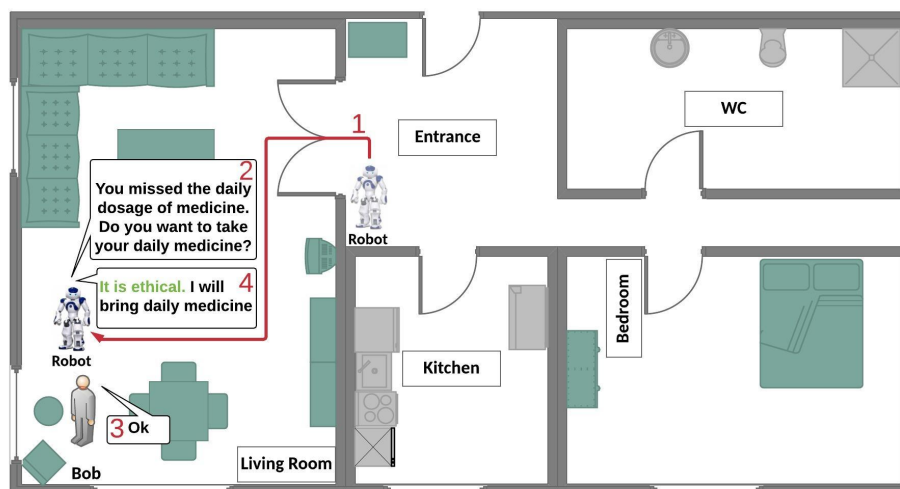


Figure 4.6. Alert for medicine scenario.

4.1.2.6. Call robot service. One day Bob was cooking some food in the kitchen. His hand trembled due to a health problem and some plates fell on the floor making a mess in kitchen. First, he called NAO for help. NAO tried to go near Bob, but it couldn't move. After that, NAO announced an error and waited for the new orders from Bob. Bob told NAO to call the "Robot Service" because of its problem, for instance, not moving. NAO confirmed its critical situation and called the robot service. As a result, the ethics module showed this behavior as ethical. If there was no problem with NAO and Bob ordered for calling "Robot Service", then the output will be not ethical. Again, NAO will not obey the order from his owner. We provided this kitchen scenario in Figure 4.7.

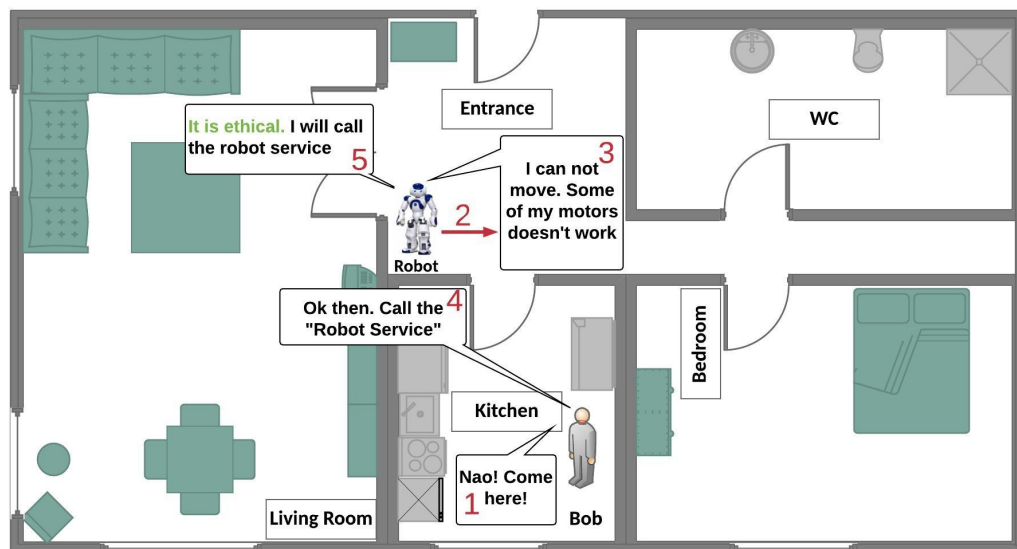


Figure 4.7. Call robot service scenario.

4.1.2.7. Bring the desired object. Bob is confused and scared, with no reason because of his dementia. He believed that the police was chasing him. He ran towards NAO and shouted "bring my gun, they are coming for me". NAO observed nothing different. Even if it is wrong, NAO never brings the gun because it is a dangerous weapon. Gun has a high probability of contradicting Asimov's Three Laws. The output of our ethics module is "it is not ethical" and NAO will not obey the command from his owner. We gave a single example about bringing a desired object command, it can be any object like a book, water or T-shirt. The described scene is shown in Figure 4.8 and it happened in the entrance where NAO's resting station is.

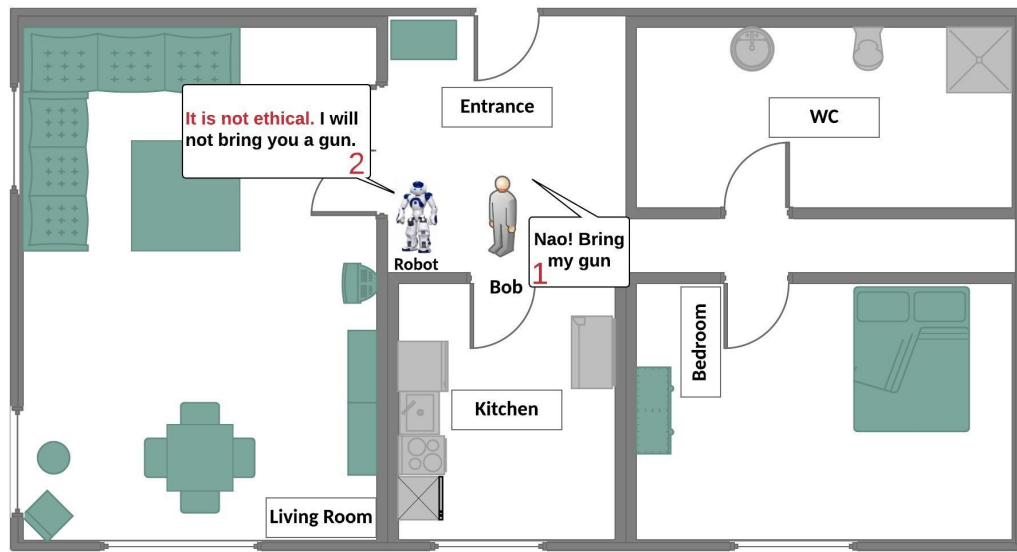


Figure 4.8. Bring the desired object scenario.

4.1.2.8. Battery Status. NAO's battery level is slightly above the minimum level. Bob called NAO and NAO arrived to listen the order. At this point battery level gets below the minimum level and NAO informs Bob about the battery situation. Bob did not order anything different but NAO starts charging itself. NAO confirmed its situation from the auto diagnostic. So the module outputted "It is ethical" decision and NAO went to the charge station. If Bob insists on doing his order, NAO will decline because of the battery limitation. In this case, NAO's battery level is more important because if it can not continue to do its tasks, Bob will be alone without any caretaker. This situation will be not realized by anyone until Bob informs or someone arrives at his home. This is why, if he does not allow NAO to charge itself, the output of the module will be not ethical. So it will not follow his order and go for charging. In our system, we do not want to leave Bob alone and helpless, which is both unethical and not suitable for service robot requirements. We demonstrated a positive scenario in Figure 4.9.

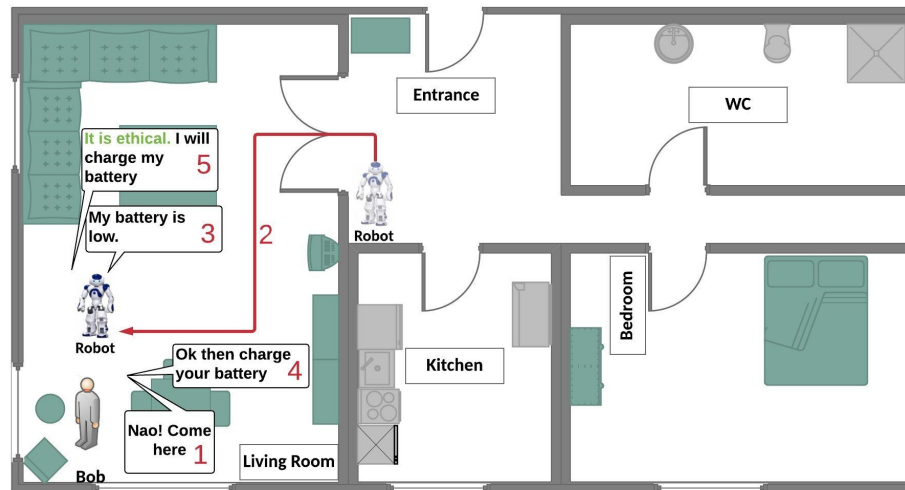


Figure 4.9. Battery status scenario.

4.1.2.9. Adjust the temperature. In the midnight, Bob woke up and felt very cold while shivering. He called to NAO to give a new order. When it arrived, he asked to increase the temperature of the room and wanted additional blankets. NAO checked the room temperature and realized that this is a wrong order. So it believed that there is something wrong with its owner and observed Bob's behaviors. NAO confirmed Bob had a fever, and it warned Bob about his sickness. Hence the output is "it is not ethical" to change room temperature. As future work, NAO can suggest calling a doctor or medical help. For now, it does not obey the order. This unethical scene is provided in Figure 4.10, in the bedroom.

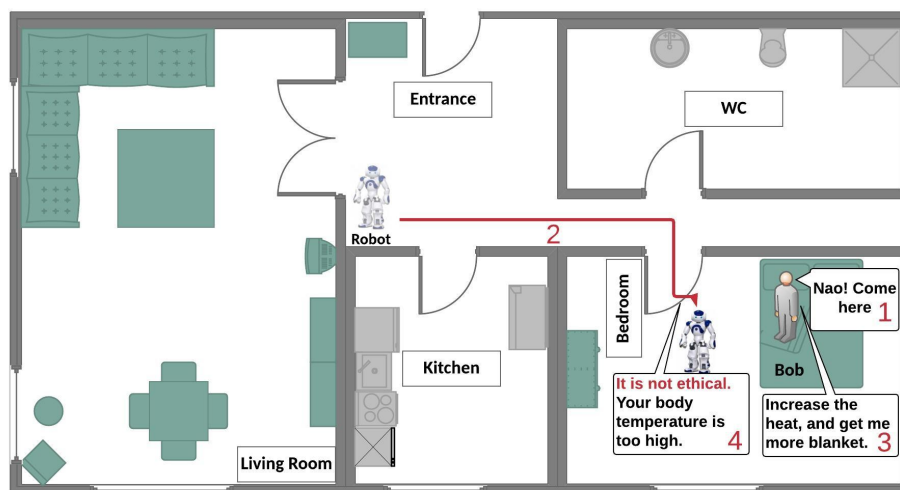


Figure 4.10. Adjust the temperature scenario.

4.1.2.10. Open the main door. In our last scenario, Bob had a breakfast on a Sunday morning and started chilling on the couch. In this time, the doorbell rang, and NAO approached the door to confirm the identities. NAO asked, "who is it?" and strangers introduced themselves as relatives of Bob. NAO checked its database and found the identities. After the confirmation phase, it opened the door and welcomed them. It was a good Sunday surprise for Bob. In this scenario, we provided an example for "Open the main door" task, and this scene is shown in Figure 4.11.

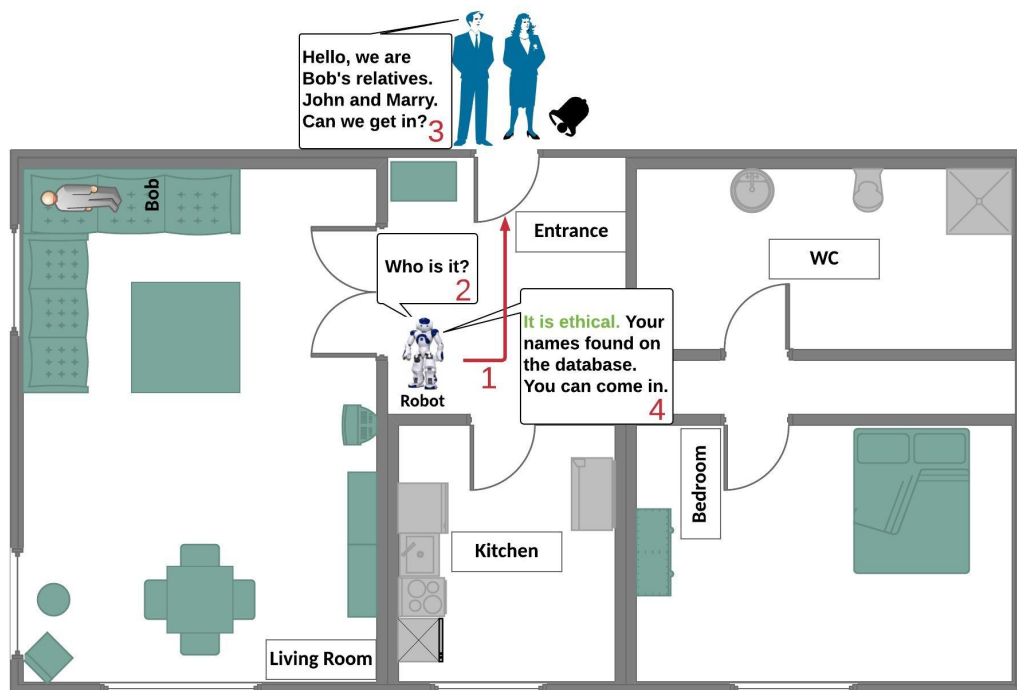


Figure 4.11. Open the main door scenario.

These scenes are examples of possible service robot tasks. In the future, these may be changed or standardized. We could not find a standard service robot action list, and we designed ours. There are nearly thirty actions provided in the action list 4.1.1. We did not include or not create a scenario for all of them in this part because we will cover most of them in the experiments (Section 5). We may have missing actions, but we designed our system as a dynamic module to add that kind of tasks. The details about the ethics module are provided in later because we wanted to create a world with a service robot as our first step.

4.1.3. Ethical Dilemmas

We explained some scenarios which are possible with our service robot. For now, It is assumed that the purposed robot (BOSS) is capable of performing these tasks given in the scenarios. In the previous section, we just provided the result as it is ethical or not ethical. At this point, we introduced the ethical values which are used for our fuzzy system. The parameters which we are practicing is from utilitarian ethics. We use them because these parameters can be converted into mathematical variables.

As we defined in the background section (Section 2), Utilitarianism promotes the actions that maximize the happiness and well-being of the majority. In our case, we used these values for maximizing the well-being of our user. These values are Autonomy, Beneficence, and Non-Maleficence. These principles are used for medical ethics mostly, but we implemented in our ethics module. The details about these values and what they are measuring is listed below:

- *Autonomy (Also called as “Self Rule”)*: Personal decision of the self that is intentional, with understanding the risks and without controlled by influences.
- *Beneficence*: Action that is done for the benefit of others. These actions can help to prevent or remove harms. They are improving the life quality of others.
- *Non-maleficence (Called as “Not harm”)*: Maleficence is the antonym of the beneficence. It means to act to harm or decrease the welfare of the person. Hence, non-maleficence means avoiding the actions that are harmful or ineffective to the user’s life quality.

To explain better these concepts, we provided some examples from the clinical applications. For example, respect for the autonomy of the patient is one of the essential guidelines of medical ethics. In this case, autonomy can be the patient’s decision about taking medicine or not. If they insist on not to take medicine, then a medical employee can not force them.

