



**UNSW**  
AUSTRALIA

**School of Computer Science and Engineering**

**Faculty of Engineering**

**The University of New South Wales**

# **Smart Home Automation Through Voice and Gesture Control**

by

**Jake Palandri**

Thesis submitted as a requirement for the degree of  
Bachelor of Engineering in Software Engineering

Submitted: December 2024

Supervisor: Prof. Claude Sammut

Student ID: z5313097

# Abstract

With the ever-increasing prevalence of smart devices in the home, the need for a more intuitive and efficient way to control these devices has become apparent. This thesis aims to develop a system that uses computer vision techniques to recognise gestures, in conjunction with voice commands, to control a variety of smart appliances.

Existing literature in the space of smart home automation and computer vision has been reviewed to inform the design and development of the system. Existing gaps within the literature consisted of gaps within the integration of smart home technologies with computer vision techniques for home automation. There is a significant amount of research into human action recognition and behavioural monitoring in ambient assisted living environments for the care of the elderly and disabled but none of which covers the integration of this technology into smart home automation.

The system has been designed to be user-friendly and accessible to all members of the household, regardless of their technical expertise, which was another gap noted in existing deployments. The system has been evaluated based on its accuracy, ease of use, and overall effectiveness in controlling smart devices. The results of this evaluation were used to inform future research and development in the field of smart home automation.

The results indicated that the system was able to accurately recognise gestures with 97% accuracy and voice commands with 68% accuracy under normal conditions. With background speaking, the voice command accuracy dropped to less than 10%. The system was able to control a variety of smart devices, including lights, and televisions, with the ability to expand to control other devices in the future without significant work from the users.

# Acknowledgements

The author acknowledges the assistance from their supervisor, Professor Claude Sammut, as being integral to their progress across the thesis. The author also acknowledges the use of this template for L<sup>A</sup>T<sub>E</sub>X, graciously provided by the School of Computer Science and Engineering.

# Abbreviations

**AAL** Ambient assisted living

**AC** Alternating Current

**AI** Artificial Intelligence

**ASCII** American Standard Code for Information Interchange

**BMS** Behaviour Monitoring System

**CSA** Connectivity Standards Alliance

**CV** Computer Vision

**DVD** Digital Video Disc

**Distro** Distribution

**ECHO IV** Electronic Computing Home Operator

**Hz** Hertz

**IAI** Institute of Artificial Intelligence, University of Bremen

**IoT** Internet of Things

**IR** Infrared

**JS** JavaScript

**LED** Light Emitting Diode

**ML** Machine Learning

**MQTT** Messaging Queuing Telemetry Transport

**PIR** Passive Infrared

**rclpy** ROS Client Library for Python

**Repo** Repository

**RGB** Red, Green and Blue

**RGB-D** Red, Green, Blue and Depth

**ROS** Robot Operating System

**STT** Speech To Text

**TS** TypeScript

**TV** Television

**USD** United States Dollars

**VCR** Videocassette Recorder

**Wi-Fi** Wireless Fidelity

**YOLO** You Only Look Once

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Literature Review</b>	<b>3</b>
2.1	Background . . . . .	3
2.1.1	History of the Smart Home . . . . .	3
2.1.2	Problems with Smart Homes . . . . .	5
2.1.2.1	Interoperability . . . . .	6
2.1.2.2	Lack of True Intelligence . . . . .	7
2.1.2.3	Requirement for Internet Access . . . . .	7
2.2	Review of Existing Literature . . . . .	8
2.2.1	Computer Vision for Smart Home Automation . . . . .	8
2.2.2	Human Action Recognition for Smart Home Automation . . . . .	10
2.2.3	Gesture Recognition for Smart Home Automation . . . . .	11
2.2.4	Behaviour Monitoring for Elderly Care . . . . .	13
2.2.5	Voice Control and Multi-Modal Smart Home Automation . . . . .	15
2.3	Gaps in Literature . . . . .	16
2.4	Problem Statement . . . . .	17
2.5	Aims and Outcomes . . . . .	18

<b>3 Methodology</b>	<b>20</b>
3.1 Proposed System . . . . .	20
3.1.1 Technology Stack . . . . .	21
3.1.1.1 Thesis A . . . . .	21
3.1.1.2 Thesis B . . . . .	21
3.1.1.3 Thesis C . . . . .	21
3.2 Hardware Components . . . . .	25
3.3 Software Components . . . . .	27
3.4 Component Interactions . . . . .	29
3.5 Code Logic . . . . .	31
<b>4 Evaluation</b>	<b>34</b>
4.1 Solved Problems . . . . .	34
4.2 System Evaluation . . . . .	36
4.2.1 Quantitative Analysis . . . . .	38
4.2.2 Qualitative Analysis . . . . .	40
<b>5 Conclusion and Future Work</b>	<b>41</b>
5.1 Future Work . . . . .	41
5.2 Conclusion . . . . .	42
<b>References</b>	<b>44</b>
<b>Appendix</b>	<b>46</b>
A.1 User Manual . . . . .	46
A.1.1 Setup Guide . . . . .	46
A.1.2 Execution Guide . . . . .	47

# List of Figures

2.1	Extracted Hand Gestures from Microsoft Kinect (Desai & Desai, 2017) . . . . .	12
2.2	Room-to-room State Transition Model (Eisa & Moreira, 2017). . . . .	14
3.1	Simplified Technology Stack . . . . .	20
3.2	Thesis A Proposed Technology Stack . . . . .	22
3.3	Thesis B Implemented Technology Stack . . . . .	23
3.4	Thesis C Final Technology Stack . . . . .	24
3.5	Home Assistant Interface . . . . .	26
3.6	Physical Setup Diagram . . . . .	27
3.7	Kinect to ROS 2 Layers . . . . .	30
3.8	Command Web Interface . . . . .	31
4.1	User 1 Testing . . . . .	37
4.2	User 2 Testing . . . . .	37
4.3	User 3 Testing . . . . .	38

# List of Tables

2.1	Abnormal Behaviours—Descriptions and related health-declines (Eisa & Moreira, 2017). . . . .	13
4.1	User Testing Results . . . . .	38
4.2	User Testing Results . . . . .	39



# Chapter 1

## Introduction

Smart home technology has seen an exponential rate of growth over the last century. The technology has evolved from simple home appliances, to automatic control systems for a wide range of home devices. In particular, each iteration of new smart home technology has evolved to be more user-friendly and more integrated with the user's daily life, freeing up time for more leisure activities, while automating menial tasks.

This report will cover existing literature and the investigation into the future of smart home technology, considering computer vision and human action recognition techniques for the purpose of home automation. Additionally, the implementation of custom voice controlled smart home devices will be discussed, as well as the potential for multi-modal control of smart home devices.

Some of the key challenges that have been identified include the interoperability of smart home devices, the requirement for constant internet access, and the need for a more accessible and customisable interface for users without technical expertise. The system that will be developed in this thesis aims to address these challenges by allowing users robust control over the voice commands and devices available to them, without the need to modify the code for the system with each new device that is added to the home.

Each chapter of the report covers the topics as follows:

- Chapter 2 explains the background for this thesis and related research.
- Chapter 3 outlines the functionality of the system and how each component interacts with the others.
- Chapter 4 analyses the problems that the implemented system solves and discusses the results of user testing.
- Chapter 5 summarises the report.
- Chapter A provides a user manual for the system so that it can be replicated by others.

## Chapter 2

# Literature Review

### 2.1 Background

#### 2.1.1 History of the Smart Home

The evolution of smart homes can be traced back to the early 1900s when the introduction of various electric household appliances promised to reduce the time spent on mundane household chores and free up more time for leisure, revolutionising domestic life. Devices such as the electric mixer and iron, the world's first refrigerator in 1913, and the pop-up toaster in 1919, emerged in this period, bringing about the beginnings of a more automated home environment, with the following decades innovating accordingly (Hertzmann, 2016).

In 1966, there was a significant leap forward in home automation technology when Jim Sutherland, an engineer from Pittsburgh, Pennsylvania created the ECHO IV. The Electronic Computing Home Operator (ECHO IV) was the first device designed specifically for home automation and was hand-crafted with surplus electronic parts (Spicer, 2016). With the ability to compute shopping lists, control home temperature, limit children's television time with questionnaires and even tell the weather, this was the first glimpse at a whole-home computerised system that was capable of automating

every element of home life.

The following year saw the introduction of Honeywell's Kitchen Computer. A computer designed for housewives to use in the kitchen as a recipe storage device with a built-in chopping board. At \$10,000 United States Dollars (USD) in 1967 (nearly \$100,000 USD today when accounting for inflation), this recipe storage device which had no display, was well described as "amazingly beautiful and hopelessly impractical" (*From SCI-FI dreams to everyday reality. Tracing the evolution of smart home through history.*, 2023; Stein, 2011; *US Inflation Calculator*, n.d.). No units were ever sold. Despite its failure on the open market, the device received a lot of attention and the public began to imagine a future where homes would be interactive.

Eight years later, in 1975, a group of engineers from Scotland released the X10 protocol. Capable of controlling up to 256 devices on a single circuit, X10 sent messages to devices through a property's existing Alternating current (AC) electrical wiring. At the time, this was an efficient way to send basic signals through large spaces before the widespread adoption of wireless technology in all electronic devices. It later advanced to be controllable from a computer, enabling users to schedule events and run decision-based sequences. The protocol formed the basis for many domestic control installations for several decades and was one of the first protocols to completely cover the home automation spectrum, with power, security and lighting (*X10*, n.d.).

Throughout the 1980s and 1990s, the idea of having robots as companions was solidified in popular culture through science fiction movies. During this same period, advancements in battery technology and the rapid decrease in the size of microprocessors meant that home robots became achievable, and by 1996, the Electrolux Trilobite was released, the world's first robot vacuum (*Invention of Robotic Vacuum Cleaners*, n.d.). This class of device is now a quintessential smart home product and was a turning point for the industry. The robot vacuum cleaner heralded the integration of the first wave of smart home devices that were not hard-wired into the building. Along with this, the launch of the first-ever internet-connected refrigerator from LG, in the year 2000, demonstrated the increasing integration of technology into everyday household

appliances (Mahajan, Nikam, Patil, & Dond, 2017).

This brings us to the current century, where, in the early 2000s, technology began to boom. The home computer had become commonplace and the internet had become more accessible and understandable to the general public. More smart devices, like speakers that speak to you about news headlines and weather or act as an alarm, began to appear on store shelves at affordable price points and home automation became a realistically achievable goal.

Today, smart homes are now a reality, they are not exclusively for the eccentric and wealthy. They provide homeowners comfort, security, energy efficiency and convenience at all times, regardless of whether they are home or not. Through the use of innovative technology, homeowners can turn their homes into state-of-the-art machines that can be controlled and monitored from anywhere in the world.

It is clear to see that with every iteration of the development of increasingly quasi-intelligent smart devices, early adopters were promised a more comfortable lifestyle with menial tasks being delegated to machines, saving more time for leisure. With the beginnings of electric home appliances reducing the time required to dedicate to cooking and cleaning, through to interactive devices capable of holding simple conversations with users, this rapidly growing technology shows no signs of slowing down and its future capabilities are near limitless.