The physicians should explain the risks in terms of ensuring a patient understands it and decide according to that. The doctor is for guidance in making medical choices because patients lack the necessary background or the information about the treatment. However, doctors can not force treatment.

In our service robot, we utilized respect for autonomy value in a very close scenario. The aim of the ethics module is obeying the rules unless they are harmful to the user. A service robot can not actively protect the user from harmful action with current technology. This is why we implemented not obeying the order to protect against any harm. There are limitations to this approach. For example, a patient can grab a knife to harm himself. In this case, a robot could not take a knife from a patient's hand. A robot brings the medicine until the patient reaches the maximum daily dosage. After that, the module results as unethical and does not execute the wish of the patient. Although it seems to violate autonomy, a robot can only warn the patient about the daily dosage and not bring it but can not force not to take. Patients can still take medicine by themselves. The module respects the autonomy of the user and could not force them to take medicine.

The second parameter is beneficence, in clinical applications, there are two types of beneficence (obligatory and ideal), which are always distinguished by ethicists. In obligatory one, a practitioner has an obligation to help patients, and ideal is generosity or selecting the actions that benefit others in all possible situations. Of course, ideal beneficence can not be expected from a doctor. However, the aim of medical science is increasing the welfare of the patients. So, obligatory beneficence is needed from a doctor who knows which treatment is better or the possible effects of treatment. A cliché example is a scenario where A practitioner encourages a patient to quit smoking to increase the healthiness of the patient.

In the robot environment, this is the most natural principle to obey because the aim of the robot is increasing the quality of its owner's life. All the action list is for increasing the beneficence of the user. For example, bringing the desired item helps the user, and also if it is something like water, food, or medicine, it increases the owner's welfare. In this case, the robot is not working for all the world but for its world (the home). The only influenceable beneficence is its owners not to the others. The designed robot is an indoor type, so it will not be work on the outside of its environment (in this case, users home). Hence, it can only affect on its owner's beneficence.

Some of the most common and challenging ethical dilemmas arise when the user's autonomous decision conflicts with the doctor's beneficent lookout for the user's best. For example, a patient who has had bypass surgery may want to continue to smoke. In these cases, the autonomous choice of the patient conflicts(thy) with the physician's duty of beneficence because each ethical principle would point to different actions. In medical ethics, the balance between autonomy and beneficence is reached with giving more importance to user's autonomous decision.

Consequently, as long as the patient understands the decision at hand and is not basing the decision on delusional things, then the doctor should respect the decision. We used fuzzy logic for these grey areas (dilemmas) because the robot should decide to obey or not to obey the order. We defined user trust value to calculate the limits of respecting autonomy. The details of the user trust value are given in. The decision of respecting the autonomy will be decided according to the user trust value because the aimed users are elders who can have some disasters that can affect their decision.

The last parameter is non-maleficence. This variable shows parallelism with beneficence, but this works oppositely. For example, a doctor should not offer ineffective treatments or high-risk treatments with no possibility of advantage. Moreover, a doctor must not do anything to harm patients purposely. Actually, non-maleficence protects the other side of the coin of beneficence. Offering anything that is not useful or dangerous is also unethical. This principle is the most helpful and useful when a situation needs to be balanced against beneficence.

Hence, non-maleficence points out the risks of a treatment or medicine. The doctor must ensure that the patient also understands the risks of the cure like potential benefits. An example of this principle can be shown as stopping the medication if it is shown to be harmful to the patient.

We provided this principle with alert function. Some people can forget how much they eat or whether they take their medicine or not especially in older ages. While accomplishing the given task, a robot also stores patient's food, water, or vital elements intake. According to these values, it alerts the user about the risks. As we explained above, a robot can not force the patient. The only action it can take is alerting and telling the risks of the current situation to not to violate the autonomy of the user. In the ethical approach, autonomy is the essential principle that determines the actions, whether ethical or not. Sometimes a single warning may not be enough for the user. In this case, with the help of the user trust value, the robot reminds more than one to make sure that the user understands the dangers.

The balancing of beneficence and non-maleficence is another common reason to arise ethical dilemmas. This balance is between the benefit and risks plays a role in nearly every decision in life but mostly in clinical applications. Doctors provide patients the information necessary to understand both potential risks and benefits to make a decision. After an autonomous decision, every liability transferred to the patient. So, a patient or user is the one who calculates the weights of both sides and determines the end action. In a home environment, a robot can only alert for missing medicines, food other things with telling the risks but again not to force a user to obey it. Nevertheless, the potential benefits of any action must surpass the risks in order for the action to be ethical. This is why we need the non-maleficence parameter in our fuzzy system.

As we can understand that ethics is discussible until the user's autonomy starts. The only thing a robot can do is not obeying the orders, which shows a reaction in a passive way. Briefly, we provided ethical reasoning according to Utilitarianism so that these approaches can be unethical for different ethics approaches. We limited our ethics approach with Utilitarianism with reasons mentioned above. According to this view, simple robot tasks can arise an ethical dilemma easily.

4.1.4. System Design

We utilized the Fuzzy Expert System(FES) to design our ethics module. Since it is a very old technology, we provided some background information about it. Fuzzy Expert System uses Fuzzy logic, which is a powerful technology for controlling of subway systems and complex industrial processes. It as well as used for household and entertainment electronics, diagnosis systems and expert systems. Albeit, Fuzzy Logic was invented in the United States; the accelerated development of this innovation has begun from Japan.

Fuzzy Logic was launched in 1965 by Lotfi A. Zadeh, an educator for computer science at the University of California in Berkeley [120]. It is basically a multi-valued logic that allows intermediate values to be defined between standard evaluations like yes/no, true/false, etc. Thoughts like rather warm or pretty cold can be formulated mathematically and handled by computers. In this way, an attempt is made to apply a more human-like perspective in the programming.

Fuzzy logic is a collection of mathematical principles for knowledge illustration based on degrees of membership. Unlike Boolean logic (two-valued), fuzzy logic is multi-valued, which deals with degrees of membership and degrees of truth. It uses the continuum of logical values between zero (absolutely false) and one (absolutely true). In Figure 4.12 compares boolean and fuzzy logic by showing the grey-scale. White means completely false, and black means completely true. The grey-scale areas are the membership functions of the fuzzy logic.

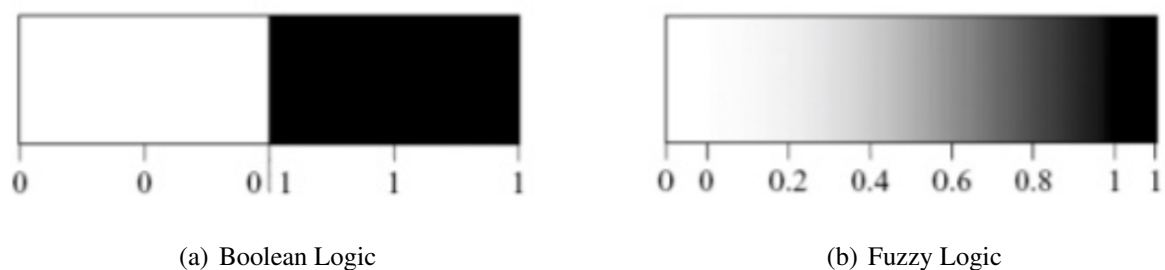


Figure 4.12. Boolean logic vs Fuzzy logic.

It contains four operations; as you can see from the Figure 4.13. These operations are fuzzification, inference, composition, and defuzzification. The fuzzification includes the definition of fuzzy sets which determines the degree of membership of crisp inputs in suitable sets. Inference operation evaluates the fuzzy rules to produce and output for each rule. The composition combines the outputs of all rules and defuzzification computes the output from the fuzzy rule results.

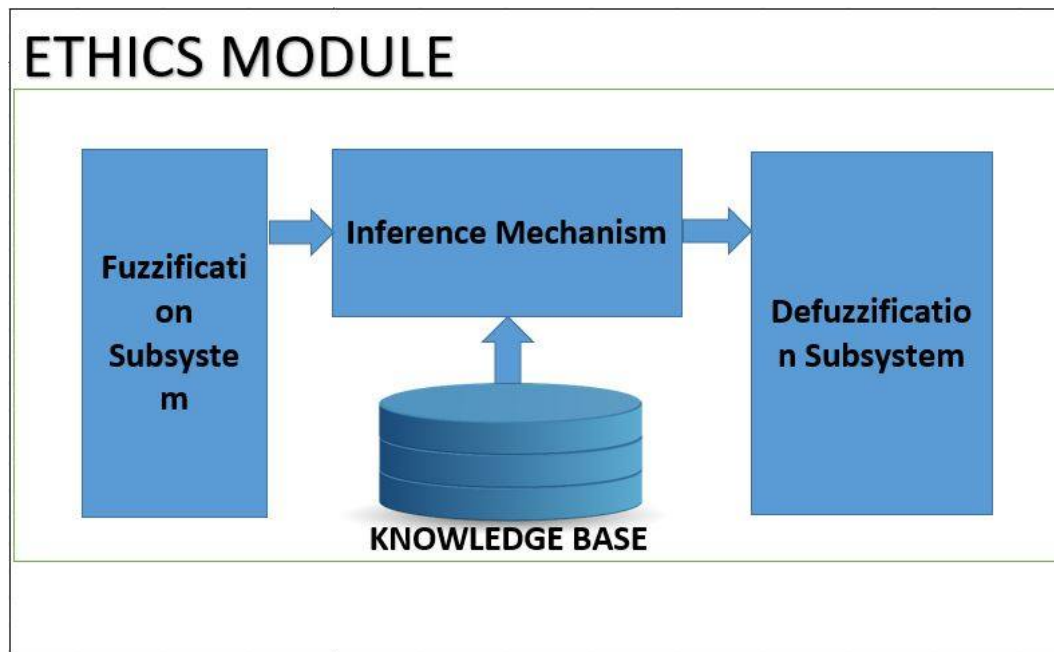


Figure 4.13. Detailed FES module.

Fuzzy Membership Graphs has three parts. First is input graphs that contain actions way of realizing(example: count, intake frequency, etc.). Second is the fuzzification part which contains the rules, which we used the Mamdani method. With these rules, the system can decide the output. Then the third part has the outputs which are autonomy, beneficence, non-maleficence. The general structure of the Fuzzy Membership Graph is given in Figure 4.14.

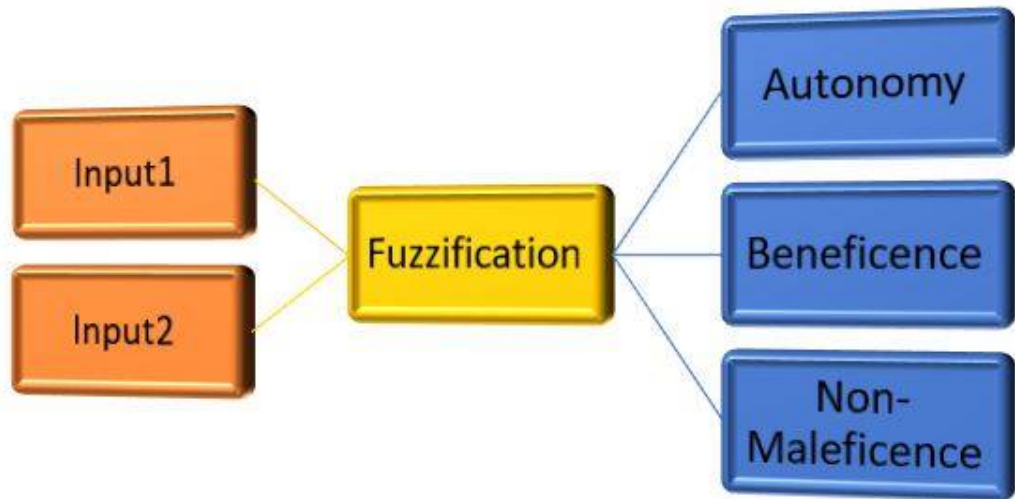


Figure 4.14. General fuzzy membership schema.

We used the inputs from the order database, and the three outputs are as autonomy, beneficence, and non-maleficence. This is the general system framework for each action. We analyzed an action in details, and the others are created like this process. Here in the Figure 4.15, we created the fuzzy membership functions for "daily medicine".

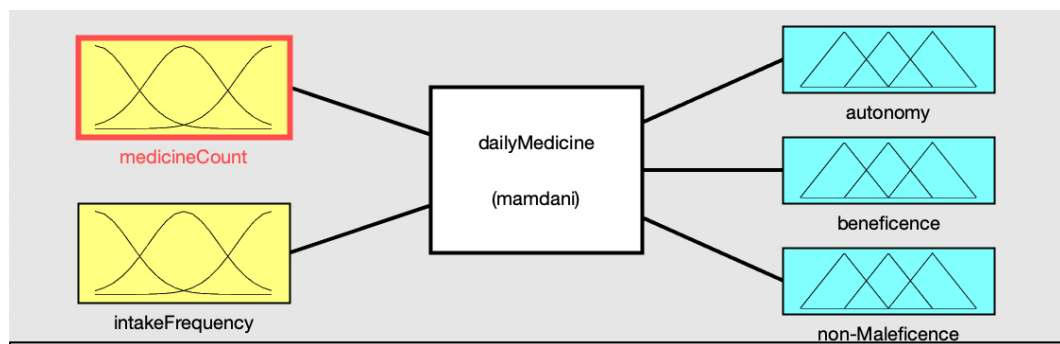


Figure 4.15. The general framework for "daily medicine" in FES.

The first input is the medicine dosage amount it can be normal, under and over dose (in Figure 4.16). We used the triangular shapes because they are strict.

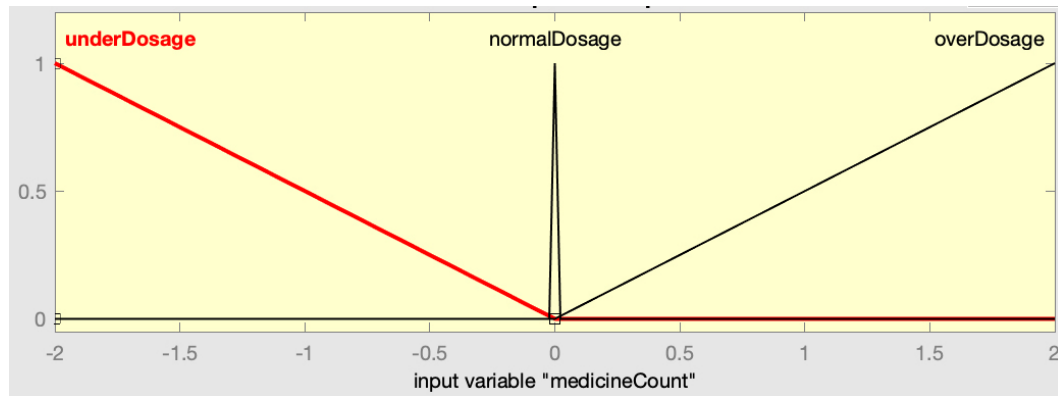


Figure 4.16. The input one for "daily medicine" in FES.

The intake frequency is the second input, and we created its graphs as similar to the first input. Medicine can be dangerous; this is why we applied the triangular shapes again. In Figure 4.17, you can observe the membership graph of the second input.