### **2.1.2 Problems with Smart Homes**

With increasingly modernised smart home technology, its uses become more versatile, however, this has resulted in ever-growing challenges that hinder its expansive potential. Consequently, whilst many smart home technologies progress to solving existing problems, this allows for newer and more complex problems to point out existing flaws within their modus operandi.

There are several factors that limit the advancement of the industry and the effectiveness of smart devices as a whole. Most of the problems with modern smart homes can

be grouped into three main categories:

- Interoperability
- Lack of true intelligence, and
- Requirement for Internet access

(Wilson, Hargreaves, & Hauxwell-Baldwin, 2015).

#### **2.1.2.1 Interoperability**

There are many different competing standards in the world of smart homes, including devices that use Zigbee, Z-Wave, Bluetooth and Wireless Fidelity (Wi-Fi) protocols. This means that any given smart home device may not be compatible with another device, even if they are from the same manufacturer. This can be frustrating for consumers who want to create a customised smart home environment that meets their specific needs and preferences when selecting their products and discourages their entry into the market.

There is currently a new protocol, called Matter, in development by the Connectivity Standards Alliance (CSA), in partnership with some of the leading smart home companies, that aims to unify each of these standards into a single, open-source standard and allow legacy devices to be able to communicate with one another. This is a step in the right direction, however, considering historical context, this may result in it becoming an additional competing standard, with lengthy certification processes driving away smaller manufacturers, and, in turn, competition in the market.

Currently, there are many competing smart home platforms and hubs, such as Amazon Alexa, Google Home, Samsung SmartThings and Apple's HomeKit. Even if the consumer does manage to choose devices all running on the same communication protocols, they may still be unable to control them all from a single app or interface. Apple and Google are both notorious for having lengthy and expensive certification processes for

third-party developers to integrate their devices into their platforms, which can be a significant barrier for smaller manufacturers, splintering the market further.

This usually leads to one of two possible outcomes for the consumer. Either consumers end up locked into a single ecosystem, unable to switch to a different platform without replacing all of their devices, which can be a significant financial burden. Or, they end up with a mix of devices from different manufacturers that are unable to communicate with one another each with its own dedicated app, which can be inconvenient and more time-consuming than not having smart devices at all.

#### **2.1.2.2 Lack of True Intelligence**

Despite their name, smart homes are not actually all that intelligent in their automation abilities. Even in some of the most advanced smart homes, the devices are limited only to simple routines and schedules, and basic decision trees relying on primitive data from a limited number of sensors in the home.

With the recent advent of Artificial Intelligence (AI) and Machine Learning (ML) technologies, it is theoretically possible to create systems that can learn from user behaviour and adapt to their preferences. However, most smart home devices are not yet capable of this level of predictive control and instead rely on the user to program them to perform specific tasks at specific times.

#### **2.1.2.3 Requirement for Internet Access**

Finally, the requirement for constant internet access for smart home devices may be one of the most significant inhibitors of the expansion of the smart home industry. Many smart devices sold today often needlessly process all of their commands remotely, meaning that they require constant internet access to communicate with their respective cloud services in order to be able to control the devices, even from within your own home. There are very few options on the market for devices that allow for local

processing of commands. For example, Google, Amazon and Apple’s respective voice assistants, Google Assistant, Alexa and Siri, all require an internet connection to process voice commands, so without one, they are unable to control any devices in the home.

So if your internet connection goes down, you may be left with a house full of smart devices that are suddenly very dumb. Not only will you not be able to access the devices remotely to control them or monitor your house, but you may also lose the ability to control them from within the same network as they cannot connect to their company’s servers.

## 2.2 Review of Existing Literature

### 2.2.1 Computer Vision for Smart Home Automation

In a paper published at the 2019 International Conference on Communication and Signal Processing, Mohammad Hasnain R. et al. set out to “develop a smart [Internet of Things] (IoT) based light control system” using computer vision (CV) and artificial intelligence (Hasnain, S, P, & G, 2019). Their objective was to reduce the wastage of electricity due to the “negligence and forgetfulness” of people using the environment (Hasnain et al., 2019).

The authors correctly identify one of the biggest issues with presence detection in common home automation deployments. Most setups use an array of infrared (IR) sensors to detect movement in a room but this poses a few challenges. Firstly, any movement in view of the sensors will trigger an action that is intended only for human presence. So if a book falls off a shelf or an animal runs past and interrupts the IR beams, then the sensor will report back a false positive result. Additionally, another limitation of this method, which was not outlined in this paper, is that if a person remains stationary in the room, for example, whilst sitting on a couch, then there will be a false negative result suggesting that there is no human presence in the room.

To solve these downfalls of existing implementations, the authors took a different approach, utilising AI to identify people through a camera in a living space, enabling them to differentiate between objects and people. However, there are a couple of areas that they did not explore that were within reach using the setup that they had implemented and that this thesis intends to expand upon.

A limitation regarding their investigation was that the authors only used the sensors as a toggle for light emitting diodes (LEDs) with the aim of reducing unnecessary energy use. Furthermore, to identify people in the camera, they used You Only Look Once (YOLO), a computer vision AI model for used object detection, classification, and segmentation, which also has a built-in “skeletonisation” feature, allowing you to track a person’s posture and make more intelligent decisions based on a person’s movements. These two ideas present an opportunity to fill a potential gap in the market making more intelligent decisions in the home using existing technology and real-world implementations. Consequently, this thesis aims to implement a system capable of executing more complex tasks such as controlling other smart devices in the house like a television or blinds with more potential states than just off and on.

Aside from this paper, there has been extensive research into the use of computer vision in smart home deployments. Nonetheless, very few of these discuss the integration of computer vision with smart home devices, and even fewer attempt to implement this integration. There is extensive existing research that focuses on computer vision in homes, but it is mostly regarding intrusion detection and security (Cucchiara, Grana, Prati, & Vezzani, 2005; García, Meana-Llorián, G-Bustelo, Lovelle, & Garcia-Fernandez, 2017; Nandhini, Rabik, Kumar, & Brahma, 2020; Patchava, Kandala, & Babu, 2015; Sefat, Khan, & Shahjahan, 2014; Zhang, Yi, & Saniie, 2019).

V. Patchava et al. implemented a front-end for a smart home with the ability to toggle devices and stream the video feed from a camera with computer vision capabilities, while C. González García et al. implemented a system that could detect the presence of people and measured the rate of true positives, false positives, true negatives and false negatives (Patchava et al., 2015; García et al., 2017). Others implemented pres-

ence detection with CV and were able to identify intruders and any weapon that they were carrying (Cucchiara et al., 2005; Nandhini et al., 2020). These papers discuss facial recognition software and the ability to alert homeowners of potential intruders. Moreover, the authors discuss integration with custom smart home dashboards and IoT platforms, but do not discuss the integration of these systems with smart home devices and automation, which is the primary focus of this thesis.

### **2.2.2 Human Action Recognition for Smart Home Automation**

Futher research exists into potential methods of person tracking, human action recognition and behaviour analysis using cameras in the home (Chaaraoui et al., 2013). Here, there is more of a focus on using this technology to support the elderly and the disabled, through ambient assisted living (AAL) to “improv[e] their quality of life and [maintain] their independence” (Chaaraoui et al., 2013).

A. Chaaraoui et al. discuss their deployment with both red, green and blue (RGB) cameras collecting information and recognising “key poses” and hand gestures by creating silhouettes of the person, and also by using red, green, blue and depth (RGB-D) cameras using the depth information to more accurately track movements. They were able to achieve this using a Microsoft Kinect camera and depth sensor “with low-cost and real-time performance”, which is the same sensor that was used for the development and deployment of the custom home automation system in this thesis, as will be covered in more detail in Chapter 3. This existing deployment shows promising signs that the proposed setup should work smoothly in a home environment. While they were able to get these systems working in real-time, recognising primitive human actions such as standing, walking, sitting and falling, there was no mention in the paper of actual integration with home devices.

In 2000, J. Krumm et al. followed very similar processes, using background subtraction to create silhouettes of people and overlayinbgv depth data over the top of these cutouts through stereo images (Krumm et al., 2000). However, due to their early adoption of the technology they were somewhat limited by the computational power of computers

during this period and had to run far more complicated and unwieldy setups that would not be suitable for a real home deployment today. This meant that their trackers ran at only 3.5 Hertz (Hz), even with the computational load shared between three computers, and often had trouble tracking more than three people at a time or people wearing similarly coloured clothing.

The authors were able to integrate their system with some home devices with custom programs such as a wireless mouse that could be carried to any table in the room and the clicks and movements of the mouse would be redirected to the nearest computer display. Another program “automatically start[ed] and stop[ped] a [Videocassette Recorder] (VCR) or [Digital Video Disc] (DVD) movie when a person sits on or stands up from a couch.” The film would also be “automatically rerouted to different displays in the room, depending on where the person [sat]” (Krumm et al., 2000).

These are similar objectives to what this thesis partially aims to achieve, however, the technology has advanced significantly since then and the proposed setup should be able to run on a single computer with a single camera and depth sensor, making it more accessible to the general public.

### **2.2.3 Gesture Recognition for Smart Home Automation**

In 2017, S. Desai and A. Desai proposed an algorithm for gesture recognition for home automation, also using the Microsoft Kinect (Desai & Desai, 2017). Their proposal was specific to only recognising hand gestures, not full-body tracking. This involved segmenting the hand from the image at a specified depth range of 250-650mm, removing the noise and background from the image and converting the RGB into a solid binary black and white image, as shown in Figure 2.1. The finger positions were then compared against a set of predefined gestures and then classified. The classification then defines the action that the system should take, sending a signal to an Arduino board to control different home appliances such as a television, charger or fan.

Using this technique, they were able to achieve an 88% accuracy rate in recognising the

Figure 2.1: Extracted Hand Gestures from Microsoft Kinect (Desai &amp; Desai, 2017)



gestures, a promising result for an early-stage development research project. However, this only worked when users were giving gestures within 65cm from the camera, and only within a 40cm range. This is not practical for day-to-day use as this would either require dozens of cameras per room, spread out every 1.3m, or it would require users to get up from wherever they are in the room and move to the camera in order to control their devices. In this case, controlling their devices manually may be equally as convenient. Due to these limitations, tracking gestures, only via hand movements, does not appear to be a viable solution for smart home automation, and full body tracking may be required to make this technology more practical for everyday use.

This paper, among others, refers also to computer vision technology recognising gestures through wearable devices (Krumm et al., 2000; Starner, Auxier, Ashbrook, & Gandy, 2000; Volety & Geethanjali, 2022). Implementing a system to track wearable devices opens up many possibilities for gesture recognition as well as the potential to integrate medical monitoring tools with it. In fact, using the wearable device, T. Starner et al. were able to track the user's gestures against control gestures, which measure continuous input from the user, with an accuracy of 95% and user-defined gestures with an accuracy of 97% (Starner et al., 2000).