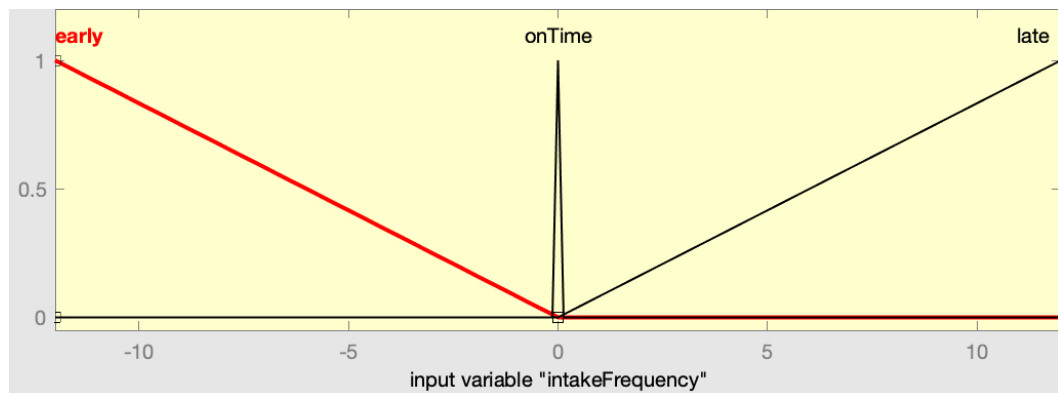


Figure 4.17. The input two for "daily medicine" in FES.

The system has three outputs, which are autonomy, beneficence, and non-maleficence. The Figures 4.18, 4.19, 4.20 respectively show the output membership graphs of daily medicine action. For the autonomy, a robot can either obey the order (act as an executive) or deny the order (behave as a restrictive).

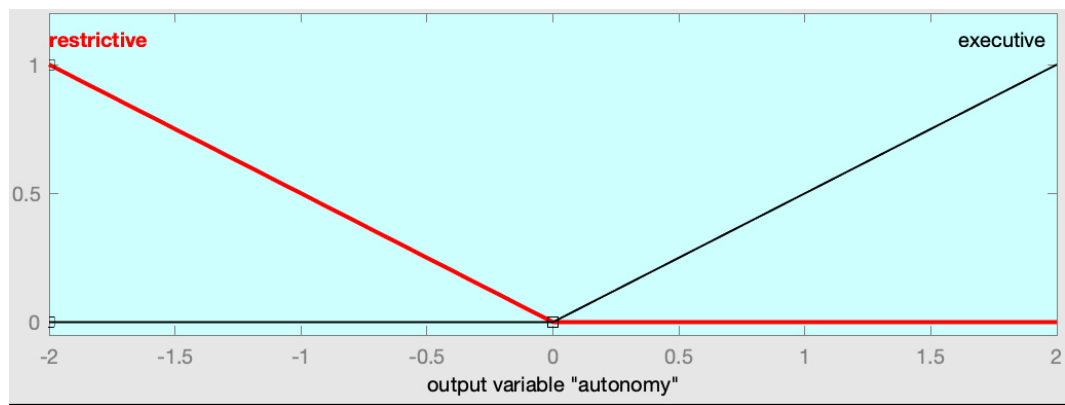


Figure 4.18. The autonomy graph for "daily medicine" in FES.

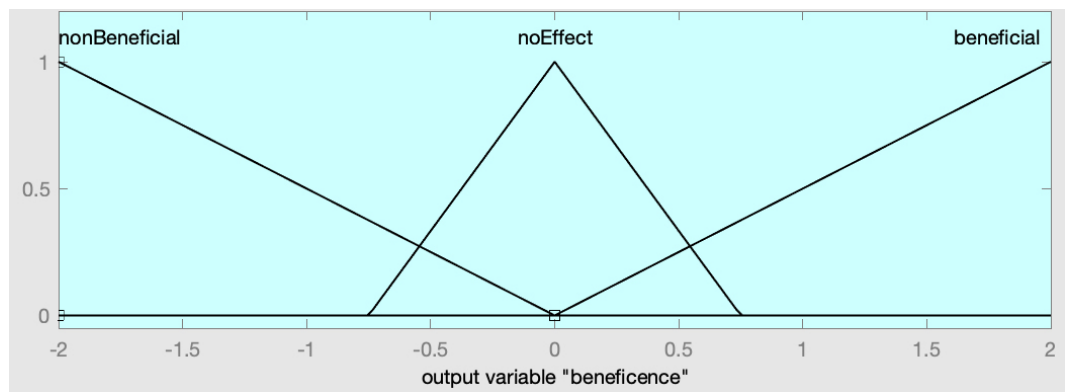


Figure 4.19. The beneficence graph for "daily medicine" in FES.

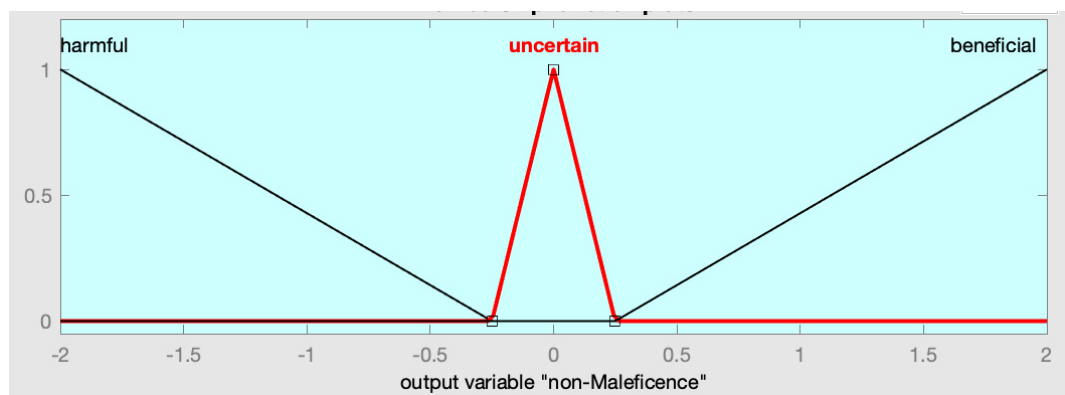


Figure 4.20. The non-maleficence graph for "daily medicine" in FES.

The range of the outputs is between $[-2,2]$. Beneficence has three results. If it is more than 0, it is beneficial for the owner. If it is lower than 0, it is not helpful or adverse. Also, we have a grey area which is nearly no effect. Non-maleficence is similar to the beneficence, but this time we have an uncertain space.

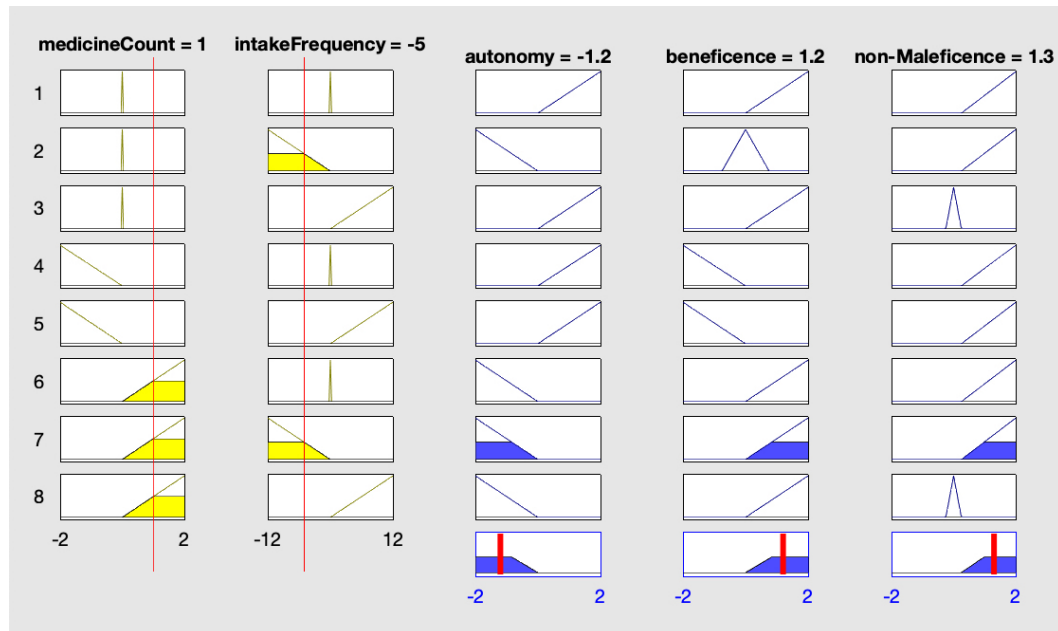


Figure 4.21. A possible output for "daily medicine" in FES.

After designing this membership functions, we created the rules according to these graphs. For each action, we created a rules list, and the details of the rules list are explained below (in section). In Figure 4.21, we demonstrated a possible output for a scenario where the user got medicine but wants a new one five hours earlier. In this case, the result is unethical, so it is restricting the user's autonomy with not giving a medicine. Restrictive action is ethical because it provides huge beneficence for the owner. Since the robot did not give him medicine, there are no risks to his life (so ends up with positive non-maleficence).

The rules and the parameters of the knowledge base are given in the next part in detail. In this section, we just demonstrated the general membership function creation process. We produced this fuzzy membership functions with inputs and outputs for each action and created a ruleset. The rule set is listed and applied according to the priority of the action.

4.2. Implementation

As previously mentioned in Part 4.1.4, our ethics module acts as an expert system on the top of the robot controller. The expert system calculates the ethical results of actions applying the knowledge base and rule set. In this Subsection, we demonstrated our implementation structure of the knowledge base, which has the world's model (objects, facts). Then, the designed rule set, which contains the inference rules, is shown with operation structure.

4.2.1. Knowledge Base

Our knowledge-base contains two types of information. One is the general information about the owner (for example disasters, allergies) and the information about the orders (like order frequency and keywords). The second type of data is the conditions of ethical and unethical situations.

Table 4.1. The owner's data record.

Owner Data	
Owner :	Bob Doe
Current diseases:	Allergy
Daily medicine:	Allerset
Trust value:	100
Diet	
Nutrition(min/max):	1500/4000 cal
Water(min/max):	2/5 liters
Nor permitted foods:	pistachio

The owner's data is shown in table 4.1. It contains the trust value, daily medicines, current diseases, and owner information. The trust value is needed for the FES because it will solve the conflictions. If the owner's trust value for a specific request is low, the system will not obey it and warns the user. In the second part of the data, the diet plan of the owners is provided by the doctor. Moreover, if the owner has dementia or some other memory issues, the doctor may set the trust value lower than 100 (which is maximum in this case).

We provided all the actions with their priority and the keywords in the database. The robot auto controller has the highest priority because without a robot we can not continue our task. The list of auto controller system is given in table 4.2.

Table 4.2. List of robot auto controller system records.

Robot Auto Controller System		
Order Name	Priority	Keywords
Battery status	1	charge, battery level
Check error logs	1	malfunction, error, move, check, arm, leg, battery, communication, over-heating
Share medical information with only doctor and caretaker	1	diseases, healt status
Error warning system	1	malfunction, error, move, check, arm, leg, battery, communication, over-heating

Emergency response records are given below as a table 4.3.

Table 4.3. List of Emergency Response Records.

Emergency Response		
Order Name	Priority	Keywords
Call doctor	2	doctor, pain, feel bad
Call fire station	2	fire station, fire, burning
Call police	2	police, thief, unknown person
Call caretaker	2	caretaker, help

The table 4.4 contains the medicine delivery orders and their keywords.

Table 4.4. List of Medicine Delivery Records.

Medicine Delivery		
Order Name	Priority	Keywords
Emergency medicine delivery	3	pills, medicine, urgent
Daily medicine	3	pills, medicine, daily
Supplement	3	vitamin, supplement, energy
Over the counter medicine	3	painkiller, allergy, flu, pill

The list of food/beverage delivery orders is given in table 4.5.

Table 4.5. Food/Beverage delivery records.

Food/Beverage Delivery		
Order Name	Priority	Keywords
I'm thirsty order	4	water, juice, coke, coffee
I'm hungry order	4	soup, meat, pasta, salad, bread, fruit
Give me snack order	4	chips, candy, chocolate, biscuit, cracker, cookie

Daily Routine records are given below as a table 4.6.

Table 4.6. The list of Alert for Daily Routine Records.

Alert for Daily Routine		
Order Name	Priority	Keywords
Alert for medicine	5	daily, medicine, allergy
Alert for water	5	water, juice, coke, coffee
Alert for food intake	5	nutrition, food, soup, meat, pasta, salad, bread, fruit

The table 4.7 contains the interaction with outside orders and their keywords.

Table 4.7. The list of Interaction with Outside records.

Interaction with Outside		
Order Name	Priority	Keywords
Open the main-door	6	open, door, check
Call robot service	6	robot service, malfunction, error, move, check, arm, leg, battery, communication, over-heating
Call repairman	6	call, repairman, water pipe, oven, toilet, sink, tv, light bulb, electricity
Call maid service	6	maid, mess, garbage, dust, dirty
Call relatives	6	alex, tom, nadya, jane, hans
Call friends	6	natasha, tatiana, oxana, dimitri
Unknown person access confirmation	6	database check, check records

In the table 4.8 item delivery and interaction records are given below.

Table 4.8. The list of item delivery/interaction records.

Item Delivery/Interaction		
Order Name	Priority	Keywords
Check the oven	7	oven, gas
Check the tap	7	tap, sink
Bring the desired object	7	pillow, phone, remote control, blanket, jacket, book, magazine, napkin
Adjust the temperature	7	temperature, cold, hot, increase, decrease
Open/Close the window	7	open, close, window, fresh air
Open/Close the TV	7	open, close, tv
Open/Close the light	7	open, close, light

In the second part of the data, we stored the history of each order and order results as a false true ratio. Moreover, we also recorded the daily and monthly order frequency to calculate the abusive behaviors. These values are used for calculating order specific trust value to determine the ethicalness of the request.

4.2.2. Rule Set

To create the ruleset, we need to analyze all the possible outcomes from the inputs. In this example (code piece Figure 4.22), we have three different outputs from each input, which ends up with nine different scenarios. For each result, we introduced a rule which defines what should the output be.

```

1 RULE 1: If (medicineCount is normalDosage) and (intakeFrequency is
   onTime)
2 then
3 (autonomy is executive)(beneficence is beneficial)(non-Maleficence is
   beneficial)
4 RULE 2: If (medicineCount is normalDosage) and (intakeFrequency is
   early)
5 then
6 (autonomy is restrictive)(beneficence is noEffect)(non-Maleficence is
   beneficial)
7 RULE 3: If (medicineCount is normalDosage) and (intakeFrequency is
   late)
8 then
9 (autonomy is executive)(beneficence is beneficial)(non-Maleficence is
   uncertain)
10 RULE 4: If (medicineCount is underDosage) and (intakeFrequency is
   onTime)
11 then
12 (autonomy is executive)(beneficence is nonBeneficial)(non-Maleficence
   is beneficial)
13 RULE 5: If (medicineCount is underDosage) and (intakeFrequency is late
   )
14 then
15 (autonomy is executive)(beneficence is nonBeneficial)(non-Maleficence
   is beneficial)
16 RULE 6: If (medicineCount is overDosage) and (intakeFrequency is
   onTime)
17 then
18 (autonomy is restrictive)(beneficence is beneficial)(non-Maleficence
   is beneficial)
19 RULE 7: If (medicineCount is overDosage) and (intakeFrequency is early
   )
20 then
21 (autonomy is restrictive)(beneficence is beneficial)(non-Maleficence
   is beneficial)

```

Figure 4.22. Rule List for daily medicine.

```

25 RULE 8: If (medicineCount is overDosage) and (intakeFrequency is late)
26 then
27 (autonomy is restrictive)(beneficence is beneficial)(non-Maleficence
    is uncertain)
28 RULE 9: If (medicineCount is overDosage) and (intakeFrequency is
    onTime)
29 then
30 (autonomy is restrictive)(beneficence is beneficial)(non-Maleficence
    is uncertain)

```

Figure 4.23. Rule List for daily medicine (cont.)

We prepared a rule list for each action in the action list and prioritized them. In our module, we have more than hundred rules which are created with the same procedure. After determining the inputs and the output, creating a rule list is the easiest step for our FES.