Using a device that interacts directly with the user's body, will significantly increase the accuracy of gesture recognition, as the device is able to measure its own movements, and hence the user, as opposed to a potentially distant camera that attempts to infer

the user's movements based on depth data. However, this creates a significant hurdle as the user must be wearing the device at all times in order to control their devices, which significantly impedes the user experience and increases the inconvenience that this technology is supposed to alleviate.

#### **2.2.4 Behaviour Monitoring for Elderly Care**

In a 2017 paper, “A Behaviour Monitoring System (BMS) for Ambient Assisted Living,” S. Eisa, and A. Moreira constructed an environment that could monitor the behaviour of elderly people in their homes (Eisa & Moreira, 2017). The system used passive infrared (PIR) sensors to monitor the movement of the elderly person between rooms and detect any deviations from their normal routine. The theory behind their practice was that “[a]n elderly person remains in good health as long as he or she can carry on his/her daily activities as usual with no significant deviations from the normal daily routine.” This meant that if the system could identify any irregularities, then it could alert a carer or family member to check in on the user to ensure they are okay and potentially catch any health issues early.

This was implemented by treating the user's location in the house (by room) as a state machine, with the user able to transition states by moving between rooms, as shown in Figure 2.2. From this, a transition matrix is created for a given time of the day, which represents the probabilities of the user either staying in the same room or moving from one room to another. Training an AI model to recognise the user's regular behaviour, the system could then recognise any deviations from this behaviour. If the behaviour was deemed abnormal, it would then be classified to determine possible causes of the deviation, such as the user oversleeping if they are in their bedroom for too long, Table 2.1 shows the abnormal behaviours that the authors listed.

Figure 2.2: Room-to-room State Transition Model (Eisa &amp; Moreira, 2017).

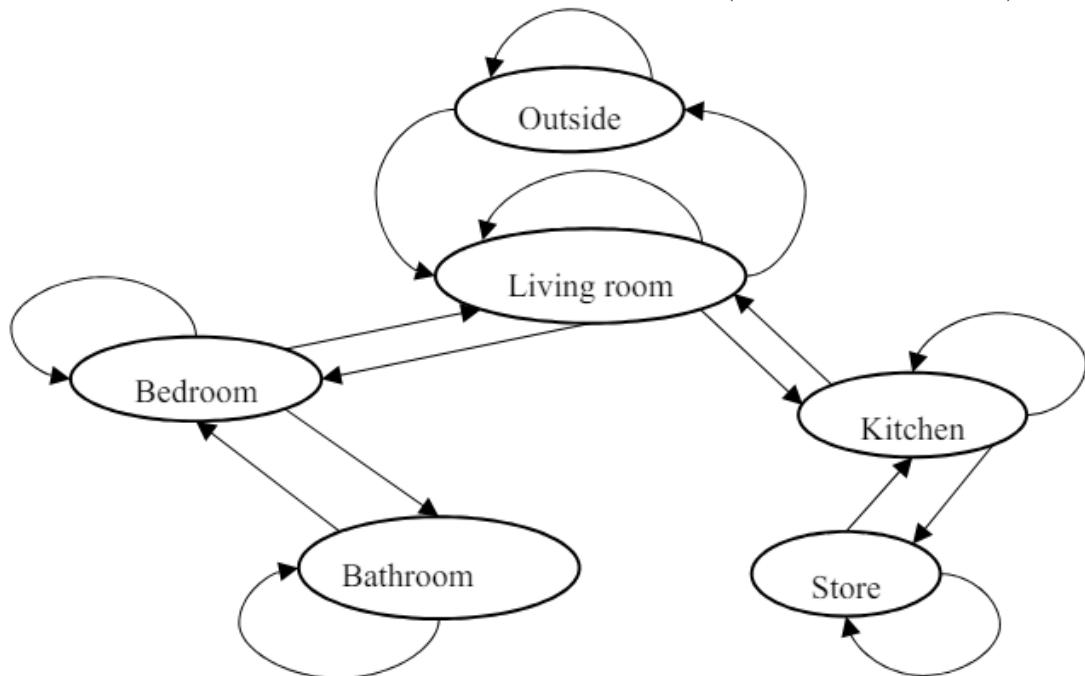


Table 2.1: Abnormal Behaviours—Descriptions and related health-declines (Eisa &amp; Moreira, 2017).

Behaviour	Description	Health-Declines
OverSleeping	An extended stay at bedroom; longer than usual.	Mobility problems, strokes.
LessSleeping	Detected motion at one of the rooms, not a bedroom, during the usual sleeping time.	Having sleepless time due to anxiety, depression or may be developing Alzheimer's diseases.
NotBackHome	Monitored person stayed outside longer than usual.	Having trouble coming back home or get lost or wandering outdoors.
Dead	No movement and long stay at one of the rooms, not bedroom nor outside;	Death.

The model they implemented also self-updated once every week so that it could adapt to the user’s “behavioural changes that are not necessarily abnormal behaviours” over time and account for seasonal changes in behaviour (Eisa & Moreira, 2017). This meant that they were able to implement a behaviour monitoring system that learned to adapt over time and consistently monitor the user’s behaviour for any irregularities

with primitive PIR sensors. If this was achievable with such simple sensors, then the possibilities with more advanced depth sensors and computer vision should allow for more advanced behaviour monitoring and prediction.

There is also some overlap in this field with the computer vision and gesture recognition technologies that have been discussed. M. Oudah et al. implemented a hand gesture recognition system, similar to that of S. Desai and A. Desai, but with the aim of assisting the elderly and disabled in their daily lives by allowing them to use gestures to get support if communication is a problem for them (Oudah, Al-Naji, & Chahl, 2021). With working deployments of both of these technologies in the AAL field, this shows promising signs for potential in the home automation field.

### **2.2.5 Voice Control and Multi-Modal Smart Home Automation**

Using voice commands for controlling devices in a home is an already well explored area of research, with many existing implementations and products on the market for over a decade. Still, in recent years there have been many developments from individual researchers into creating custom open-source voice control systems for smart phones to control devices in the home. Papers from P. J. Rani et al. and J. R. Aluru et al. both discuss the implementation of mobile apps that take voice commands from the user and send them to a miniature computer, an Arduino board or Raspberry Pi respectively, to control devices in the home (Aluru, Kadapa, & Kumar, 2021; Rani, Bakthakumar, Kumaar, Kumaar, & Kumar, 2017). Their implementations are both very versatile and could allow for the user to control a wide range of devices in the home, simplifying home device control. However, both of these systems are limited to only voice commands and do not integrate with any other systems, meaning adding devices requires manual programming and setup, making it far less user-friendly and accessible to the general public. They also are reliant on the use of a smart phone application to record the voice commands which is not always practical for the user to have on them at all times. An improved implementation may have static microphones in the home that are always listening for commands, similar to the Amazon Echo or

Google Home devices, so that users are not burdened by the extra step of having to retrieve their phone and open an app, at which point they could have already manually controlled the device. Additionally, they both require an active internet connection, which is a significant limitation for smart home devices, as discussed earlier.

In 2020, E. Dutt et al. explored the prospect of multi-modal inputs for controlling a smart home (Dutt, Maduri, Gupta, Rathour, & Kushagra, 2020). They proposed a system that could control a smart home using both voice and gesture recognition, with the aim of making the system completely hands free, with no physical touch required, making it more accessible to people with physical disabilities. Their system employed IR sensors to detect a gesture or obstruction of the sensor which triggered the master switch to allow voice commands to be processed. Once activated, the user can then use 15 different voice commands (as specified by the users in advance) to send signals to devices to control them. There are several issues with this approach, including the limited functionality available and the primitive gesture detection methods. Only allowing 15 voice commands is a significant limitation as this restricts the user to only a handful of actions that they can perform with their devices. Furthermore, the gesture detection system is incredibly primitive, only detecting an IR sensor obstruction, which essentially does not detect unique gestures at all but acts as a replacement for a wake word. Despite this, the idea to combine multiple inputs to control a smart home is promising and could be expanded upon to create a more intelligent system, which is one of the goals this thesis aims to achieve.

N. Ganji et al., however, do integrate more complex gesture recognition into their system, using an RGB camera to detect hand gestures, with four pre-defined gestures that determine which device to control (Ganji, Gandreti, & Krishnaiah, 2022). This allowed users to use more complex gestures to control their devices, expanding the possibilities for more robust control of devices in the home. This, in conjunction with voice commands, allows for a more natural interaction with the smart home. While this is a significant improvement over the previous implementations, it still has limitations in the number of gestures that can be recognised and the lack of integration with other systems. The gestures and voice commands are pre-defined so the user cannot customise

them to their needs without significant programming knowledge. The system also requires an internet connection at all times to process the voice and gesture commands, prohibiting the system from functioning when the internet connection is lost.

### 2.3 Gaps in Literature

After a thorough review of existing literature in the space that is being explored by this thesis, it is clear that there has been extensive research into many of the facets of the proposed system in consideration. Computer vision is a widely researched field and has been applied to many different areas of smart home technology, from security to gesture recognition. Human action recognition and gesture recognition have both also been studied for decades, and more recently there has been a significant amount of research into voice controls and multi-modal input recognition.

Despite this, there is a significant gap in the literature when it comes to the actual integration of these technologies with smart home devices and automation. Most of the literature either discusses the technology in isolation or discusses implementation into smart homes on a theoretical level. Other papers discuss solving some of the singular problems with smart homes individually, but none cover the full scope of the problems that are faced by the industry. For example, there is vast research that discusses the use of computer vision for intrusion detection and security in a smart home, but the possibility of using this technology for automating tasks and controlling devices in the home is rarely considered. Almost none of the papers have actually implemented the technology into a real-world smart home environment and tested the results. This presents a significant opportunity for this thesis to contribute to the field by providing a real-world implementation of these technologies and evaluating their performance and user experience. Each of the critiques made in the literature review will be goals to address in this thesis, with the aim of creating a more intelligent and autonomous smart home environment that is more accessible to the general public.

## 2.4 Problem Statement

Smart homes today promise enhanced convenience, safety, and security through the connectivity of more and more devices and the automation of regular tasks. However, several critical challenges prove to be a hindrance to their effectiveness as well as widespread adoption and, in turn, the growth of the industry.

Interoperability issues between devices and existing platforms, the lack of genuine intelligence in these automated systems, and the constant reliance on internet access pose significant obstacles to realising the true potential of smart home technology.

In an attempt to combat these challenges, this thesis aims to develop an entirely custom and customisable, intelligent environment where a user's movements will be able to accurately predict how to control devices within the home, through gesture recognition and behaviour prediction. Utilising local processing and sensor data, the system will enhance convenience, safety, and security without relying on internet connectivity or support from third-party companies to allow interoperability.

Moreover, with several gaps in the existing research and literature in this area, this thesis aims to be able to integrate computer vision techniques and depth sensor technology into a smart home environment, utilising this data to assist in automation, not just home monitoring and security. This will provide a real-world implementation of these technologies and a live demonstration of the home's capabilities, as opposed to the theoretical implementations that are currently discussed in the recent literature. Of those papers with existing implementations, there will be significant expansion on the capabilities of the system, with the aim of controlling more devices and giving the user the ability to easily customise the system to their needs.

Overall, this research endeavours to contribute to the evolution of smart homes by introducing innovative approaches to enhance their intelligence and autonomy, ultimately improving the quality of life for users.

## 2.5 Aims and Outcomes

The research objectives of this thesis include designing and implementing artificial intelligence and computer vision for movement prediction and gesture control, integrating these with concurrent voice controls and into existing smart home infrastructure, finally evaluating the performance and user experience of the environment.

With entirely locally processed data, the system will be able to control devices within the home without the need for internet connectivity, enhancing convenience, safety, and security for users. Computer vision techniques will be used to identify users in the home and track their movements and gestures to identify devices, with voice commands being used to specify an action for the device to take.

The anticipated outcomes of this research include advancements in smart home technology, improved user experiences, and insights into addressing the broader challenges in the field as outlined in the literature review. While the research focuses on addressing specific aspects of smart home functionality, certain limitations, such as device compatibility and privacy concerns, will be acknowledged and considered throughout the study.

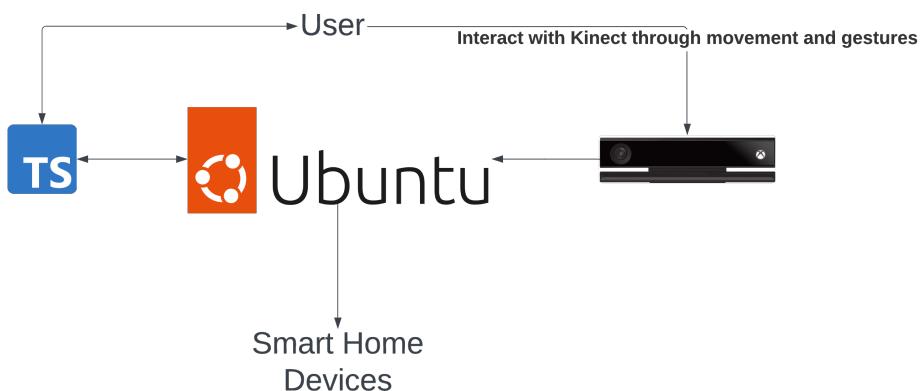
## Chapter 3

# Methodology

### 3.1 Proposed System

To solve the issues outlined in Chapter 2, in Thesis A, the following system was proposed. Using a Microsoft Kinect v2 depth sensor, user's gestures can be detected and processed on an Ubuntu machine. Utilising Robot Operating System (ROS) 2, once the gestures are processed, signals can be sent over the ROS network to control devices throughout the home. Finally, a TypeScript (TS) frontend would be developed to enable users to customise their home to make the system adaptable to user needs. This led to the initial simplified technology stack as shown in Figure 3.1.

Figure 3.1: Simplified Technology Stack



### **3.1.1 Technology Stack**

#### **3.1.1.1 Thesis A**

Throughout the development of this system, the technology stack gradually changed over time in order to accommodate growing goals and to overcome challenges. In Thesis A, there were many challenges in getting the Kinect camera to communicate with ROS 2. A very specific environment setup is required in order for these systems to function together reliably as shown in Figure 3.2.

#### **3.1.1.2 Thesis B**

Following this, during the complete development of the system in Thesis B, the technology stack was modified slightly in order to accomodate a more robust and customisable system for managing smart devices in the home as in Figure 3.3. Using the open source home automation platform Home Assistant allows for far more versatility in automation abilities, choices for devices and provides it's own user interface, eliminating the need for a separate custom TS frontend.

#### **3.1.1.3 Thesis C**

During Thesis C, one of the stretch goals that was set, implementing voice commands and multi-modal device control, was developed. This, again, made some minor changes to the technology stack, incorporating a new TS frontend for users to define their own custom voice commands, and a Vosk speech recognition model integrated as a node on the ROS network. The final technology stack for this thesis as it is currently implemented is as shown in Figure 3.4

Figure 3.2: Thesis A Proposed Technology Stack

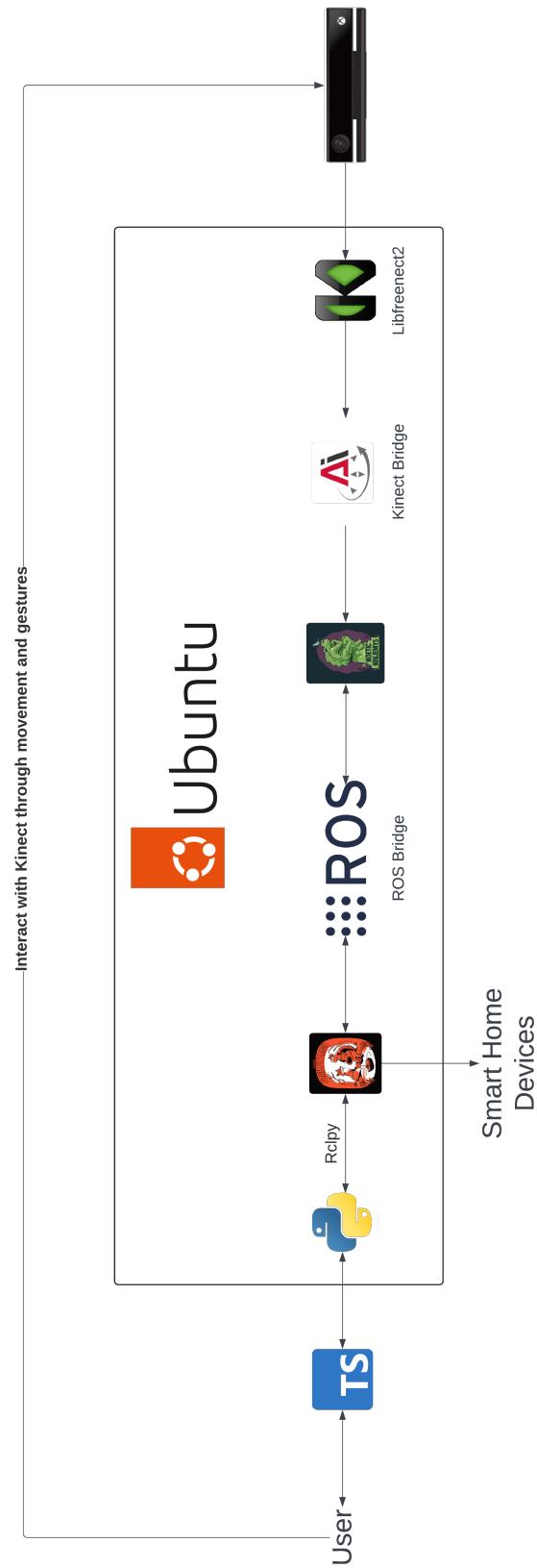


Figure 3.3: Thesis B Implemented Technology Stack

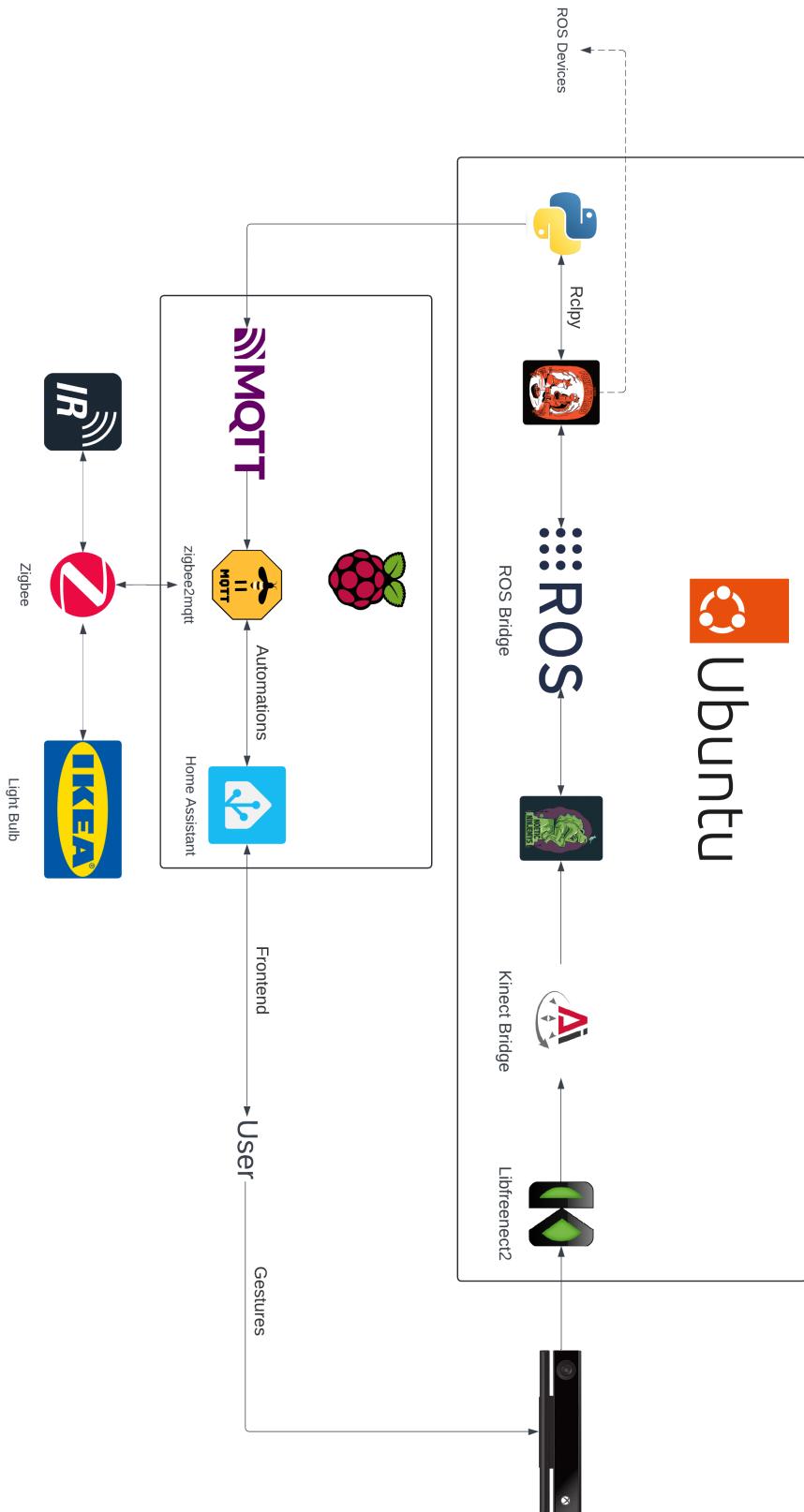
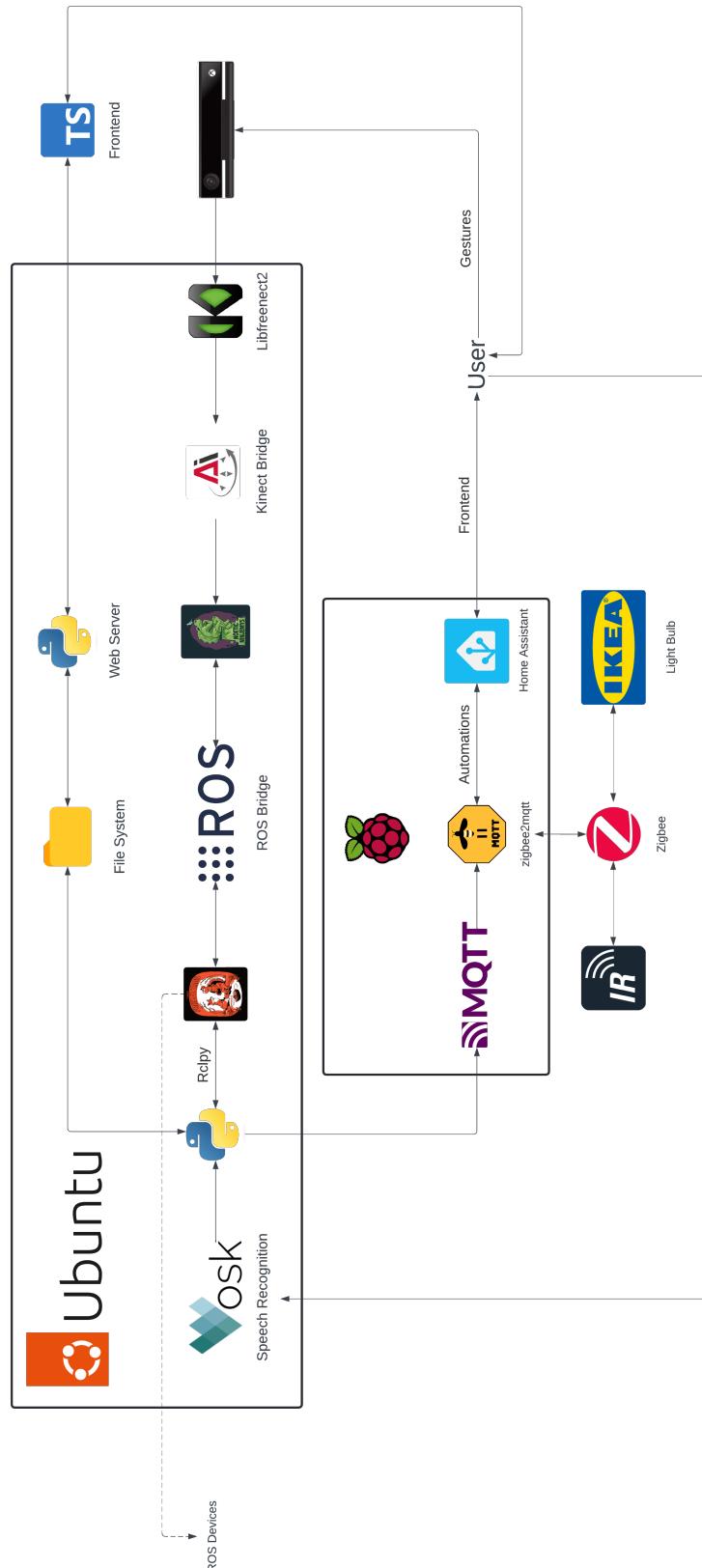