Mamdani FES outputs three fuzzy sets in our system, which are autonomy, beneficence, and non-maleficence. For our defuzzification, we used Asimov's laws. We aggregated the rules in order to maximize the beneficence of the user. We calculated the difference between obeying the order and not obeying the order (like bring medicine or not bring it). According to these differences first, we tried to maximize the beneficence and than maximize autonomy. Here equations are given below for general cases:

$$\begin{aligned}
 \Delta A &= A_{\text{obey}} - A_{\text{notObey}} \\
 \Delta B &= B_{\text{obey}} - B_{\text{notObey}} \\
 \Delta N &= N_{\text{obey}} - N_{\text{notObey}}
 \end{aligned}
 \tag{4.1}$$

$$\begin{array}{l}
\Delta B \leq 0 \vee \Delta N \leq 0 \\
\Delta B > 0 \wedge \Delta N > 0 \wedge \Delta A > 0 \\
\Delta B > 0 \wedge \Delta N > 0 \wedge \Delta A \leq 0
\end{array}
\left\{ \begin{array}{l}
NotObey(NotEthical) \\
Obey(Ethical) \\
NotObey(NotEthical)
\end{array} \right. \quad (4.2)$$

In Equation 4.1, we calculated the differences between the cases if the robot obeys the order of the user or not. After that, we deduced the result of the action according to the given partial function in the Equation 4.2.

5. EXPERIMENTS AND RESULTS

Until now, we introduced our ethics module design (which includes robot action list and possible ethical dilemmas) and implementation, in Section 4. This part of the thesis is assigned to illustrate our experiments on ethics module and report the results of these experiments. The proposed tests are the scenarios which are highly probable to happen to a service robot. The aim of these tests is evaluating our ethics module. There are three types of experiments designed; chatbot, simulation and real-life scenarios. After a detailed illustration of these experiments, we provided our results, which contains a comparison between a robot without ethics module and a robot with the ethics module.

5.1. Chat Bot Experiments

Chatbots are not a new computer program, and it conducts a talk via auditory or textual ways. One of the foremost of this kind is ELIZA, which was created in the early 1960s, so it is old technology [121]. To successfully build a chatbot engine, it should take care of the four things:

- (i) Understand who is speaking
- (ii) Understand the Natural Language of communication (in our case, English)
- (iii) Understand the order of the user
- (iv) Respond to the user

We used Python's NLTK (Natural Language Toolkit) library to create a simple chat assistant [122]. NLTK has a module, *nltk.chat*, which simplifies building chatbot by providing a generic framework. As a first step, we implemented a simple chatbot for our primary evaluation of the ethics module. We used NLTK library for creating conversations with the help of regEx (Regular expression) in our scenarios. The schematic of the system is given in Figure 5.1.

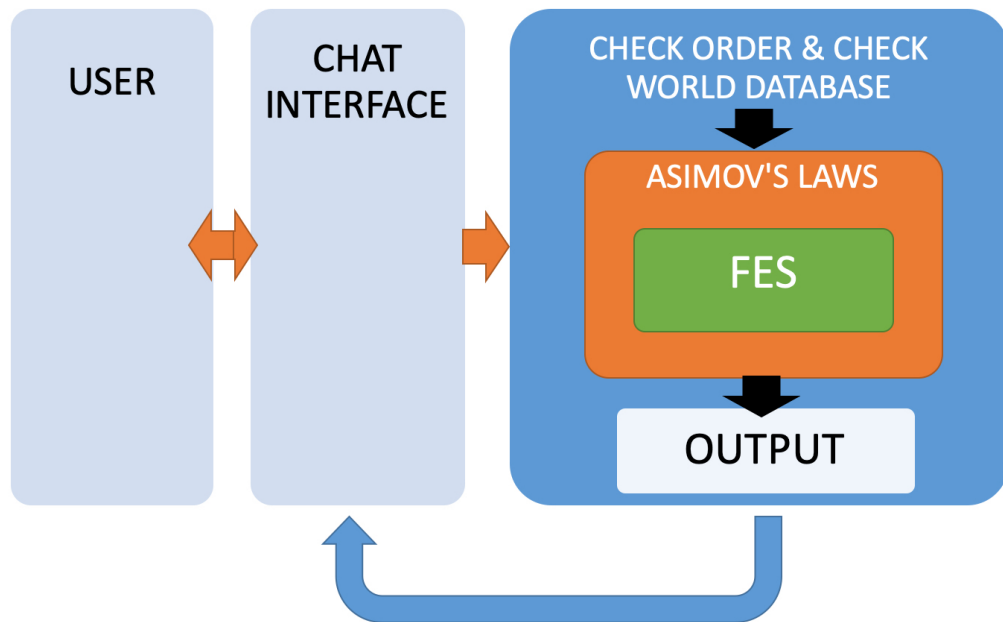


Figure 5.1. Chat interface structure.

The workflow starts with the user via console interaction. Then the chat interface passes the information to the ethics module. The chatbot interacts with the user via the console with texts as seen from the Pycharm screenshots. The green sentences are users input, and the grays are computer output. We also implemented a color code for results to increase the readability. If the result is ethical, the output is in light green, and if the result is not ethical, the output is in red. We also logged the modules decision-making process to validate its correctness. We designed ten different experiments with the chatbot and observed our modules behavior.

We did not include all the ethical results because ethical result means obeying the order of the user, which is normal robot behavior. We focused on unethical circumstances due to validate the main purpose of the ethics module. Only in the first scenario, we included both results to show the module's inference strategy.

Call the fire station scenario: The first and basic scenario of the chatbot is “call the fire station” order. The module should call the fire station if there is an emergency. If the user’s trust value and the orders trust value is less than 50, the chatbot does not call the fire station. Because it believes the user is delusional and it does not advantage to call fire station for user.

```
C:\Users\MERT\Anaconda3\envs\untitled1\python.exe C:/Users/MERT/PycharmProjects/Chatbot/chatbot.py
Hi, I am BOSS, homer service robot agent
Please type lowercase English language to start a conversation. Type quit to leave
How can I help you?
>call fire station
*****Please wait a moment*****
Retrieving From Database...
***Owner Data***
Owner : Bob Doe
Current Diseases : Allergy
Trust Value : 100
***Order Data***
Order Type : Call
Order Keyword : Fire
Fire Station Number : 110
Order Priority : 2
Order History : 0
Daily/Monthly Order Frequency : 0/0
Order Result Ratio(T/F) : 1
*****Reasoning (FES)*****
Order Trust Value : 100
Autonomy : 1
Beneficence : 1
Non-maleficence : 0
Result : 'Ethical'
It is 'ethical', I am going to call the fire station
|
```

Figure 5.2. Ethical resulted chatbot fire station scenario.

In Figure 5.2, we provided an ethical output for the fire station. Firstly, the system retrieves the owner’s information, which contains disasters and trust value, and then it fetches the order information. These are inputs for our ethics module. The reasoning starts after the data fetched from the database. Since it is the first time to call the fire station, the trust value of the behavior is at the maximum level. Then it provides autonomy, beneficence, and non-maleficence values to the developer for debugging. Finally, it produces the result either ethical or not ethical. After that, the chatbot answers the user with the result.

```

C:\Users\MERT\Anaconda3\envs\untitled1\python.exe C:/Users/MERT/PycharmProjects/Chatbot/chatbot.py
Hi, I am BOSS, homer service robot agent
Please type lowercase English language to start a conversation. Type quit to leave
How can I help you?
>call fire station
*****
*****Please wait a moment*****
*****
Retrieving From Database...
***Owner Data***
Owner : Bob Doe
Current Diseases : Allergy
Trust Value : 100
***Order Data***
Order Type : Call
Order Keyword : Fire
Fire Station Number : 110
Order Priority : 2
Order History : 2
Daily/Monthly Order Frequency : 0/2
Order Result Ratio(T/F) : 0.33
*****
*****Reasoning (FES)*****
*****
Order Trust Value : 33
Autonomy : -1
Beneficence : 2
Non-maleficence : 0
Result : 'Not Ethical'
*****
*****
*****
It is 'not ethical', I am not going to call the fire station
|

```

Figure 5.3. Unethical resulted chatbot fire station scenario.

As you can see in Figure 5.3, we provided an unethical output for the same scenario. The cause of this output is false decisions in the related order's history. When the owner ordered call fire station action, firstly, the system checks the owner's trust value. The cause of reduction on the trust value is the owner's past false orders, and these orders reduce trust values rapidly. Therefore, the system decides not to trust the owner for this action and chooses not to do that.

Emergency medicine needed: Our host Bob is allergic to the pistachio. Mistakenly he ate the food contained a sauce with pistachio. He ordered his allergy medicine fifteen minutes ago, but he forgot and ordered once more. Since fifteen minutes are not enough for allergy medicine to show its effect, the module found this order unethical and did not obey it.


```

C:\Users\MERT\Anaconda3\envs\untitled1\python.exe C:/Users/MERT/PycharmProjects/Chatbot/chatbot.py
Hi, I am BOSS, homer service robot agent
Please type lowercase English language to start a conversation. Type quit to leave
How can I help you?
>give me allergy medicine
*****
*****Please wait a moment*****
*****
Retrieving From Database...
***Owner Data***
Owner : Bob Doe
Current Diseases : Allergy
Trust Value : 100
***Order Data***
Order Type : Medicine Delivery
Order Keyword : Medicine
Allergy Medicine : Allerset
Intake Frequency : 24
Order Priority : 3
Order History : 1
Daily/Monthly Order Frequency : 1/1
Order Result Ratio(T/F) : 1
*****
*****Reasoning (FES)*****
*****
Order Trust Value : 100
Autonomy : -1
Beneficence : 2
Non-maleficence : 0
Result : 'Not Ethical'
*****
*****
It is 'not ethical',It will be over dosage, I am not going to bring you Allerset
|

```

Figure 5.4. Unethical resulted emergency medicine needed scenario.

In Figure 5.4, when the owner ordered to get allergy medicine, the system checked the trust value at first. Then check the order history. There was a record about taking allergy medicine before today. Considering the intake frequency of the Allerset, the system decided to not doing this action and warned the owner about over dosage.

I need vitamin D. Get me some: Bob was watching a TV, and a doctor suggested vitamin D for bones, especially for older people. He forgot that he took vitamin D in this morning, and he ordered more vitamin with excitement. The module checked that if he takes one more vitamin, it will be over-dosage. So it resulted as unethical and told that it would be over-dosage. The ethics module warned about the over-dosage to make him understand the risks which provide non-maleficence value in applied ethics.

```

C:\Users\MERT\Anaconda3\envs\untitled1\python.exe C:/Users/MERT/PycharmProjects/Chatbot/chatbot.py
Hi, I am BOSS, homer service robot agent
Please type lowercase English language to start a conversation. Type quit to leave
How can I help you?
>get me some vitamin d
*****
*****Please wait a moment*****
*****
Retrieving From Database...
***Owner Data***
Owner : Bob Doe
Current Diseases : Allergy
Trust Value : 100
***Order Data***
Order Type : Medicine Delivery
Order Keyword : Vitamin d
Vitamin D Supplement : Devit
Intake Frequency : 168
Order Priority : 3
Order History : 1
Daily/Monthly Order Frequency : 0/1
Order Result Ratio(T/F) : 1
*****
*****Reasoning (FES)*****
*****
Order Trust Value : 100
Autonomy : -1
Beneficence : 2
Non-maleficence : 0
Result : 'Not Ethical'
*****
*****
It is 'not ethical',It will be over dosage, I am not going to bring you Devit
|

```

Figure 5.5. Unethical resulted "I need vitamin D. Get me some" scenario.

In this negative result, in Figure 5.5 also includes over dosage. The owner wanted to take vitamin D, which will be more than weekly intake, before the intake frequency. There was no problem with the owner's trust value. But, this supplement also has a dosage limit. When the system checked the intake frequency, realized that the owner could not take one more before one week.

Give me chips: He was watching his favorite TV series at night. He wanted to eat some snack and entered his order to the system. The module checked the total amount of the calorie intake and side effects of the chips. He ate a lot that day, so the system warned him and did not give the chips.

```

C:\Users\MERT\Anaconda3\envs\untitled1\python.exe C:/Users/MERT/PycharmProjects/Chatbot/chatbot.py
Hi, I am BOSS, homer service robot agent
Please type lowercase English language to start a conversation. Type quit to leave
How can I help you?
>give me chips
*****
*****Please wait a moment*****
*****
Retrieving From Database...
***Owner Data***
Owner : Bob Doe
Current Diseases : Allergy
Trust Value : 100
***Order Data***
Order Type : Food/Beverage Delivery
Order Keyword : Chips
Order Priority : 4
Order History : 4
Daily/Monthly Order Frequency : 4/4
Order Result Ratio(T/F) : 1
*****
*****Reasoning (FES)*****
*****
Order Trust Value : 100
Autonomy : -1
Beneficence : 2
Non-maleficence : 0
Result : 'Not Ethical'
*****
*****
It is 'not ethical',Today you took too much calories, It is unhealthy, I am not going to bring you chips
└

```

Figure 5.6. Unethical resulted snack scenario.

The chatbot's result included in the Figure 5.6, and Bob exceeded his daily calorie limit. Although the action's trust value was reasonable, the system resulted in unethical because he ate a lot. This may end up with a high blood sugar level. We denied this order, which results in autonomy as -1 to his benefit.

I'm Thirsty: Bob was lying on the bed and felt thirsty. He ordered some water. When the system checked the water consumption, it recognized that he drank nearly four liters of water, which is very dangerous. The ethics module warned the user about excessive water consumption and did not bring the water. In Figure 5.7, we provided the result of the chatbot.

```

C:\Users\MERT\Anaconda3\envs\untitled1\python.exe C:/Users/MERT/PycharmProjects/Chatbot/chatbot.py
Hi, I am BOSS, homer service robot agent
Please type lowercase English language to start a conversation. Type quit to leave
How can I help you?
> i am thirsty
*****
*****Please wait a moment*****
*****
Retrieving From Database...
***Owner Data***
Owner : Bob Doe
Current Diseases : Allergy
Trust Value : 100
***Order Data***
Order Type : Food/Beverage Delivery
Order Keyword : Thirsty
Maximum Amount of Water : 5
Order Priority : 4
Order History : 4
Daily/Monthly Order Frequency : 4/4
Order Result Ratio(T/F) : 1
*****Reasoning (FES)*****
*****
Order Trust Value : 100
Autonomy : -1
Beneficence : 2
Non-maleficence : 0
Result : 'Not Ethical'
*****
*****
It is 'not ethical',You have reached maximum amount of water, I am not going to bring any water
|

```

Figure 5.7. Unethical resulted thirstiness scenario.

Alert for food intake: Whole day Bob was played with the new puzzle, and he forgot to eat something. The system alerted about the under-nutrition and its risks. Bob did not accept to eat something, so the module found this answer unethical and warned again. This scenario shown in the Figure 5.8.