Figure 3.4: Thesis C Final Technology Stack



## 3.2 Hardware Components

The hardware components used in the development of this system include the following:

- Dell XPS15 9510 Laptop
- Microsoft Kinect v2
- Raspberry Pi 3
- Sonoff Zigbee 3.0 USB Dongle Plus
- Ikea Trådfri Zigbee Bulbs
- Zigbee IR Emitter

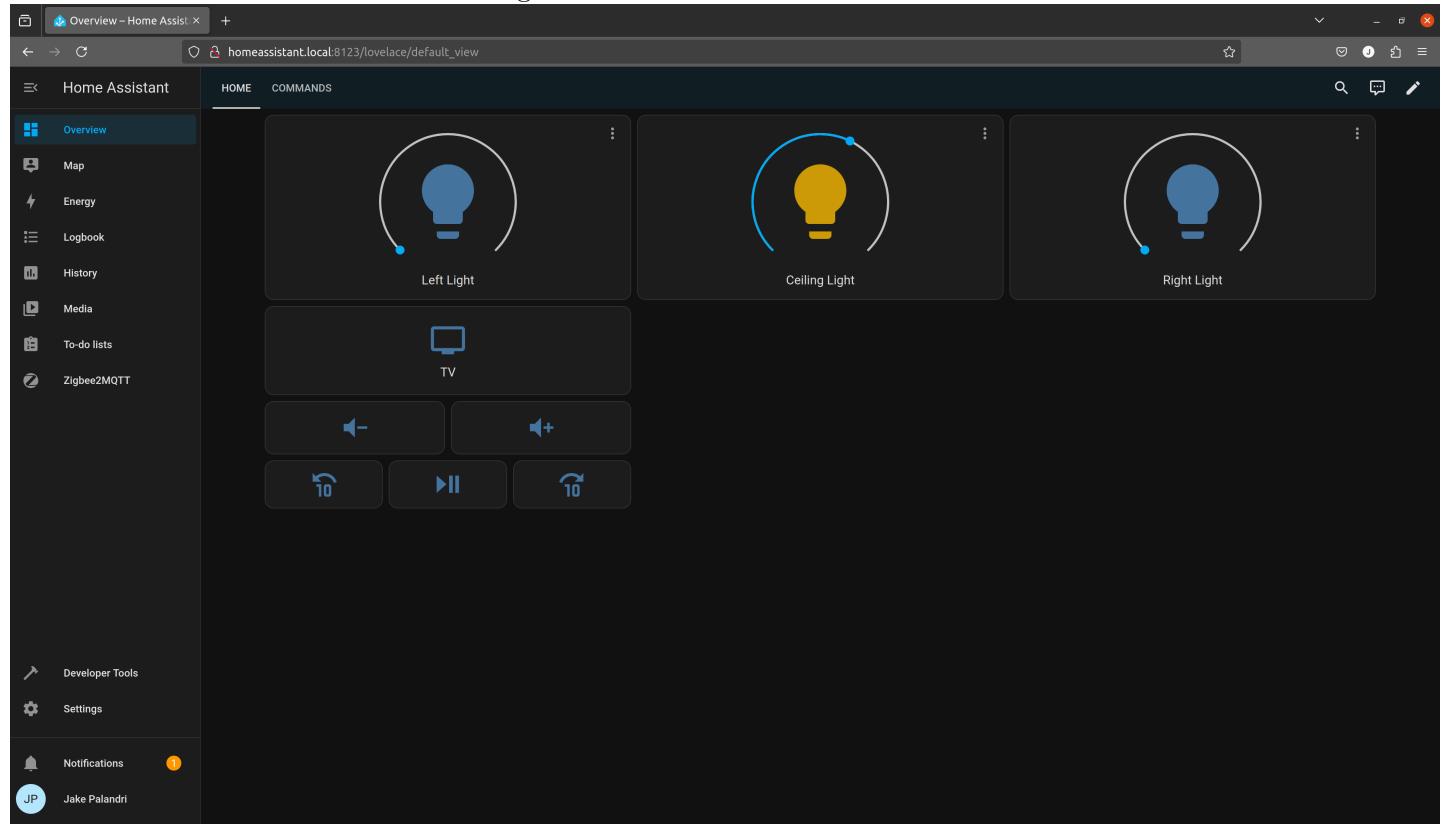
The selection of the Dell XSP15 Laptop was due to its immediate availability as it was already owned. This is a powerful laptop with a discrete graphics card and was always sure to be capable of running the software required. However, a laptop of this calibre is not completely necessary and most modern hardware should be sufficient.

The Microsoft Kinect v2 was also selected due to it being readily available. This, however, is a necessary piece of equipment for this system as it is currently designed. In theory, other depth cameras could be suitable for alternative deployments but this would require significant work to modify the existing environment setup and potentially large portions of the codebase.

The Raspberry Pi 3, was selected for multiple reasons. Home Assistant is relatively lightweight so the choice of hardware was not limited by processing power, therefore a miniature computer such as a Raspberry Pi was appealing due to its low cost, small footprint and ease of use with Home Assistant. The Pi 3 specifically was chosen as it was available on-hand, though other single-board computers capable of running Home Assistant would also be fitting choices.

The Sonoff Zigbee dongle, which is attached to the Raspberry Pi, along with the Ikea light bulbs and IR emitter were primarily chosen for their low price points. Zigbee is

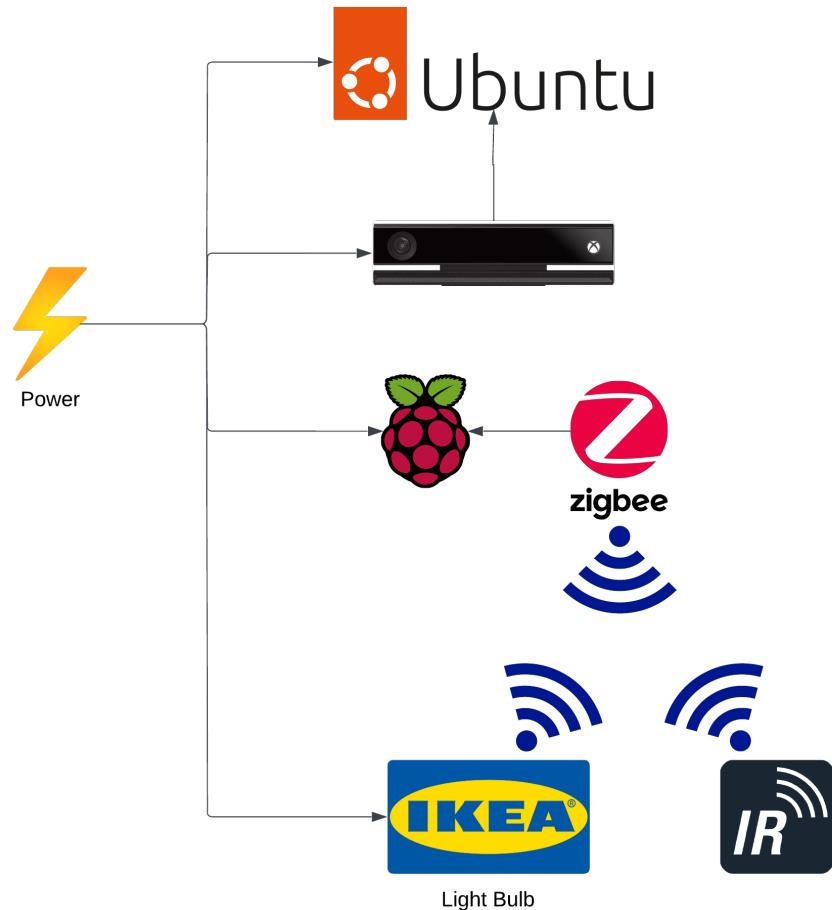
Figure 3.5: Home Assistant Interface



an open source communication standards with lower fees for accreditation than other competing standards such as Z-Wave, therefore their products tend to be cheaper than their Z-Wave counterparts. The dongle is required in order to communicate between Home Assistant and the smart devices. The light bulbs and IR emitter were purely chosen as two examples to demonstrate the functionality of the system. Any competing standard devices or different device types entirely would also be valid choices for the system depending on user needs.

The components should be set up as shown in Figure 3.6.

Figure 3.6: Physical Setup Diagram



### 3.3 Software Components

The key software components used in the development of this system include the following:

- ROS
- ROS Client Library for Python (rclpy)
- Home Assistant

- Messaging Queuing Telemetry Transport (MQTT)
- Vosk Speech Recognition
- YOLO
- Python Web server
- TypeScript frontend

The proposal to use ROS to control devices in the smart home was made very early in the research stages. This was because many of the robots and devices existing in the lab already operated on ROS. Building the smart home on the same system would allow for further research and development into getting these projects working cohesively with one-another, as there are some operations that can be performed by robots in the lab, such as pointing to and retrieving an item, that could be useful in a smart home. In its current state, the commands to the smart home are accessible by other devices on the ROS network, however there is no integration with the robotics work being done, but this does allow for future expansion of the project.

Rclpy is the ROS library for Python and was one of only two options for development. The alternative was a C++ library, however, due to the nature of the project, utilising AI and the lack of C++ skills upon embarkment of the project, rclpy was chosen to write the code for custom ROS nodes. This is where the gesture recognition software operates as well as the speech to text model and voice command recognition nodes.

Home Assistant was chosen because, as an open source home automation platform, it allows users to customise their environment to their liking to a much higher degree than other mainstream platforms such as Apple HomeKit or Google Home. This allowed a simple integration with the ROS components of our system as there was much more choice for methods to communicate between the platforms, nothing was locked down.

MQTT is a lightweight publisher/subscriber machine to machine messaging service. This, together with zigbee2mqtt, allows the ROS nodes to send an MQTT message over the network and control Zigbee devices in Home Assistant directly. This was

selected due to its lightweight messaging service, cheaper products and simplistic implementation. It is worth noting that using MQTT does not limit users to using Zigbee devices as the MQTT messages can be used to trigger automations in Home Assistant. If a user determines that bluetooth, Wi-Fi or Z-Wave devices are more suitable for their requirements, then they are still able with no extra work required.

Vosk's speech-to-text (STT) models were selected for two main reasons. Firstly, these are the same models used by other robots in the Human Robot Interaction lab where this system was deployed. Maintaining a smaller number of tools makes the environment far easier to manage as functionality expands. Additionally, Vosk supports streaming audio directly from a microphone to the model, whereas many other local STT models only support processing entire audio files. Other models such as OpenAI's local Whisper models were tested in development, however the inability to stream audio made reliably interpreting the audio challenging. YOLO was also selected as this was already in use by other systems in the lab and it provides a skeletonisation feature which allows the extraction of the coordinates of key points on a users body to track their poses and gestures.

Python and TypeScript were selected for their simplicity in development with TS chosen over JavaScript (JS) for it's strict typing and as an extension challenge for the developers, previously unexperienced with TS. Other languages for the backend web server and the frontend web app are also suitable choices.

### **3.4 Component Interactions**

The process of enabling communication between ROS 2 and the Kinect camera took many iterations of testing throughout Thesis A, and resulted in a somewhat convoluted setup as shown in Figure 3.7. However, this was necessary for the goals of this thesis. The Institute of Artificial Intelligence at the University of Bremen (IAI) has developed a set of drivers to allow communication between the Kinect and ROS 1, and another user created a fork of this repository to support OpenCV 4.0, an open source CV library,

which was required for this deployment. Since there is no existing Kinect bridge for ROS 2, the data from the Kinect must be processed first by ROS 1 and passed through to ROS 2. Because of this, the operating system required is Ubuntu 20.04, Focal Fossa, as this Linux distribution (distro) supports multiple versions of both ROS 1 and ROS 2, allowing communication between one another. OpenKinect's Libfreenect2 drivers are required in order for IAI's Kinect 2 Bridge to access the data from the sensors on Ubuntu. Finally, Open Robotics, the developers of ROS, provide a bridge between ROS 1 and ROS 2 that allows the translation of the topics, services, and actions between nodes in either direction. These tools, in conjunction with ROS 1 distro Noetic Ninjemys, and ROS 2 distro Foxy Fitzoy, ROS 2 is able to access the data from the Kinect.

Figure 3.7: Kinect to ROS 2 Layers

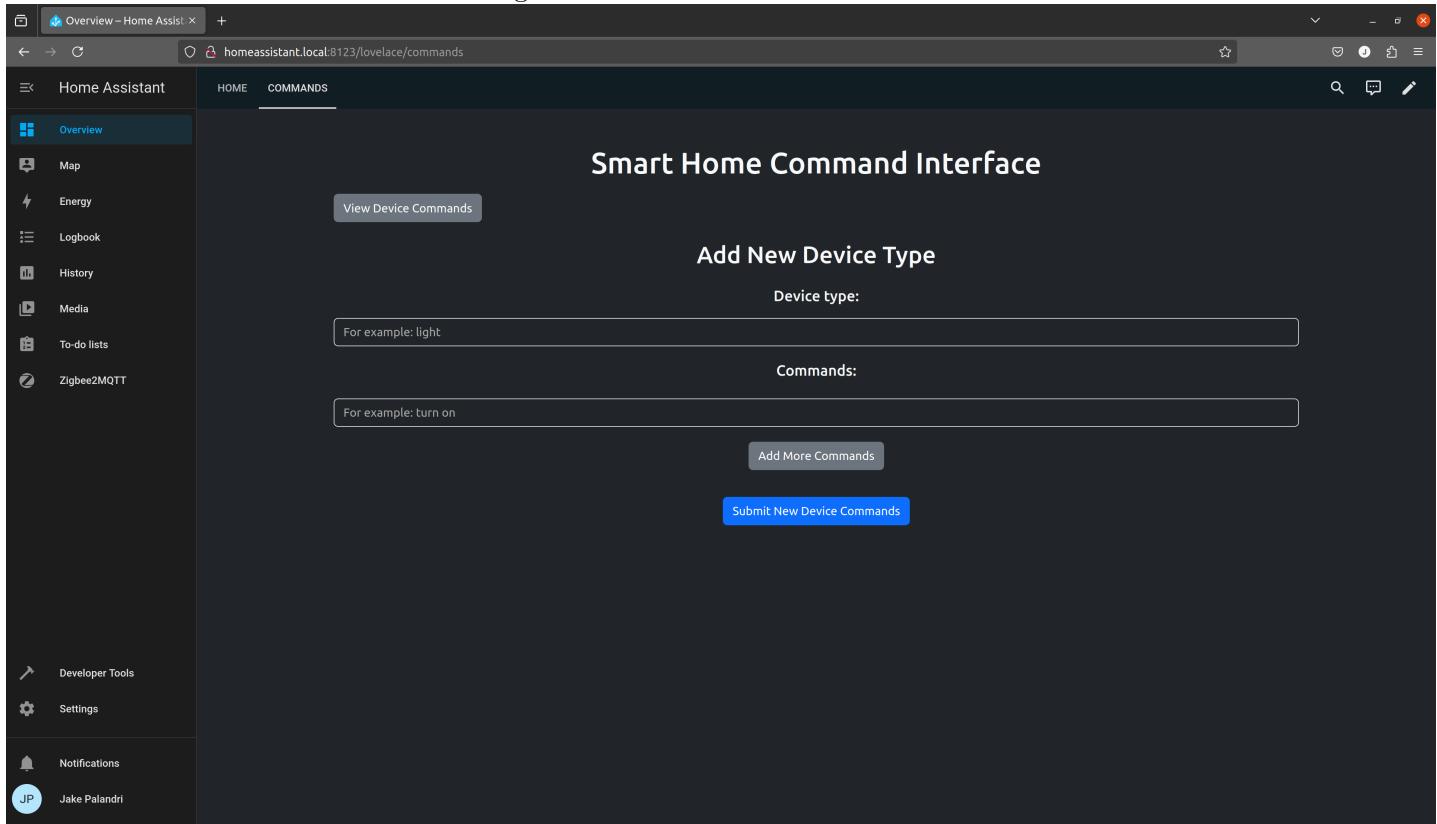


To integrate this with the current system, rclpy is used to process and register the gestures from the users and signals are then sent to the Home Assistant server on the network via MQTT. A Raspberry Pi 3, is used to run the Home Assistant server along with an MQTT broker, and zigbee2mqtt, a service to bridge mqtt signals directly to zigbee smart devices. From there, Home Assistant is able to control the smart devices over the zigbee network through manual control from the user or automations.

Using rclpy, there are two ROS nodes created to process the speech and the user's gestures. The STT node takes audio directly from the computer's built-in microphone and streams it to Vosk to process the speech. The result of the Vosk model is then published to a topic on the ROS network. The second node subscribes to both the STT topic published by the first node, and the topics created by the Kinect sensors through the layers of translation from the camera to ROS 2. This then processes the voice commands at the same time as the gestures and will send the MQTT message to trigger devices.

Finally, there is a web app that allows users to customise which voice commands are available to control their devices. This is available on the network through a local IP address, or within Home Assistant as shown in Figure 3.8. The web app accesses a JavaScript Object Notation (JSON) file on the system that contains all device types available in the home and their associated commands. From here the user can add, modify or remove commands and device types entirely depending on the needs of their home. This modifies the JSON file on the system, which in turn triggers an update in the ROS node to ensure that only valid commands are being checked for.

Figure 3.8: Command Web Interface



### 3.5 Code Logic

Once the Kinect data and STT result has reached the final ROS node, there are primarily two processes running concurrently in order to trigger device commands.

First, with each frame from the Kinect’s camera, the image callback is triggered which determines a users gesture and then stores it for use later by the voice commands process. To determine a user’s gesture, eight points on the users body are extracted using YOLOv8. The hip, shoulder, elbow and wrist of both the left and right sides. Then, using the depth data from the Kinect, the pixel coordinates from the image can be converted to 3-Dimensional world coordinates. In case of any errors in YOLO selecting a point just outside the users silhouette, the depth value is taken as the point closest to the camera within a five pixel radius of the point. Using the wrist and elbow coordinates, a vector can be created to determine which direction a user is pointing. This is done by extending this vector and intersecting it with the bounding box that is the dimensions of the room. Then it can be determined which wall, floor or ceiling is being pointed at. A gesture is only registered if the distance from the users wrist to their hip or their shoulder surpasses a pre-defined threshold. The gesture is then stored in a list as a history with an associated timestamp with a gesture being stored at 250 milliseconds intervals and storing the last 10 seconds of gestures.

At the same time, whenever the STT node sends a new string the speech processing callback is triggered. The output from Vosk provides multiple possible interpretations of the speech in order to minimise misinterpretations. Each of these alternatives is checked to see if the sentence begins with the chosen wake word, set to “home” by default. The wake word is not required to trigger a command but is required for the user to receive feedback for an invalid command. This was done so that if a user forgets to use the wake word then if a valid command is given then it will still trigger a result. However, if regular speech is picked up not intended to control a device, then the invalid command feedback is not triggered.