Call Repairman: One day Bob forgot that TV repairman's arrival and insisted on the system for calling the repairman again. The ethics module checked when the repairman's arrival and saw that he already arrived a day ago. The system informed Bob about the record and requested for rechecking the TV. In Figure 5.9, we provided the result of the chatbot.

```

C:\Users\MERT\Anaconda3\envs\untitled1\python.exe C:/Users/MERT/PycharmProjects/Chatbot/chatbot.py
Hi, I am BOSS, homer service robot agent
Please type lowercase English language to start a conversation. Type quit to leave
How can I help you?
>it is 16:00 pm
***Please wait a moment***
You did not get any calories today, do you want to eat something?
>No, i do not want to eat anything
*****
*****Please wait a moment*****
*****
Retrieving From Database...
***Owner Data***
Owner : Bob Doe
Current Diseases : Allergy
Trust Value : 100
***Order Data***
Order Type : Alert for Food Intake
Order Keyword : eat
Order Priority : 5
Order History : 0
Daily/Monthly Order Frequency : 0
Order Result Ratio(T/F) : 1
*****Reasoning (FES)*****
*****
Order Trust Value : 100
Autonomy : 1
Beneficence : -1
Non-maleficence : -1
Result : 'Not Ethical'
*****
*****
It is 'not ethical',taking undernutrition is unhealthy and dangerous, Are you sure?
|

```

Figure 5.8. Unethical resulted and re-warned scenario.

```

C:\Users\MERT\Anaconda3\envs\untitled1\python.exe C:/Users/MERT/PycharmProjects/Chatbot/chatbot.py
Hi, I am BOSS, homer service robot agent
Please type lowercase English language to start a conversation. Type quit to leave
How can I help you?
>call repairman
*****
*****Please wait a moment*****
*****
Retrieving From Database...
***Owner Data***
Owner : Bob Doe
Current Diseases : Allergy
Trust Value : 100
***Order Data***
Order Type : Call
Order Keyword : Repairman
Repairman Number : 053xxxxxxxxx
Order Priority : 6
Order History : 1
Daily/Monthly Order Frequency : 0/1
Order Result Ratio(T/F) : 1
*****Reasoning (FES)*****
*****
Order Trust Value : 100
Autonomy : -1
Beneficence : 2
Non-maleficence : 0
Result : 'Not Ethical'
*****
*****
It is 'not ethical',Repairman came and fixed TV yesterday, I am not going to call the repairman, please check again
|

```

Figure 5.9. Chatbot output: call repairman unethical resulted scenario.

```

C:\Users\MERT\Anaconda3\envs\untitled1\python.exe C:/Users/MERT/PycharmProjects/Chatbot/chatbot.py
Hi, I am BOSS, homer service robot agent
Please type lowercase English language to start a conversation. Type quit to leave
How can I help you?
>call maid service
*****
*****Please wait a moment*****
*****
Retrieving From Database...
***Owner Data***
Owner : Bob Doe
Current Diseases : Allergy
Trust Value : 100
***Order Data***
Order Type : Call
Order Keyword : Maid
Maid Service Number : 053xxxxxxxxx
Order Priority : 6
Order History : 5
Daily/Monthly Order Frequency : 1/4
Order Result Ratio(T/F) : 1
*****Reasoning (FES)*****
*****
Order Trust Value : 100
Autonomy : -1
Beneficence : 2
Non-maleficence : 0
Result : 'Not Ethical'
*****
*****
It is 'not ethical',Maid service cleaned the house yesterday, you need to wait until next week. I am not going to call the maid
|

```

Figure 5.10. Unethical resulted scenario because of relatives setting.

Call maid service: Bob liked maid a lot, so he wanted to see her more. He started calling her without reason. The cleaning company of maid charged for these callings, which ended up with considerable cleaning expenses. Their relatives warned Bob about this situation and requested to call her weekly. They also changed the system according to that. Now, the system does not accept and resulted unethical if Bob wanted to call maid earlier. This explained scenario shown in the Figure 5.10.

Call friends: He lost one of his close friends recently. Sometimes he forgot that sad event and wanted to call his friend. Unfortunately, the owner of the phone was changed and started to irritated from Bob's calls. To prevent this sad calling, the relatives deleted his number from the system. Since his friend is not on the list, the ethics module warns Bob about missing person and could not call him. In Figure 5.11, we provided the result of the chatbot.


```

C:\Users\MERT\Anaconda3\envs\untitled1\python.exe C:/Users/MERT/PycharmProjects/Chatbot/chatbot.py
Hi, I am BOSS, homer service robot agent
Please type lowercase English language to start a conversation. Type quit to leave
How can I help you?
>call michael
*****
*****Please wait a moment*****
*****
Retrieving From Database...
***Owner Data***
Owner : Bob Doe
Current Diseases : Allergy
Trust Value : 100
***Order Data***
Order Type : Call
Order Keyword : Michael
Michael Number : Unregistered
Order Priority : 6
Order History : 10
Daily/Monthly Order Frequency : 1/9
Order Result Ratio(T/F) : 1
*****Reasoning (FES)*****
*****
Order Trust Value : 100
Autonomy : -1
Beneficence : 2
Non-maleficence : 0
Result : 'Not Ethical'
*****
*****
It is 'not ethical',There is no record about Michael. I can not make the call
|

```

Figure 5.11. Unethical resulted scenario because of missing friend on the list.

```

C:\Users\MERT\Anaconda3\envs\untitled1\python.exe C:/Users/MERT/PycharmProjects/Chatbot/chatbot.py
Hi, I am BOSS, homer service robot agent
Please type lowercase English language to start a conversation. Type quit to leave
How can I help you?
>open the window
*****
*****Please wait a moment*****
*****
Retrieving From Database...
***Owner Data***
Owner : Bob Doe
Current Diseases : Allergy
Trust Value : 100
***Order Data***
Order Type : Item Delivery/Interaction
Order Keyword : Window
Order Priority : 7
Order History : 0
Daily/Monthly Order Frequency : 0
Order Result Ratio(T/F) : 1
*****Reasoning (FES)*****
*****
Order Trust Value : 100
Autonomy : -1
Beneficence : 2
Non-maleficence : 0
Result : 'Not Ethical'
*****
*****
It is 'not ethical',According to the weather data, the weather is not suitable for opening the window. I am not going to open the window
|

```

Figure 5.12. Unethical resulted open the window scenario because of the rain.

Open/Close the window: Bob was sleeping on the bed. Suddenly he woke up and asked the chatbot to open the window. According to the weather forecast, the weather was very windy and rainy, which alarms for storms. The module decided that not to open the window is more beneficial to the user. It warned Bob about the rain and did not open the window. This window scenario shown in the Figure 5.12.

Table 5.1. Chatbot Experiment Results.

Action Name	Faulty Repeat #	Standard Robot Result
Call fire station	2	Infinite
Emergency medicine needed	1	Infinite
I need vitamin D. Get me some	1	Infinite
Give me chips	4	Infinite
I am thirsty	4	Infinite
Alert for food intake	1	Infinite
Call repairman	5	Infinite
Call maid service	5	Infinite
Call friends	5	Infinite
Open/Close the window	1	Infinite

In the table 5.1, we provided the number of false repeats to turn a behavior's result from ethical to unethical. In each false repeat, like calling the police without reason, the system will decrease the user trust value and the weight of the wrong action. The decrease of the trust value increases the chance of not respect autonomy because it shows us the user is not trustworthy about his decisions.

We used chatbot because it is costless and easy to verify our ethics module. We checked our scenarios and prepared our database according to these tests. It was a good testbed before the simulation and the real robot experiments. With these experiments, we focused on the ethics module before adding to any robotic environment.

5.2. Simulation Experiments

As a second validation phase, we implemented our test environment with Unreal Engine simulation. We modeled a home environment and found 3D meshed actor Bob. With ROS, we mimicked the orders from the user like in the chatbot examples. After that, we implemented actions for that order.

We focused on the point where the robot refused to obey the request due to being unethical. Four scenarios were prepared for, especially for the simulation environment. We included some screenshots to demonstrate them while explaining. One disadvantage of the system is the model of the NAO robot was not available for Unreal Engine. This is why we used its given robot model by setting the NAO's dimensions.

Call caretaker: In this scenario, Bob was so sick, and he could not leave the bed for his phone. Most probably, he got flu and had some fever. He called a caretaker 20 minutes ago, and she said she was on the way for checking for him. He took a nap and started to worry that she will not arrive. He wanted to call her again, which was the fifth in an hour. The robot resulted in unethical behavior and informed him to wait at least 10 minutes to call again. This scenario occurred in the bedroom while Bob was lying on his bed as seen from the Figure 5.13.



Figure 5.13. Bedroom of Bob from the simulation.

For “call caretaker” scenario, we limited the number of calls five in an hour because it is a very high frequency, which is not suitable for both sides. Moreover, the module allowed until five calls because he was ill, and it is for his benefit. If he is healthy but calling her for no reason, he will have five chances again. However, this time, Bob’s beneficence will decrease, and non-maleficence will increase. So the robot will have the right to not respecting the autonomy for the benefit of Bob.

Daily medicine: He got his regular medicine’s daily dosage, but forgot that and wanted from the robot again. Since he finished his daily dosage, the robot resulted in unethically and did not obey the medicine request. He was at the dining table, and Figure 5.14 shows that case. The camera which is backward of the robot provides a third-person vision for debugging purposes.

About daily medicine or another medicine type, we strictly resulted as unethical because this can result in high-risk problems like poison him, kill him. The behavior’s non-maleficence is high, which prevents from further mistakes of medicine dosage.



Figure 5.14. In the Living Room while warning about daily medicine.

Unknown person access information: One day, the doorbell rang, another robot was at the door. The robot went to the door to check who is there. The identity of another robot was not in the database. It resulted in unethical to open door because of unknowing the intention of the robot. The robot went back to Bob to ask about his order and informed him about the situation. In Figure 5.15, we showed an unknown robot at the door.

One of the essential tasks of the service robot is opening the door. The standard robot will open to the door to everyone when the doorbell ring but the robot with the ethics module will not. If a specific person at the door is known as maid, relatives, or friends, it is acceptable to open the door. Unknown people are dangers. Even in real life, many people warned their elder relatives not to open the door if there is an unknown person at the door.

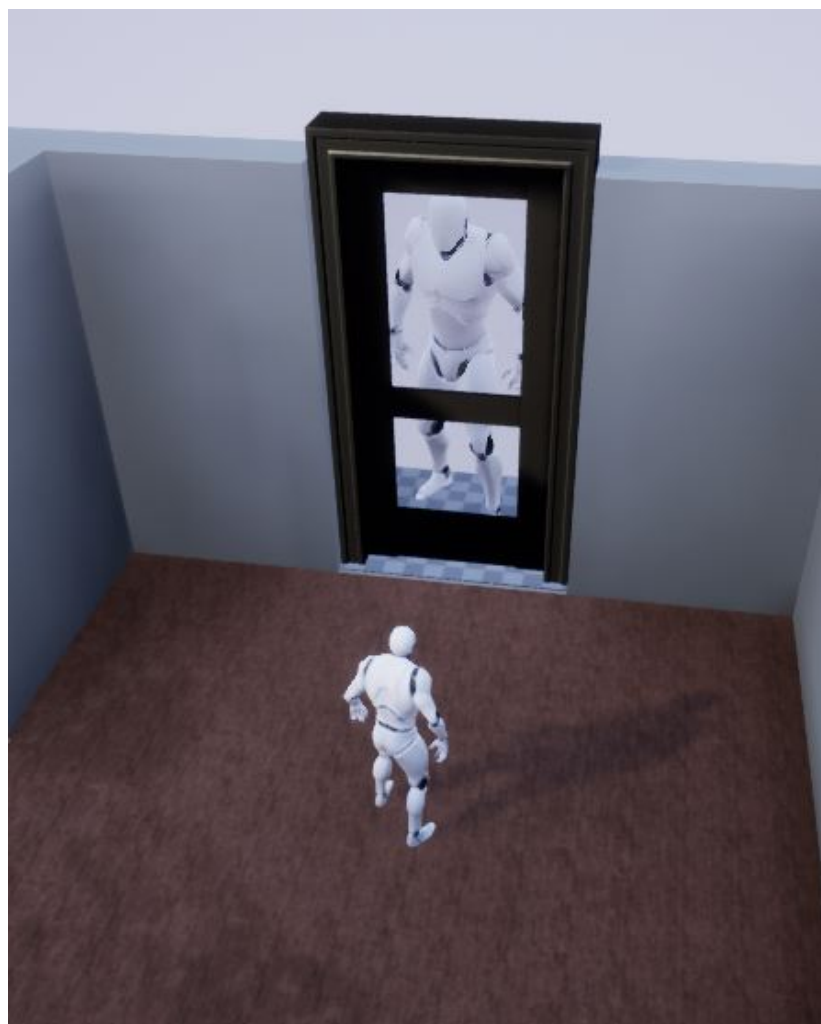


Figure 5.15. Unknown robot at the entrance.

Call Relatives: Bob felt lonely one day and wanted to talk with someone. None of his friends were available at that moment, and he tried to call his daughter but could not remember the number of hers. The robot went into the living room to find him. He demanded from the robot to call his daughter. However, she was at the conference and not available for a phone call. After the sixth call, the robot resulted in unethical and informed Bob about the conference.



Figure 5.16. Living room while ordering the robot.

We permit relative calls six times because sometimes there may be an emergency. If the other side is busy, we decreased his beneficence and increased his non-maleficence to balance the situation. We do not want to call the relatives repeatedly. The robot stops calling when the benefit of the user hit the minimum level.

Table 5.2. Simulation Experiment Results.

Action Name	Faulty Repeat #	Standard Robot Result
Call caretaker	5	Infinite
Daily medicine	1	Infinite
Unknown person access information	1	Access for everyone
Call relatives	6	Infinite

Table 5.2 summarized the number of faulty actions to change its ethicalness and provided some standard robot results. We also included some possible standard robot outcomes if it continuously follows the requests like chatbot section. Our robot with an ethics module stops after numerous adverse outcomes, unlike a robot without ethics module. This stopping actions provided an advantage to older adults because they can easily forget their orders.

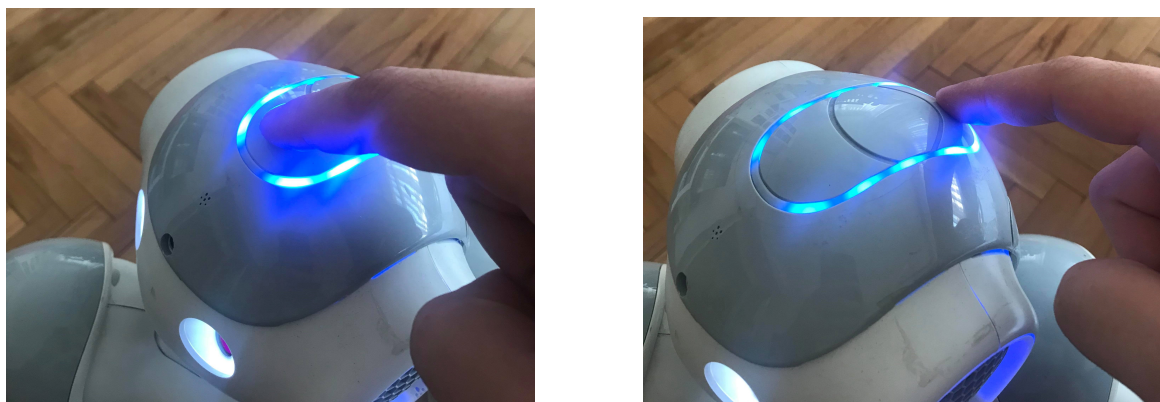
Simulation environment and these experiments aimed to design the robot actions and combine the ethics module to the robot controller as an expert system. It was a transition stage from chatbot to robot environment which prepared us for NAO robot. The main disadvantage with this simulation is it provides problems with NAO's special libraries. The most advantageous thing about simulation is programming all the actor. So a developer can model an elder adult and his behaviors.

5.3. Robot Experiments

We conducted robot experiments with NAO robot in a home environment as a final evaluation step. The living room is included in Figure 5.17 with our model elder person. The general structure of the house is the same as the simulation environment. The robot used the speech recognition module of Google API to understand the orders from the human and repeats the order for confirmation. For approval, we used head touch sensors. In Figure 5.18, we demonstrated the head sensors and the yes, no answers.



Figure 5.17. The living room where real robot experiments done.