Following this, each of the alternative interpretations from Vosk are checked against the JSON object containing all the valid device types and commands. These are checked against a regex string in order to allow users to speak in a more natural manner and still trigger the correct device with other superfluous words in the command.

The regex string for matching a command for a single device is:

```
rf".*\b{command}\b.*\bthat\b.*\b{device}\b.*"
```

For the case where the command is “turn on” and the device is a “light,” this gives users the ability to say the simple command “turn on that light”, or speak more naturally and say “could you turn on that light over there for me please.”

The regex string for matching a command for all devices of a single type is:

```
rf".*\b{command}\b.*\b(all|every)\b.*\b{device}s?\b.*"
```

This gives users the option to say the commands “turn on all lights,” or “turn on every light,” or more naturally “can you turn on all of the lights please.”

Once a command is matched, it is added to a list, and once each of the device types and commands is checked against the command with the latest match is selected. This is done so that a user can correct themselves mid-sentence and still trigger their intended command. For example, a user might say “can you turn on that light, I mean that fan.” This would trigger the automation to turn on the fan. Once a command is recognised, the timestamp at which the word “that” is said is compared to the history of gestures and the gesture with the closest timestamp to the keyword is selected.

If no match is found for a given alternative interpretation from Vosk, the process loops. If no match is found at all, and a wake word was supplied, then an audio cue is triggered to let the user know that their command was not recognised and to try again. Additionally, if the user gives a valid command but does not gesture to point to their device, another audio message is relayed to the user to ask them to gesture when giving commands. Finally, if a match is found, the command message is sent over MQTT and captured on Home Assistant to trigger custom automations set up by the user.

# Chapter 4

## Evaluation

### 4.1 Solved Problems

The three main problems with existing smart homes, as outlined in Chapter 2, are a constant requirement for internet access, interoperability between smart devices and systems, and the lack of true intelligence. The goal of this thesis was to attempt to tackle some of these challenges. By utilising local AI models for pose estimation and speech recognition, as well as using Home Assistant as the platform to control the devices in the home, this solves the issue as all computation is done on the local network, eliminating the need for an internet connection. The system is still accessible from outside the home over the internet, but in the case that the network drops its connection, users within the house do not lose their ability to control their devices. Using Home Assistant also solves the interoperability issue. As an open source home automation platform, devices from previously incompatible communication standards can all be brought together and controlled under one system, whether they are compatible with Google, Apple or Amazon's platforms or if they use Bluetooth, Wi-Fi, Z-Wave, Zigbee or another proprietary standard. Tackling the issue of true intelligence to the extent that was described in the literature review, by utilising AI to predict user behaviour and patterns, was outside the scope of this thesis. However, using AI for pose estimation and speech

recognition does provide an added layer of quasi-intelligence to improve the landscape of the smart home industry.

One of the key rules that should be considered when developing any smart home or smart device is that functionality should always be added to, and never replaced. A smart light that is no-longer controllable from the wall switch due to it's new smart features is not an improvement on the old system. Technology is not perfect and fails every day. If the smart features fail to work then the house is left completely inoperable. In the deployment of this thesis, all devices had functionality added, and did not interrupt existing operation. The light bulbs were still controllable by a physical switch, and in addition they had a new remote to control them, they could be controlled through the home assistant app or website and finally by voice commands and gestures. The IR emitter allowed the television (TV) to be controlled by voice and gestures but the original TV remote was still available for regular use.

There were also some problems that had not been addressed in the existing literature that were solved in this thesis. Hasnain et al. used YOLO for person detection but did not utilise the pose estimation capabilities, and additionally the automations they implemented were primitive, only allowing an on and off state of an LED (Hasnain et al., 2019). This thesis uses pose estimation for more advanced gesture detection and is capable of controlling devices with many more states than just on and off. The TV can be controlled to turn it on and off, turn the volume up and down, and fast-forward and rewind the content being viewed. The lights could also be controlled to toggle their power and change their brightness up and down. Krumm et al. had employed a similar system, however were limited to 3.5Hz refresh rate of their cameras, which in this thesis has been increased to 30Hz (Krumm et al., 2000).

Aluru et al. and Rani et al. both implemented voice control systems for the home through mobile apps. However, as outlined in the literature review, this limits the convenience gained by requiring users to have their phone with them at all times to control devices as there was no stationary microphones throughout the home (Aluru et al., 2021; Rani et al., 2017). Utilising the microphone of the computer used to run

the software required, or another microphone plugged into it, means that the system is always ready and listening for commands. Dutt et al. faced the challenge of only having a total 15 possible voice commands that occupants of the home could use to control their devices (Dutt et al., 2020). This is not very user friendly and does not allow for a modular system with room to expand as more devices are added to the home. With the web app implemented in this thesis, users are able to add their own commands and devices as they see fit. Finally, Ganji et al. had the most advanced system implemented of those covered in Chapter 2, with more advanced gesture detection and device control. However, due to the complicated nature of their implementation, it is still not a very accessible system for end users. Devices and gestures in theory could be added to the system, but this would require significant technical knowledge and reprogramming the software to accommodate any new devices. Utilising Home Assistant as the platform for the smart devices allows users to simply add new devices to the system themselves and modify which gesture and voice command combos control which devices. Additionally, with the web app that was developed, users are able to add to and customise their list of allowed voice commands and devices. This added level of flexibility makes it far more appealing for a real world deployment.

## 4.2 System Evaluation

After development was completed, three users were chosen to test the reliability of the system as shown in Figure 4.1 and Figure 4.2. Each was given a script of 13 commands to read from, and were given the option to add their own words to the commands if they saw fit and attempt other related commands.

- Turn on that light (With no gesture)
- Turn on that light
- Turn off that light
- Turn up that light

- Turn down that light
- Toggle that light
- Turn on all/every light(s)
- Turn on that TV
- Turn off that TV
- Turn up that TV
- Turn down that TV
- Fast forward that TV
- Rewind that TV

Figure 4.1: User 1 Testing



Figure 4.2: User 2 Testing



#### 4.2.1 Quantitative Analysis

Following these tests, the results were as shown in Table 4.1

Result	User 1	User 2	User 3
True Positive	12	13	19
False Positive	0	0	0
False Negative (Speech)	5	7	7
False Negative (Gesture)	0	0	2
Success (%)	71%	65%	68%

Table 4.1: User Testing Results

In all cases, commands given with no gesture were correctly identified and the user was notified. Additionally, on only two occasions was there a false negative result due to a gesture being misidentified. Both of these were with user 3 and likely due to the fact that she was wearing a very loose jumper causing YOLO to select the incorrect location for the elbow as shown in Figure 4.3. It could also possibly have been correctly selected by YOLO, but when the depth data at this coordinate was calculated, the loose fabric caused a smaller depth to be read. There was also never any false positive results, meaning the system will never trigger a device automation without explicit input from the user. Overall, the success rate was an average of 68%, almost entirely due to verbal commands being misinterpreted and commands took approximately two to three seconds from the end of the sentence to the device being triggered.

Figure 4.3: User 3 Testing



OpenAI's Whisper STT model was explored for interpreting voice commands, however it does not support streaming to the model directly, meaning that in order to get it

to work, chunks of audio had to be sent in five second snippets. This causes a few problems, most notably that a command can be cut off mid sentence. To solve this, the audio was sent in rolling 10 second clips every 5 seconds to maintain a full sentence. The drawback with this is that it creates a significant delay in processing the commands making the system significantly slower. Additionally, the timings with the Whisper model were far less consistent, meaning that aligning the timings with the keyword “that” when gesturing became an impossible task as it was non-deterministic. Despite Whisper having much more accurate and consistent interpretations, because of these reasons it was unfeasible for use in this project.

In addition to the three users testing the system under standardised conditions, there were tests conducted under varying conditions with different types of background noise. The three tests were with background talking, with background music with lyrics, and with background music without lyrics.

Result	Background Talking	Background Music (with Lyrics)	Background Music (Instrumental Only)
True Positive	1	1	13
False Positive	0	0	0
False Negative (Speech)	12	15	7
False Negative (Gesture)	0	0	0
Success (%)	8%	6%	65%

Table 4.2: User Testing Results

With any amount of words being spoken in the background, whether it is conversation or lyrics in a song, even at a low volume, the accuracy of the model drops significantly to 8% and 6% respectively. With instrumental only music playing in the background the accuracy of the model drops only marginally to 65%. With an approximate 90% drop in accuracy, this solution is not suitable for use in an environment where there will be constant conversation, unless users are in a controlled environment and can pause discussions in order to trigger their commands. In an environment with gentle music only a 5% drop is recorded in preliminary testing and would be appropriate for deployment.

#### 4.2.2 Qualitative Analysis

Users were surveyed after their testing, and the responses were relatively consistent. All users considered that this system would overall provide an increase in convenience if it had been a little more polished. When it worked, it was responsive and felt like a step towards a more advanced intelligent home environment. However, the accuracy of the STT model was a big enough drawback that most would not want this in their homes yet. With time for technology to improve and a better STT model operating this could be a promising prospect.

In environments with background speech, the system also experienced a significant slow down in processing speed of the STT model. This created large delays in transcription, even once the speaking had stopped, until it had time to catch back up to itself during periods of quiet. This could affect its use in busy homes if commands are not registered in time after long periods of high volume conversation.

## Chapter 5

# Conclusion and Future Work

### 5.1 Future Work

At the conclusion of this thesis, there are several areas that require further research and development to improve the system's accuracy and functionality. The following is a list of potential future work that could be undertaken to improve the system:

Firstly, implementing typeless voice commands would allow users to give the command ‘turn on that’ with a gesture, without specifying a device. This would only be possible for the case when the direction they point has one device with that matching command. This would require the system to store the locations of devices in the room and match the direction of the gesture to the location of the device as the command is given, not on the Home Assistant end.

Next, implementing specific locations with vectors from gestures would allow users to point to a specific location in the room, as opposed to the four walls, ceiling and floor that are currently possible, to give more robust control over devices. This would also require the system to store the locations of devices in the room in order to accurately identify the selected device, and would allow users to have more flexibility in the positioning of their devices.

When the structure for the MQTT commands was chosen to be `{"kinect_pose": "direction_device.command"}`, it was chosen to be as close as possible to the structure of existing Home Assistant commands. In future, it would save users time during setup if the system could automatically trigger a command on the Home Assistant end when a command is received from MQTT. This would mean sending a command such as `light.ceiling.turn_on` and having the system automatically trigger the corresponding action without requiring the user to set up an automation in Home Assistant.

The system currently uses an RGB image for pose estimation, which would not work at night with the lights off. Therefore users are not able to control their devices in the dark. The Kinect has an infrared sensor that could potentially be used to detect poses in the dark, by detecting the heat signature of the user and translating this into an RGB image.

Another quality of life feature that could be implemented, is the ability to have more advanced commands to include variables such as brightness levels or volume. This would allow a user to have more granular control over their devices, with a command such as ‘turn on that light to 50%’.