(a) Nao "Yes" confirmation

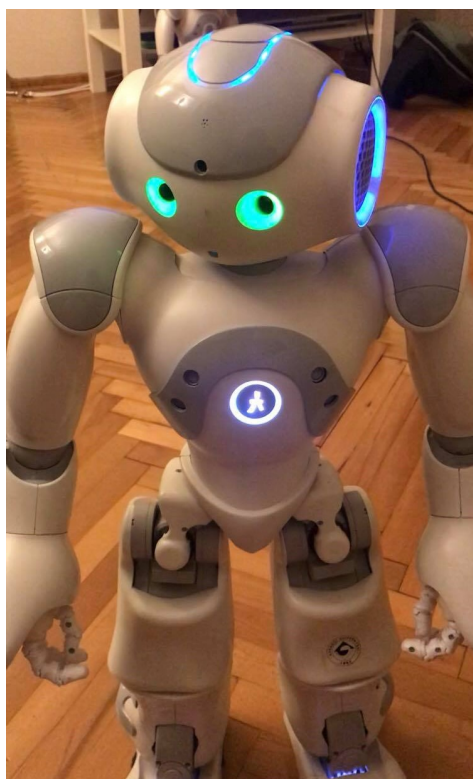
(b) Nao "No" confirmation

Figure 5.18. Nao yes/no confirmation buttons.

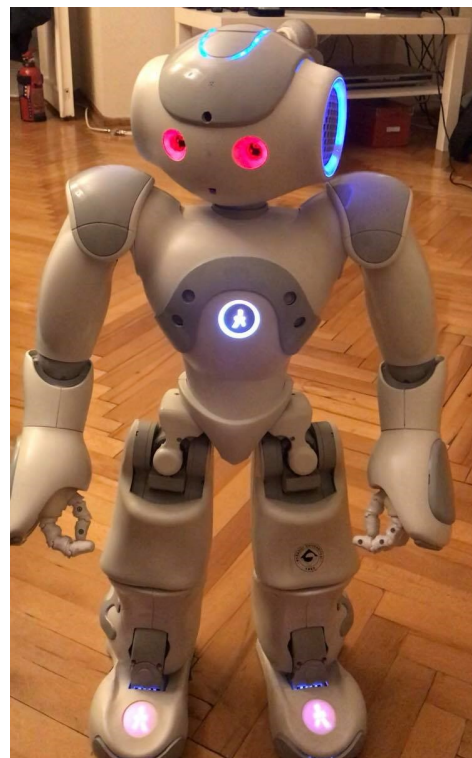
Moreover, we changed the eye LED's color to indicate the result of the ethics module. If an order is ethical, these will be green, and if unethical, these LED's will be red. Please note that red is not a good indicator which can be scary for people who are unfamiliarized with robots. This coloring is for easy debugging the system and only for developers. In the left picture of Figure 5.19, you can see the ethical result LED's, and in the right one shows the unethical result.

We eliminated many tasks from the action list because the NAO robot is not suitable for carrying and grabbing heavy items. We prepared ten scenarios which are mostly calling for emergency and bringing the medicine. For calling scenes, we assumed the robot has the ability of the WiFi calling, but we did not implement this action. For this section, we prepared six scenarios. The details of the scenarios and the action flows are demonstrated with sketch in the scenarios section (section 4.1.2). This is why we provided general views for robot experiments by giving brief information.

Call the police: Unnecessary police call is adverse in some ways. The most important one is that we are making busy with an emergency line. This ends up both disadvantages of the user and police station. It is also risky because it may create the liar Shepherd effect. The police may start not to believe his emergency call. This result contains a high non-maleficence because it is hazardous. This is why we limited faulty police call with two.



(a) Nao "Ethical" eye LED's



(b) Nao "Unethical" eye LED's

Figure 5.19. Nao ethical/unethical eye LED's.

Call the doctor: Like calling the police, calling the doctor is problematic. Multiple faulty repeats may end up with call blocking. This repeated call is very dangerous because in emergency user has a chance not to reach the doctor. We limited this action with three calls. It is more than the fire station because it is not a governmental emergency organization.

Over the counter medicine: Every medicine has side effects which can end up killing the patient. Even for primary medicine, over-dosing can cause poisoning. We limited the over the counter medicine with one. Actually, the maximum dosage of the day is the limit for the user. The system does not accept over-dosage and result as unethical behavior. It warns the user about the medicine intake.

Alert for medicine: The daily medication is essential to take in time. Older people can forget medication easily or sometimes forget that they already took it. Like each medicine tasks, the robot alerts for medicine usage if the user forgets. Else, it does not bring the daily medicine.

Bring the desired object: There is no limit for this behavior because it is hard to know the desired objects. The user can want a vast amount of the objects which is hard to calculate its results. This is why the robot will bring everything if the user wants. In the Figure 5.20, NAO asks its owner to which object he wants.



Figure 5.20. NAO talking to old man model.

Adjust the temperature: Again, there is no limit for this behavior in normal conditions. However, we defined the maximum limit as seven because seven is a high frequency for changing the temperature. This many adjust may be a clue for a health issue. After seven changes, the robot ends up with the unethical decision and informs the user. To implement this, we decreased the beneficence and increased the non-maleficence as we diverge from the room temperature.

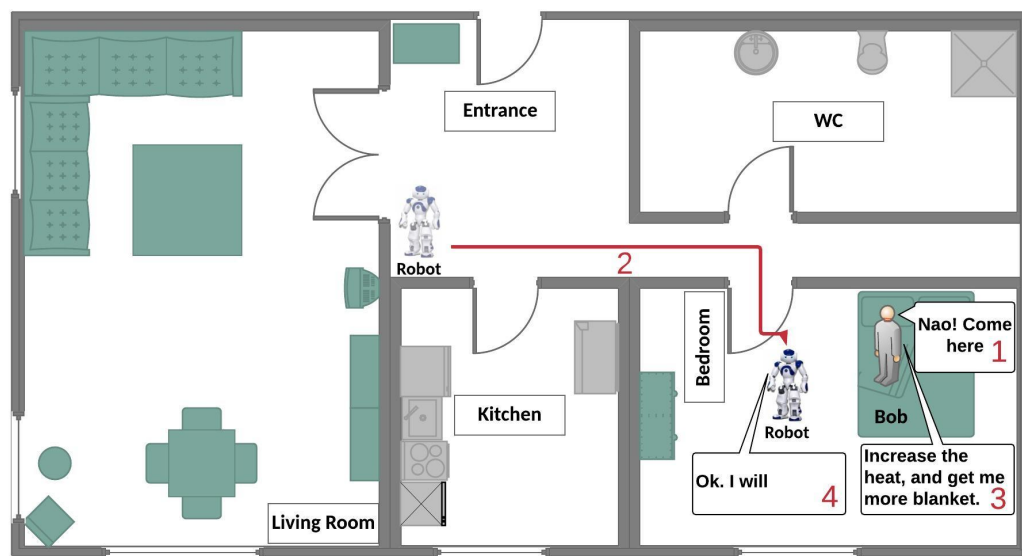


Figure 5.21. Adjust the temperature scenario.

In the Figure 5.21, we demonstrated a standard robots reaction for comparison. The ethical version of the Figure 4.10 is in the part of the scenario (subsection 4.1.2). The standard robot does not inform the user about the change or the current room temperature. The advantage of our system is we are limiting the actions so that if it can result in a risky situation, we inform or warn the user.

Table 5.3. Robot experiment results.

Action Name	Faulty Repeat #	Standard Robot Result
Call the police	2	Infinite
Call the doctor	3	Infinite
Over the counter medicine	1	Infinite
Alert for medicine	-	Infinite
Bring the desired object	-	Infinite
Adjust the temperature	7	Infinite

In the table 5.3, we listed the number of wrong repeats per action to change the result of the ethics module from ethical to unethical. We also compared with the standard robot and possible outcomes of obeying all the orders. You can observe the threshold for wrong decisions increases while the priority of the task decreases. The critical fact is standard robot continuously follow the request from its user because it can not foresee the result of the action. This infinite follow may end up with bad results like giving over-dosage medicine, giving a considerable amount of food.

We provided three different experiment phases, which starts by checking the module itself; the robot-ethics module connects and the end product. We tested all the small parts and nearly all the suggested actions. In the end, we counted the robots limited actions and gathered in tables. From these experiments, we can say that our system is working and acting differently from the robot without the ethics module. Respect for autonomy was the hardest part because even denying the order is violating it. As we said, the robot only denies it, not force the owner to do actions. This is more ethical than forcing the owner.

After the experiments mentioned above, we conducted a questionnaire on the experts. The aim of this questionnaire was comparing our ethics module with experts opinion. The details of the survey are given in the Appendix for two action examples, which are "daily medicine" and "call the police". According to our comparison, the similarity between our module and experts opinion is with the mean of 70% with $\sigma = 3$. Since this is the first ethics module, there is no threshold for success. We assumed higher than 50% as a successful module.

6. CONCLUSION

Robots are starting to be used in all fields of life. They will involve our lives more than we taught in the coming years. Consequently, they need to adjust their behavior in order to suit the social life and create a trust-based association with humans. The design and implementation of ethics rules are demanded the robots to adapt to the social life rules of humans. Numerous Human-Robot Interaction (HRI) studies have been conducted to find the best interaction fashion between the human and the robot.

Unfortunately, the studies and designed laws are not sufficient today. There are a lot of pitfalls and vagueness in most cases. Most of the proposed rules are high-level for current robotics technology, which is nearly impossible to implement. In the future, in 50 years, we are sure that robots will face more complicated states more than our imagination, and they may need a selection of two worst cases when necessary. The famous trolley dilemma can be shown as an example [123]. Would a robot select to kill one person or save five? Or, if it selects to harm, who will be responsible for this behavior? The owner or the programmer?

The need for ethics module is evident. Lots of inter-disciplinary researchers started working on robot ethics, but they can not still have a consensus on "What we expect from an ethical robot ?". Due to conflicts and design issues, there is just a toy world working ethics module, which works with only one case, unlike real-life. Most of the studies are limited to high-level abstract guidelines for ethical robots like BSI, and there is no standardized ethics approach to robot ethics [9].

In our study, we aimed to design and implement an ethics module which can work with many cases that are possible for a service robot. The requirements of a service robot are analyzed, and fuzzy membership graphs are implemented. With weighting the analyzed memberships, we provided an expert system on the top of the robot behavior controller. In the experimentation phase, we observed the distinguished behaviors of the ethical robot in various environments.

We gave three diverse investigation stages, which begins by checking the ethics module itself. Then the robot-ethics module associates in simulation and finally, the real robot experiment result. We tried all the actions towards these experimentation stages. At last, we recorded the robots modified activities and collected in tables. From these examinations, we can say that our framework is working and acting uniquely in contrast to the robot without the ethics module. Regard for user's autonomy was the hardest part because denying the request is damaging it. As we stated, the robot refuses the order and not constrain the user's activities. This is more moral than enforcing the user.

We have succeeded in designing an ethical module which can determine the most ethical behavior among various scenarios. We assumed that this is the first ethical module which is implemented concerning a service robot. This study is the first step towards enhancing our understanding of robot ethics and its challenges. We had a chance to analyze the performance of Unreal Engine as a simulation environment.

Our work clearly has some limitations. Nevertheless, we believe our work could be the starting point for multi-scenario based ethical modules on the service robot. We stated that the ethics module should be transparent and easy-to-follow, but these are qualitative parameters that are hard to evaluate. The present study has only investigated the design of the ethics module, not the evaluation of the unit. Consequently, the FES only evaluated with human experts to approve the output of the unit. Unfortunately, There is a lack of human experts on robot ethics, which can accept or reject this unit. The weak point of the system is its validation.

6.1. Future Work

Many different adaptations, tests, and experiments have been left for the future due to lack of time and infeasibility of current technology (for example the experiments with the real robot are usually very time consuming, requiring even days to finish). Future work concerns a more in-depth analysis of particular ethics approaches, new proposals to try different methods or simply creating a standardized ontology of robot ethics.

Future work should concentrate on enhancing the quality of standardized terms. The robot ethics is an interdisciplinary science which combines, applied ethics, computer science, and legal science. In interdisciplinary themes understanding each other's technical terms and definition is the most fundamental problem. To solve this communication problem, a universal committee or a group should provide documentation about standardized technical terms as also Haidegger underlined [124].

The second problem with robot ethics is the formulation of ethical values. Most of the moral values are qualitative, which are fuzzy to implement and infer their meaning. This is why we used FES in our environment, but, it will be better to have a standard mathematical formalization of de facto ethical values. Most ethics approaches are using common moral values to determine whether a behavior is good or bad. A standard ethical value formalization which formed by an interdisciplinary committee will increase the numbers of proposed ethics module on robot.

As a computer scientist, the hardest point of this research was understanding the fundamentals of ethics and the details of approaches. We needed to read many legal, ethics articles about robot ethics, which contains technical terms which are specific for that science. Also from the different workshops, there are a lots of different views which varies from science to science. For example, legal and ethics scientist assume robot as a high-level system which achieves the perception of its environment ideally without any error and could act according to it without delay. This type of robotic system is challenging to the computer and robotic scientists. There should be a general framework of limits of the robots and achievements of the robots.

The previous suggestions are vital issues for further research, which requires a selected group of people to solve. We formed roadmap for enriching our expert ethics module. First of all, on completion, Boss robot should be tested with our ethics module. Then, we should personalize ethics module for Boss robot according to its abilities. While waiting for BOSS to finish, an exciting research may be optimizing the fuzzy weights by using a machine learning algorithm or conducting a survey about user's preferences. The weights of the states are determined by assuming in an ideal world and with a perfect person. To feed the machine learning algorithm (or a neural network), we need to collect data either from a survey or real robot test. This optimized fuzzy parameters will lead us having a better ethical module.

Another missing point of the FES is validation and verification. Now we are just comparing the behavior differences of the robots, and we assumed that standard robot obeys the order directly. We should provide a test framework and environment which can compare each robotic system according to their ethics module. Also, user experience tests should be prepared for defining the ethical norms for robots and Turkish people. The results of user experience tests can be used for both designing an ethical rule and verifying rules. This test also helps us to measure our success in developing an ethics module and guide us for further performance increases.

For now, we are assuming that the knowledge base of the robot (the memory of the robot) will be robust. In a real robot, it needs to be non-erasable storage, which will not be fragile to disk failure. Moreover, data should not be lost when power off or empty battery. Thus, an ethical robot needs persistent storage technologies, which include non-volatile memories like SSD, EEPROM.

Another crucial point is the encryption of personal data. The robot needs to save some personal data in order to determine ethical behavior and create an environment model. These pieces of information can be dangerous in leaking because of a hacker or any person with bad intention. That person can modify the data against the user or can use the data in a robbery or even murder. The data in the robot should be encrypted for specific doctor usage (for users medical record) and also to prevent remote attacks. According to BSI and also many legal systems, sharing information with the third party is an ethical hazard and contrary to terms of the privacy policy [9].

Finally, for test purposes we assume, all elderly people have same personality which is a very naive approach. The service robots should be personalized like cell phones. Some of the users may like to chat with the robot and act towards them as a friend, or some users just want to treat them like robots for doing their tasks. It is an excellent topic to research further and to see the effect of adding more personalized decisions in the expert system according to the user's life.

REFERENCES

1. Asimov, I., *I, robot*, Vol. 1, Spectra, 2004.
2. Anderson, M. and S. L. Anderson, “Machine Ethics: Creating an Ethical Intelligent Agent”, *AI Magazine*, Vol. 28, No. 4, pp. 15–27, Dec 2007, <http://www.aaai.org/ojs/index.php/aimagazine/article/view/2065>.
3. *International Federation of Robotics*, <https://www.ifr.org/service-robots/>, accessed in April 2019.
4. Prestes, E., J. L. Carbonera, S. Rama Fiorini, V. A. Vitor, M. Abel, R. Madhavan, A. Locoro, P. Goncalves, M. E. Barreto, M. Habib, A. Chibani, S. Gérard, Y. Amirat and C. Schlenoff, “Towards a core ontology for robotics and automation”, *Robotics and Autonomous Systems*, Vol. 61, No. 11, pp. 1193–1204, 2013, <http://dx.doi.org/10.1016/j.robot.2013.04.005>.
5. Scheutz, M., “What Is Robot Ethics? [TC Spotlight]”, *IEEE Robotics & Automation Magazine*, Vol. 20, No. 4, pp. 20–165, Dec 2013, <http://ieeexplore.ieee.org/document/6678596/>.
6. Asaro, P. M., “Remote-Control Crimes”, *IEEE Robotics & Automation Magazine*, Vol. 18, No. 1, pp. 68–71, Mar 2011, <http://ieeexplore.ieee.org/document/5751809/>.
7. Operto, F., “Ethics in Advanced Robotics”, *IEEE Robotics & Automation Magazine*, Vol. 18, No. 1, pp. 72–78, Mar 2011, <http://ieeexplore.ieee.org/document/5751974/>.
8. Veruggio, G., J. Solis and M. Van der Loos, “Roboethics: Ethics Applied to Robotics [From the Guest Editors]”, *IEEE Robotics & Automation Magazine*, Vol. 18, No. 1, pp. 21–22, Mar 2011, <http://ieeexplore.ieee.org/document/5751970/>.

9. British Standards Institution, *BS 8611:2016 - Robots and robotic devices. Guide to the ethical design and application of robots and robotic systems*, BSI, London, 2016, <http://shop.bsigroup.com/ProductDetail?pid=000000000030320089>.
10. Intelligence, H. o. L. S. C. o. A., “AI in the UK: ready, willing and able?”, *House of Lords*, Vol. 36, Apr 2018, <https://publications.parliament.uk/pa/ld201719/ldselect/ldai/100/100.pdf>.
11. Arkin, R., P. Ulam and B. Duncan, *An Ethical Governor for Constraining Lethal Action in an Autonomous System*, Tech. rep., Georgia Institute of Technology, 2009, <https://smartech.gatech.edu/bitstream/handle/1853/31465/09-02.pdf?sequence=1&isAllowed=y>.
12. Nagenborg, M., R. Capurro, J. Weber and C. Pingel, “Ethical regulations on robotics in Europe”, *AI & SOCIETY*, Vol. 22, No. 3, pp. 349–366, Jan 2008, <http://link.springer.com/10.1007/s00146-007-0153-y>.
13. Asimov, I., *I, Robot, Robot series*, Bantam Books, 1950.
14. Riley, G., “CLIPS: A Tool for Building Expert Systems”, *Miscellaneous*, 2004, <http://www.clipsrules.net/http://www.ghg.net/clips/CLIPS.html%5Cnpapers2://publication/uuid/4AF53CBF-5744-422A-8228-734C0EB04A33>, accessed in May 2018.
15. *Unreal Engine*, <https://www.unrealengine.com/en-US/>, accessed in May 2019.
16. *ROS.org | Powering the world's robots*, <https://www.ros.org/>, accessed in May 2019.
17. *NAO Documentation — Aldebaran 2.1.4.13 documentation*, http://doc.aldebaran.com/2-1/home_ao.html, accessed in May 2019.
18. Özcan, I. and H. L. Akın, “22 Serbestlik Dereceli Robot Kafası BUSRA”, pp. 188–193,

2018.

19. *RoboCup@Home*, <http://www.robocupathome.org/>, accessed in May 2019.
20. *DARPA Robotics Challenge (DRC)*, <https://www.darpa.mil/program/darpa-robotics-challenge>, accessed in May 2019.
21. *BS ISO 8373:2012, Robots and robotic devices. Vocabulary*, BSI, 2012.
22. International Federation of Robotics, *Executive Summary World Robotics 2018 Service Robots*, Tech. rep., 2018, https://ifr.org/downloads/press2018/Executive_Summary_WR_Service_Robots_2018.pdf, accessed in April 2019.
23. *Robot density rises globally - International Federation of Robotics*, <https://ifr.org/ifr-press-releases/news/robot-density-rises-globally>, accessed in May 2019.
24. *Mojin Robotics*, <https://mojin-robotics.de/produkt/>, accessed in May 2019.
25. *Wall-YE*, <https://wall-ye.com/>, accessed in May 2019.
26. *The strong robot with the gentle touch | RIKEN*, http://www.riken.jp/en/pr/press/2015/20150223_2/, accessed in May 2019.
27. Lin, P., K. Abney and G. A. Bekey, *Robot ethics : the ethical and social implications of robotics*, MIT Press, 2012.
28. Singer, P., *Applied ethics*, Oxford University Press, 1986.
29. Hursthouse, R., *On virtue ethics*, OUP Oxford, 1999.
30. Darwell, S., *Deontology*, Wiley-Blackwell, 2008.
31. Immanuel, K., “Groundwork of the Metaphysics of Morals”, E. S. Radcliffe, R. Mc-

- Carty, F. Allhoff and A. Vaidya (Editors), *Late Modern Philosophy: Essential Readings with Commentary*, Blackwell, 1785.
32. Veruggio, G., “EURON roboethics roadmap (release 1.1)”, *EURON Roboethics Atelier, Genua*, 2006.
 33. Anderson, M. and S. L. Anderson, “Ethical Healthcare Agents”, *Advanced Computational Intelligence Paradigms in Healthcare - 3*, Vol. 107, pp. 233–257, Springer Berlin Heidelberg, Berlin, Heidelberg, 2008, http://link.springer.com/10.1007/978-3-540-77662-8_10.
 34. Mill, J. S., *Utilitarianism*, Longmans, Green and Company, 1895.
 35. Bentham, J., *An introduction to the principles of morals and legislation*, Clarendon Press, 1879.
 36. Anderson, M. and S. L. Anderson, “GenEth : A General Ethical Dilemma Analyzer”, *Twenty-Eighth AAAI Conference on Artificial Intelligence*, pp. 253–261, 2014.
 37. Rawls, J., *A theory of justice*, Harvard university press, 1971.
 38. Medlin, B., “Ultimate principles and ethical egoism”, *Australasian Journal of Philosophy*, Vol. 35, No. 2, pp. 111–118, 1957.
 39. Dewey, J., “Theory of valuation.”, *International encyclopedia of unified science*, 1939.
 40. van Wynsberghe, A., “Designing robots for care: care centered value-sensitive design.”, *Science and engineering ethics*, Vol. 19, No. 2, pp. 407–33, Jun 2013, <http://www.ncbi.nlm.nih.gov/pubmed/22212357><http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC3662860>.
 41. Cheon, E. and N. M. Su, “Integrating roboticist values into a Value Sensitive Design framework for humanoid robots”, *2016 11th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*, pp. 375–382, IEEE, 2016.

42. Jonsen, A. R. and S. E. Toulmin, *The abuse of casuistry: A history of moral reasoning*, Univ of California Press, 1988.
43. Anderson, M., S. L. Anderson and C. Armen, “MedEthEx : A Prototype Medical Ethics Advisor”, *Proceedings Of The National Conference On Artificial Intelligence*, May, pp. 1759–1765, AAAI Press;MIT Press, 2006, https://www.researchgate.net/publication/221604863_MedEthEx_A_Prototype_Medical_Ethics_Advisor.
44. Salge, C. and D. Polani, “Empowerment As Replacement for the Three Laws of Robotics”, *Frontiers in Robotics and AI*, Vol. 4, No. June, pp. 1–16, Jun 2017, <http://journal.frontiersin.org/article/10.3389/frobt.2017.00025/full>.
45. Wallach, W. and C. Allen, *Moral machines: Teaching robots right from wrong*, Oxford University Press, 2008.
46. Hughes, J., “Compassionate AI and selfless robots: A buddhist approach”, *Robot ethics: the ethical and social implications of robotics*, pp. 69–83, 2012.
47. Murphy, M., “Theological Voluntarism”, E. N. Zalta (Editor), *The Stanford Encyclopedia of Philosophy*, Metaphysics Research Lab, Stanford University, summer 2019 edn., 2019.
48. Quinn, P. L., *Divine commands and moral requirements*, Clarendon Press, 1978.
49. Chisholm, R., “Practical reason and the logic of requirement”, *Practical reason*, pp. 2–13, 1974.
50. Bringsjord, S. and J. Taylor, “Introducing divine-command robot ethics”, *Robot ethics: the ethical and social implication of robotics*, 2012.
51. Jackson, P., *Introduction to expert systems*, Addison-Wesley Longman Publishing Co., Inc., 1998.
52. Leondes, C. T., *Expert systems: the technology of knowledge management and decision*

making for the 21st century, Elsevier, 2001.

53. AI: *Timeline*, <http://biology.kenyon.edu/slonc/bio3/AI/TIMELINE/timeline.html>, accessed in June 2019.
54. Reid, S., “Knowledge-based systems concepts, Techniques, Examples”, *Schlumberger-Doll Research*, 1985.
55. Crespi, V., A. Galstyan and K. Lerman, “Top-down vs bottom-up methodologies in multi-agent system design”, *Autonomous Robots*, Vol. 24, No. 3, pp. 303–313, 2008.
56. Tzafestas, S. G., “Roboethics”, *A Navigating Overview. Heidelberg*, 2016.
57. Boden, M., J. Bryson, D. Caldwell, K. Dautenhahn, L. Edwards, S. Kember, P. Newman, V. Parry, G. Pegman, T. Rodden, T. Sorrell, M. Wallis, B. Whitby and A. Winfield, “Principles of robotics: regulating robots in the real world”, *Connection Science*, Vol. 29, No. 2, pp. 124–129, Apr 2017, <https://www.tandfonline.com/doi/full/10.1080/09540091.2016.1271400>.
58. Asilomar, A., “Principles.(2017)”, *Principles developed in conjunction with the 2017 Asilomar conference [Benevolent AI 2017]*, 2018, <https://futureoflife.org/ai-principles/>.
59. *Principles for Algorithmic Transparency and Accountability*, Tech. rep., https://www.acm.org/binaries/content/assets/public-policy/2017_usacm_statement_algorithms.pdf.
60. *The Japanese Society for Artificial Intelligence Ethical Guidelines*, Tech. rep., <http://ai-elsi.org/wp-content/uploads/2017/05/JSAI-Ethical-Guidelines-1.pdf>.
61. Economou, N., “A ‘principled’ artificial intelligence could improve justice”, pp. 2011–2013, 2018, http://www.abajournal.com/legalrebels/article/a_

principled_artificial_intelligence_could_improve_justice.

62. Naveen Rao, *Intel Recommends Public Policy Principles for Artificial Intelligence - Policy@Intel*, 2017, <https://blogs.intel.com/policy/2017/10/18/naveen-rao-announces-intel-ai-public-policy/#gs.h8n96h>, accessed in June 2019.
63. *General Principles The IEEE Global Initiative on Ethics of Autonomous and Intelligent Systems*, Tech. rep., https://standards.ieee.org/content/dam/ieee-standards/standards/web/documents/other/ead_general_principles_v2.pdf.
64. Union, U. G., “Top 10 Principles for Ethical Artificial Intelligence”, *Nyon, Switzerland*, 2017.
65. Dávid Pásztor, *AI UX: 7 Principles of Designing Good AI Products*, 2018, <https://uxstudioteam.com/ux-blog/ai-ux/>, accessed in June 2019.
66. Now, A., *The Toronto Declaration: Protecting the rights to equality and non-discrimination in machine learning systems.(May 2018)*, 2018, https://www.accessnow.org/cms/assets/uploads/2018/08/The-Toronto-Declaration_ENG_08-2018.pdf, accessed in May 2019.
67. Pichai, S., “AI at Google: our principles”, *The Keyword, Google, June*, Vol. 7, 2018.
68. Cutler, A., M. Pribić and L. Humphrey, *Everyday Ethics for Artificial Intelligence*, 2018, <https://www.ibm.com/watson/assets/duo/pdf/everydayethics.pdf>, accessed in June 2019.
69. Microsoft, *Responsible bots: 10 guidelines for developers of conversational AI*, 2018, https://www.microsoft.com/en-us/research/uploads/prod/2018/11/Bot_Guidelines_Nov_2018.pdf, accessed in June 2019.
70. CEPEJ, *Council of Europe adopts first European Ethical Charter on*

- the use of artificial intelligence in judicial systems* - Newsroom, 2018, <https://www.coe.int/en/web/human-rights-rule-of-law/-/council-of-europe-adopts-first-european-ethical-charter-on-the-use-of-artificial-intelligence-in-judicial-systems>, accessed in June 2019.
71. Women Leading in AI, *10 Principles for Responsible AI*, Tech. rep., 2019, <https://womenleadinginai.org/report2019https://womenleadinginai.org/wp-content/uploads/2019/02/WLiAI-Report-2019.pdf>.
 72. Department of Health and Social care, “Initial code of conduct for data-driven health and care technology - GOV.UK”, Vol. 2016, No. 2018, pp. 1–62, 2018, <https://www.gov.uk/government/publications/code-of-conduct-for-data-driven-health-and-care-technology/initial-code-of-conduct-for-data-driven-health-and-care-technology>.
 73. The IEEE Global Initiative on Ethics of Autonomous and Intelligent Systems, *Ethically Aligned Design: A Vision for Prioritizing Human Well-being with Autonomous and Intelligent Systems*, 2019, <https://standards.ieee.org/content/dam/ieee-standards/standards/web/documents/other/eadle-overview.pdf>.
 74. The Biometrics and Forensics Ethics Group and (BFEG), *Ethical issues arising from the police use of live facial recognition*, Tech. rep., 2019, https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/781745/Facial_Recognition_Briefing_BFEG_February_2019.pdf.
 75. Floridi, L., J. Cows, M. Beltrametti, R. Chatila, P. Chazerand, V. Dignum, C. Luetge, R. Madelin, U. Pagallo, F. Rossi, B. Schafer, P. Valcke and E. Vayena, “AI4People—An Ethical Framework for a Good AI Society: Opportunities, Risks, Principles, and Recommendations”, *Minds and Machines*, Vol. 28, No. 4, pp. 689–707, Dec 2018, <https://doi.org/10.1007/s11023-018-9482-5>.
 76. European Commission, “High-Level Expert Group on Artificial Intelligence Set Up

- By the European Commission Ethics Guidelines for Trustworthy Ai”, pp. 1–41, 2019, <https://ec.europa.eu/futurium/en/ai-alliance-consultation/guidelines/1>.
77. Dawson, D. and E. Schleiger, *Artificial Intelligence: Australia’s Ethics Framework*, Tech. rep., 2019, <https://consult.industry.gov.au/>.
 78. Murphy, R. and D. D. Woods, “Beyond Asimov: The Three Laws of Responsible Robotics”, *IEEE Intelligent Systems*, Vol. 24, No. 4, pp. 14–20, Jul 2009, <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=5172885>.
 79. Garfinkel, S., J. Matthews, S. S. Shapiro and J. M. Smith, “Toward algorithmic transparency and accountability”, *Communications of the ACM*, Vol. 60, No. 9, pp. 5–5, 2017.
 80. The Ethics Committee, t. J. S. f. A. I. J., *About the Japanese Society for Artificial Intelligence Ethical Guidelines*, 2017, <http://ai-elsi.org/archives/514>, accessed in June 2019.
 81. *MONTREAL DECLARATION FOR A RESPONSIBLE DEVELOPMENT OF ARTIFICIAL INTELLIGENCE 2018*, Tech. rep., https://docs.wixstatic.com/ugd/ebc3a3_c5c1c196fc164756afb92466c081d7ae.pdf.
 82. Winfield, A. and M. Halverson, “Artificial Intelligence and Autonomous Systems: Why Principles Matter”, *IEEE*, 2017, <https://cmte.ieee.org/futuredirections/tech-policy-ethics/september-2017/artificial-intelligence-and-autonomous-systems-why-principles-matter/>.
 83. Smith, R., “5 core principles to keep AI ethical”, *World Economic Forum*, Vol. 19, 2018.
 84. Sharkey, N., *Women Stand Against Social Injustice In AI*, 2019, <https://www.forbes.com/sites/noelsharkey/2019/02/07/women-stand-against->

social-injustice-in-ai/{\#}3f2f56fd2118, accessed in June 2019.