Finally, the most significant improvement that could be made to the system would be to implement a more advanced Speech To Text engine that supports streaming. There has been a significant amount of research into this area, and it is likely that a more advanced open source STT engine would be able to improve the accuracy of the system significantly.

## 5.2 Conclusion

This thesis has successfully developed a system that uses computer vision techniques to recognise gestures, and speech recognition tools for voice commands, to control any conceivable smart device in a home. Many of the challenges outlined in the literature review with existing smart homes and the current literature were addressed in the development of this system, including interoperability between devices, the requirement

for internet connectivity, and the need for a more intuitive and accessible interface for users without technical expertise.

The results of the evaluation concluded that although the system was incapable of recognising voice commands with a high degree of accuracy, it was able to recognise gestures nearly flawlessly. User feedback indicated that the system was easy to use and intuitive, and with a more accurate STT engine, the system could be a viable alternative to traditional smart home interfaces.

# References

- Aluru, J. R., Kadapa, S. K., & Kumar, G. A. E. S. (2021). Voice control iot home automation using voice assistant raspberry pi. *Annals of the Romanian Society for Cell Biology*, 25, 5819-5825. Retrieved from <https://wwwproxy1.library.unsw.edu.au/login?url=https://www.proquest.com/scholarly-journals/voice-control-iot-home-automation-using-assistant/docview/2564182405/se-2>
- Chaaraoui, A., Padilla-López, J., Ferrandez, J., García-Chamizo, J., Nieto-Hidalgo, M., Romacho-Agud, V., & Flórez-Revuelta, F. (2013). *A vision system for intelligent monitoring of activities of daily living at home* (Vol. 8277). doi: 10.1007/978-3-319-03092-0\_14
- Cucchiara, R., Grana, C., Prati, A., & Vezzani, R. (2005, 4). Computer vision system for in-house video surveillance. *IEE Proceedings: Vision, Image and Signal Processing*, 152, 242-249. doi: 10.1049/IP-VIS:20041215
- Desai, S., & Desai, A. (2017). Human computer interaction through hand gestures for home automation using microsoft kinect. In N. Modi, P. Verma, & B. Trivedi (Eds.), (p. 19-29). Springer Singapore.
- Dutt, E., Maduri, P. K., Gupta, A., Rathour, S., & Kushagra. (2020). Touch-less home automation system with voice and gesture control. In (p. 781-785). doi: 10.1109/ICACCCN51052.2020.9362829
- Eisa, S., & Moreira, A. (2017). A behaviour monitoring system (bms) for ambient assisted living. *Sensors*, 17. Retrieved from <https://www.mdpi.com/1424-8220/17/9/1946> doi: 10.3390/s17091946
- From sci-fi dreams to everyday reality. tracing the evolution of smart home through history.* (2023, 11). Retrieved from <https://www.linkedin.com/pulse/from-sci-fi-dreams-everyday-reality-tracing-evolution-smart-mypmf/>
- Ganji, N., Gandreti, S., & Krishnaiah, T. R. (2022). Home automation using voice and gesture control. In (p. 394-400). doi: 10.1109/ICCES54183.2022.9835832
- García, C. G., Meana-Llorián, D., G-Bustelo, B. C. P., Lovelle, J. M. C., & Garcia-Fernandez, N. (2017, 11). Midgar: Detection of people through computer vision in the internet of things scenarios to improve the security in smart cities, smart towns, and smart homes. *Future Generation Computer Systems*, 76, 301-313. doi: 10.1016/J.FUTURE.2016.12.033
- Hasnain, M. R., S, R., P, M., & G, S. (2019). Smart home automation using computer vision and segmented image processing. In (p. 429-433). doi: 10.1109/ICCSP.2019.8697997

- Hertzmann, P. (2016). The refrigerator revolution.
- Invention of robotic vacuum cleaners.* (n.d.). Retrieved from <http://www.vacuumcleanerhistory.com/vacuum-cleaner-development/history-of-robotic-vacuum-cleaner/>
- Krumm, J., Harris, S., Meyers, B., Brumitt, B., Hale, M., & Shafer, S. (2000). Multi-camera multi-person tracking for easyliving. In (p. 3-10). doi: 10.1109/VS.2000.856852
- Mahajan, M. P., Nikam, R. R., Patil, V. P., & Dond, R. D. (2017, 3). Smart refrigerator using iot. *International Journal of Latest Engineering Research and Applications*.
- Nandhini, M., Rabik, M. M., Kumar, K. S., & Brahma, A. (2020). Iot based smart home security system with face recognition and weapon detection using computer vision. *Regular*.
- Oudah, M., Al-Naji, A., & Chahl, J. (2021, 6). Computer vision for elderly care based on deep learning cnn and svm. *IOP Conference Series: Materials Science and Engineering*, 1105, 012070. doi: 10.1088/1757-899X/1105/1/012070
- Patchava, V., Kandala, H. B., & Babu, P. R. (2015). A smart home automation technique with raspberry pi using iot. In (p. 1-4). doi: 10.1109/SMARTSENS.2015.7873584
- Rani, P. J., Bakthakumar, J., Kumaar, B. P., Kumaar, U. P., & Kumar, S. (2017). Voice controlled home automation system using natural language processing (nlp) and internet of things (iot). In (p. 368-373). doi: 10.1109/ICONSTEM.2017.8261311
- Sefat, M. S., Khan, A. A. M., & Shahjahan, M. (2014). Implementation of vision based intelligent home automation and security system. In (p. 1-6). doi: 10.1109/ICIEV.2014.6850818
- Spicer, D. (2016, 5). *The echo iv home computer: 50 years later*. Retrieved from <https://computerhistory.org/blog/the-echo-iv-home-computer-50-years-later/?key=the-echo-iv-home-computer-50-years-later>
- Starner, T., Auxier, J., Ashbrook, D., & Gandy, M. (2000). The gesture pendant: a self-illuminating, wearable, infrared computer vision system for home automation control and medical monitoring. In (p. 87-94). doi: 10.1109/ISWC.2000.888469
- Stein, J. A. (2011, 7). Domesticity, gender and the 1977 apple ii personal computer. *Design and Culture*, 3, 193-216. doi: 10.2752/175470811X13002771867842
- Us inflation calculator.* (n.d.). Retrieved from <https://www.usinflationcalculator.com/>
- Volety, R., & Geethanjali, P. (2022). Smart home automation using wearable technology. In G. D. Gargiulo & G. R. Naik (Eds.), (p. 259-279). Springer Singapore. Retrieved from [https://doi.org/10.1007/978-981-16-5324-7\\_11](https://doi.org/10.1007/978-981-16-5324-7_11) doi: 10.1007/978-981-16-5324-7\_11
- Wilson, C., Hargreaves, T., & Hauxwell-Baldwin, R. (2015, 2). Smart homes and their users: a systematic analysis and key challenges. *Personal and Ubiquitous Computing*, 19, 463-476. doi: 10.1007/s00779-014-0813-0
- X10. (n.d.). Retrieved from <https://www.x10.com/>
- Zhang, X., Yi, W. J., & Saniie, J. (2019, 5). Home surveillance system using computer vision and convolutional neural network. *IEEE International Conference on Electro Information Technology, 2019-May*, 266-270. doi: 10.1109/EIT.2019.8833773

# Appendix

## A.1 User Manual

### A.1.1 Setup Guide

To set up your machine to run this system, you must first install Ubuntu 20.04 on your machine. During installation, the default settings can be followed, and only a “minimal installation” is required when prompted. There is no need to check either of the boxes “Download updates while installing Ubuntu,” or “Install third-party software for graphics and Wi-Fi hardware and additional media formats.” Upon restart, if you are prompted to update to Ubuntu 22.04, decline the upgrade.

Once Ubuntu is installed, you can then follow the install guide in the file `install.sh` of the GitHub repository (repo) found at <https://github.com/jakepalandri/ROSHome>. These commands in the bash file must not be run as a bash script as it will not succeed. They must be executed line by line. Due to L<sup>A</sup>T<sub>E</sub>X document formatting, the instructions cannot be reliably included in this document, however the bash script can be found here.

This includes all of the instructions on how to install based on the setup of the system used in development, however, depending on the exact machine used, your requirements may differ. For example, CUDA drivers are included in the installation script, however these are only applicable to computers with an Nvidia graphics card. The core steps that apply to all users are:

- Install Git and clone the repo
- Install ROS 1 and configure the environment
- Install graphics drivers for your respective graphics card
- Install Libfreenect2
- Install IAI’s Kinect2 Bridge for OpenCV 4
- Install ROS 2

- Install ROS Bridge
- Install modules required to run the code
- Download an STT model from Vosk
- Restart your computer
- Set up physical hardware according to Figure 3.6
- Set up your selected smart devices within Home Assistant and their corresponding automations

For your Home Assistant installation, you must first set up any required devices in the web app. Once this is completed, you will need to set up automations for each device that you wish to control and for each command that they should respond to. The structure of the MQTT messages that are sent to Home Assistant are

```
{"kinect_pose": "direction_device.command"}
```

For example, the voice command “turn on that light” with a gesture to the ceiling would be sent as

```
{"kinect_pose": "ceiling_light.turn_on"}
```

You must create an automation in Home Assistant that listens for this message on MQTT and triggers the corresponding action.

### A.1.2 Execution Guide

Once your system is set up it is time to run the code, but first there are a few variables that are specific to a user’s needs that could require some modification. Each of these is marked with a comment in the repository as “CUSTOMISABLE”. The elements that can be customised are as follows.

The dimensions of the room are defined as a bounding box in `kinect_pose.py`. These should be modified to match the dimensions of the room, with respect to the position of the Kinect.

There is a function remaining from Thesis B called `send_gesture()` which sends gesture commands to Home Assistant without a voice command. This can be uncommented to enable this functionality.

The threshold distances for a users wrist from their hip and shoulder are defined in `kinect_pose.py`. These may be modified to adjust the sensitivity of the system.

Finally, in `kinect_pose.py`, the variable `wake_word` can be modified to change the wake word that the system listens for. By default, this is set to “home”.

In `MQTTClient.py`, the IP address of the MQTT broker must be set to the IP address of the Home Assistant server. This should be set to a static IP address in your router

settings to ensure the system works reliably, and then the IP address should be set correspondingly in the code.

If you did not install the repository into the home directory, , then the path to the Vosk STT model must be updated in `speech_to_text.py`.

In `web_server.py`, the port that the web server runs on can be modified as required, and this should also be modified in the web app's TS code to match, in `main.ts`. Additionally, in `main.ts`, the IP address of the web server must be set to the IP address of the machine that is running all of the code. This should also be set to a static IP address in your router settings.

Once these modifications have been made, the system can be run by executing the commands as outlined in `ros_instructions.sh`. There are eight separate terminal windows that must be opened to run the system. These will run the following services:

- ROS 1 Core
- Kinect2 Bridge
- ROS Bridge
- Kinect Gesture and Voice Command Recognition
- Speech to Text
- Voice Command Web Server, and
- Voice Command Web App