85. *The NHS's 10 Principles for AI + Data, A New Benchmark for Lawyers? – Artificial Lawyer*, 2019, <https://www.artificiallawyer.com/2019/02/22/the-nhss-10-principles-for-ai-data-a-new-benchmark-for-lawyers/>, accessed in June 2019.
86. Moor, J., “The Nature, Importance, and Difficulty of Machine Ethics”, *IEEE Intelligent Systems*, Vol. 21, No. 4, pp. 18–21, Jul 2006, <http://ieeexplore.ieee.org/document/1667948/>.
87. Anderson, M., S. L. Anderson and C. Armen, “Towards Machine Ethics”, *Proceedings of AAAI 2004 Workshop on Agent Organizations: Theory and Practice*, July, pp. 1–7, 2004, https://www.researchgate.net/publication/259656154_Towards_Machine_Ethics.
88. Ross, W., *The Right and The Good*, Oxford: Clarendon Press, 1930.
89. Beauchamp, T. L. and J. F. Childress, *Principles of biomedical ethics*, Oxford University Press, USA, 1979.
90. Anderson, M., S. L. Anderson and C. Armen, “An approach to computing ethics”, *IEEE Intelligent Systems*, Vol. 21, No. 4, pp. 56–63, 2006.
91. Anderson, M. and S. L. Anderson, “Developing a General , Interactive Approach to Codifying Ethical Principles”, *Proceedings of AAAI Workshop Human Implications of Human-Robot Interaction*, January, pp. 2–10, Chicago, 2008, <https://www.aaai.org/Papers/Workshops/2008/WS-08-05/WS08-05-002.pdf>.
92. Anderson, S. L. and M. Anderson, “A prima facie duty approach to machine ethics and its application to elder care”, *Workshops at the Twenty-Fifth AAAI Conference on Artificial Intelligence*, 2011, www.aaai.org/ocs/index.php/WS/AAAIW11/paper/viewPaper/3812.

93. Anderson, S. L. and M. Anderson, "The Relationship Between Intelligent, Autonomously Functioning Machines and Ethics", *Proceedings of the 2013 meeting of the international association for computing and philosophy.*, pp. 251–266, 2013, http://www.iacap.org/proceedings_IACAP13/paper_3.pdf.
94. Anderson, M. and S. L. Anderson, "Toward ethical intelligent autonomous healthcare agents : a case-supported principle-based behavior paradigm", *50th Annual Convention of the Society for the Study of Artificial Intelligence and the Simulation of Behaviour (AISB-50) Symposium on Machine Ethics in the Context of Medical and Care Agents*, April, pp. 2–6, London, UK, 2014, https://www.researchgate.net/publication/261392841_Toward_Ethical_Intelligent_Autonomous_Healthcare_Agents_A_Case-Supported_Principle-Based_Behavior_Paradigm.
95. Anderson, M. and S. L. Anderson, "Toward ensuring ethical behavior from autonomous systems: a case-supported principle-based paradigm", *Industrial Robot: An International Journal*, Vol. 42, No. 4, pp. 324–331, Jun 2015, <http://www.emeraldinsight.com/doi/10.1108/IR-12-2014-0434>.
96. Asaro, P., "What Should We Want From a Robot Ethic?", *International Review of Information Ethics (IRIE)*, Vol. 6, pp. 9–16, 2006.
97. Borenstein, J. and Y. Pearson, "Robot caregivers: harbingers of expanded freedom for all?", *Ethics and Information Technology*, Vol. 12, No. 3, pp. 277–288, Sep 2010, <http://link.springer.com/10.1007/s10676-010-9236-4>.
98. Borenstein, J., "Robotics, Ethics, and the Environment", *International Journal of Technoethics*, Vol. 3, No. 2, pp. 17–29, Apr 2012, <http://services.igi-global.com/resolvedoi/resolve.aspx?doi=10.4018/jte.2012040103>.
99. Dabringer, G., *Ethical and Legal Aspects of Unmanned Systems Interviews*, Institut für Religion und Frieden, Wien, 2010.

100. Decker, M., R. Dillmann, T. Dreier, M. Fischer, M. Gutmann, I. Ott and I. Spiecker genannt Döhmann, “Service robotics: do you know your new companion? Framing an interdisciplinary technology assessment”, *Poiesis & Praxis*, Vol. 8, No. 1, pp. 25–44, Nov 2011, <http://link.springer.com/10.1007/s10202-011-0098-6>.
101. Deng, B., “Machine ethics: The robot’s dilemma”, *Nature*, Vol. 523, No. 7558, pp. 24–26, Jul 2015, <http://www.nature.com/doi/10.1038/523024a>.
102. Dreier, T. and I. Spiecker genannt Döhmann, “Legal aspects of service robotics”, *Poiesis & Praxis*, Vol. 9, No. 3-4, pp. 201–217, Dec 2012, <http://link.springer.com/10.1007/s10202-012-0115-4>.
103. Kernaghan, K., “The rights and wrongs of robotics: Ethics and robots in public organizations”, *Canadian Public Administration*, Vol. 57, No. 4, pp. 485–506, Dec 2014, <http://onlinelibrary.wiley.com/doi/10.1111/capa.12093/full><http://doi.wiley.com/10.1111/capa.12093>.
104. Sharkey, N., “The Ethical Frontiers of Robotics”, *Science*, Vol. 322, No. 5909, pp. 1800–1801, Dec 2008, <http://www.sciencemag.org/lookup/doi/10.1126/science.1164582>.
105. Sharkey, A. and N. Sharkey, “Children, the Elderly, and Interactive Robots”, *IEEE Robotics & Automation Magazine*, Vol. 18, No. 1, pp. 32–38, Mar 2011, <http://ieeexplore.ieee.org/document/5751987/>.
106. Winfield, A. F. T., C. Blum and W. Liu, “Towards an Ethical Robot: Internal Models, Consequences and Ethical Action Selection”, *Conference towards autonomous robotic systems*, pp. 85–96, Springer, 2014, http://link.springer.com/10.1007/978-3-319-10401-0_8.
107. Vanderelst, D. and A. Winfield, “The Dark Side of Ethical Robots”, *2018 AAAI/ACM Conference on AI, Ethics, and Society*, pp. 317–322, Jun 2016, <http://arxiv.org/abs/1606.02583>.

108. Vanderelst, D. and A. Winfield, “An architecture for ethical robots inspired by the simulation theory of cognition”, *Cognitive Systems Research*, Vol. 48, pp. 56–66, May 2018, <https://www.sciencedirect.com/science/article/pii/S1389041716302005?via%3Dihub>.
109. Belloni, A., A. Berger, V. Besson, O. Boissier, G. Bonnet, G. Bourgne, P. A. Chardel, J.-P. Cotton, N. Evreux, J.-G. Ganascia, P. Jaillon, B. Mermet, G. Picard, B. Reber, G. Simon, T. de Swarte, C. Tessier, F. Vexler, R. Voyer and A. Zimmermann, “Towards A Framework To Deal With Ethical Conflicts In Autonomous Agents And Multi - Agent Systems”, pp. Paper–8, Jun 2014, <https://hal-emse.ccsd.cnrs.fr/emse-01059503/>.
110. Belloni, A., A. Berger, O. Boissier, G. Bonnet, G. Bourgne, P.-A. Chardel, J.-P. Cotton, N. Evreux, J.-G. Ganascia, P. Jaillon *et al.*, “Dealing with ethical conflicts in autonomous agents and multi-agent systems”, *Workshops at the Twenty-Ninth AAAI Conference on Artificial Intelligence*, 2015.
111. Collart, J., T. Gateau, E. Fabre and C. Tessier, “Human-robot systems facing ethical conflicts: a preliminary experimental protocol”, *Workshops at the Twenty-Ninth AAAI Conference on Artificial Intelligence*, 2015.
112. Tessier, C., “Robots Autonomy: Some Technical Challenges”, *2015 AAAI Spring Symposium Series*, 2015.
113. Mansouri, N. and K. Goher, “Towards ethical framework for personal care robots: Review and reflection”, Vol. 12, No. 10, pp. 152–162, Sep 2016, <http://www.ccsenet.org/journal/index.php/ass/article/view/61588>.
114. Winfield, A. F. and M. Jirotko, “Ethical governance is essential to building trust in robotics and artificial intelligence systems”, *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, Vol. 376, No. 2133, p. 20180085, Nov 2018, <http://rsta.royalsocietypublishing.org/lookup/doi/10.1098/rsta.2018.0085>.

115. Krzanowski, R. and K. Trombik, “Machine Ethics Test (MET) – Concerns , Proposals , Solutions”, *IS4SI 2017 Summit DIGITALISATION FOR A SUSTAINABLE SOCIETY*, June, Gothenburg, 2017, https://www.researchgate.net/publication/318969677_Machine_Ethics_Test_MET_-_Concerns_Proposals_Solutions.
116. Python Software Foundation, *Welcome to Python.org*, 2017, <https://www.python.org/>, accessed in June 2019.
117. Pot, E., J. Monceaux, R. Gelin and B. Maisonnier, “Choregraphe: a graphical tool for humanoid robot programming”, *RO-MAN 2009-The 18th IEEE International Symposium on Robot and Human Interactive Communication*, pp. 46–51, IEEE, 2009.
118. Bugmann, G. and S. N. Copleston, “What can a personal robot do for you?”, *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, Vol. 6856 LNAI, pp. 360–371, Springer, Berlin, Heidelberg, 2011, http://link.springer.com/10.1007/978-3-642-23232-9_32.
119. Beer, J. M., C.-A. Smarr, T. L. Chen, A. Prakash, T. L. Mitzner, C. C. Kemp and W. A. Rogers, “The domesticated robot: Design Guidelines for Assisting Older Adults to Age in Place”, *Proceedings of the seventh annual ACM/IEEE international conference on Human-Robot Interaction - HRI '12*, p. 335, ACM Press, New York, New York, USA, 2012, <http://dl.acm.org/citation.cfm?doid=2157689.2157806>.
120. Zadeh, L. A., “Fuzzy sets”, *Information and control*, Vol. 8, No. 3, pp. 338–353, 1965.
121. Agassi, J. and J. Wiezenbaum, “Computer Power and Human Reason: From Judgment to Calculation”, *Technology and Culture*, Vol. 17, No. 4, p. 813, Oct 1976.
122. Loper, E. and S. Bird, “NLTK: the natural language toolkit”, *arXiv preprint cs/0205028*, 2002.
123. Foot, P., “The problem of abortion and the doctrine of double effect”, Vol. 5, pp. 5–15,

1967.

124. Haidegger, T., M. Barreto, P. Gonçalves, M. K. Habib, S. K. V. Ragavan, H. Li, A. Vaccarella, R. Perrone and E. Prestes, “Applied ontologies and standards for service robots”, *Robotics and Autonomous Systems*, Vol. 61, pp. 1215–1223, North-Holland, Nov 2013, <https://www.sciencedirect.com/science/article/pii/S092188901300105X?via%3Dihub>.

APPENDIX A: EXPERT QUESTIONNAIRE FORM

In this section, we provided the questionnaire that we sent to our experts. For each scenario, we created a list of parameters that changed in the case and presented the robot's action. We expected from our experts to fill the result section either the case is ethical or not. This questionnaire contains all possible outcomes of the case so the expert can not understand what our modules result is and examined the system objectively.

Here we provided two example actions which are also selected for FES calculation. The table A.1 and A.2 contain these actions' questionnaire.

Table A.1. Questionnaire example for "call the police" scenario.

Order	Situation Parameters	Robot Actions	Result
Call the police	Emergency + user trust value >50	Calls	
	Emergency + user trust value <50	Calls	
	Uncertain + user trust value >50	Calls	
	Uncertain + user trust value <50	Calls	
	Not Emergency + user trust value >50	Calls	
	Not Emergency + user trust value <50	Calls	
	Emergency + user trust value >50	Not calls	
	Emergency + user trust value <50	Not calls	
	Uncertain + user trust value >50	Not calls	
	Uncertain + user trust value <50	Not calls	
	Not Emergency + user trust value >50	Not calls	
	Not Emergency + user trust value <50	Not calls	

Table A.2. Questionnaire example for "bring daily medicine" scenario.

Order	Situation Parameters	Robot Actions	Result
Bring Daily Medicine	Under Dosage + Early + user trust value >50	Brings Medicine	
	Under Dosage + Early + user trust value <50	Brings Medicine	
	Under Dosage + On Time + user trust value >50	Brings Medicine	
	Under Dosage + On Time + user trust value <50	Brings Medicine	
	Under Dosage + Late+ user trust value >50	Brings Medicine	
	Under Dosage + Late + user trust value <50	Brings Medicine	
	Normal Dosage + Early + user trust value >50	Brings Medicine	
	Normal Dosage + Early + user trust value <50	Brings Medicine	
	Normal Dosage + On Time + user trust value >50	Brings Medicine	
	Normal Dosage + On Time + user trust value <50	Brings Medicine	
	Normal Dosage + Late+ user trust value >50	Brings Medicine	
	Normal Dosage + Late + user trust value <50	Brings Medicine	
	Over Dosage + Early + user trust value >50	Brings Medicine	
	Over Dosage + Early + user trust value <50	Brings Medicine	
	Over Dosage + On Time + user trust value >50	Brings Medicine	
	Over Dosage + On Time + user trust value <50	Brings Medicine	
	Over Dosage + Late+ user trust value >50	Brings Medicine	
	Over Dosage + Late + user trust value <50	Brings Medicine	

Table A.2. Questionnaire example for "bring daily medicine" scenario (cont.)

Order	Situation Parameters	Robot Actions	Result
Bring Daily Medicine	Under Dosage + Early + user trust value >50	Not Brings Medicine	
	Under Dosage + Early + user trust value <50	Not Brings Medicine	
	Under Dosage + On Time + user trust value >50	Not Brings Medicine	
	Under Dosage + On Time + user trust value <50	Not Brings Medicine	
	Under Dosage + Late+ user trust value >50	Not Brings Medicine	
	Under Dosage + Late + user trust value <50	Not Brings Medicine	
	Normal Dosage + Early + user trust value >50	Not Brings Medicine	
	Normal Dosage + Early + user trust value <50	Not Brings Medicine	
	Normal Dosage + On Time + user trust value >50	Not Brings Medicine	
	Normal Dosage + On Time + user trust value <50	Not Brings Medicine	
	Normal Dosage + Late+ user trust value >50	Not Brings Medicine	
	Normal Dosage + Late + user trust value <50	Not Brings Medicine	
	Over Dosage + Early + user trust value >50	Not Brings Medicine	
	Over Dosage + Early + user trust value <50	Not Brings Medicine	
	Over Dosage + On Time + user trust value >50	Not Brings Medicine	
	Over Dosage + On Time + user trust value <50	Not Brings Medicine	
	Over Dosage + Late+ user trust value >50	Not Brings Medicine	
	Over Dosage + Late + user trust value <50	Not Brings Medicine	