

401I — Final Year Project

Imperial College London

Voice Recognition RPG

Author: Baron Khan

Supervisor: Dr. Mark Wheelhouse

Contents

1 Abstract	1
2 Acknowledgements	2
3 Introduction	3
3.1 Role-Playing Games	3
3.2 Problem Statement	3
3.3 Motivation	4
3.4 Objectives	4
4 Background	6
4.1 Voice Recognition in Games	6
4.1.1 In Verbis Virtus	6
4.1.2 Skyrim with Kinect	7
4.1.3 Star Trek: Bridge Crew	8
4.1.4 Classic Zork on Alexa	8
4.2 Limitations of Zork	9
4.3 Voice Recognition Implementations	10
4.3.1 Tazti - Speech Recognition for PC Games	10
4.3.2 PocketSphinx for Unreal Engine 4	10
4.3.3 Houndify	10
4.3.4 Watson Conversation Engine	11
4.4 Natural Language Processing	13
4.4.1 Stages of Natural Language Processing	13
4.4.2 Part-of-Speech Tagging	15
4.4.3 Slot Filling	15
4.4.4 Synonyms and Hypernyms	16
4.4.5 Semantic Similarity	17
4.5 WordNet	18
4.6 Generating Objects from Descriptions	18
4.7 Generating Descriptions from Objects	19
4.8 Object Properties	19
5 Design	21
5.1 Target Platform	21
5.2 Java Programming Language	21
5.3 System Overview	21
5.3.1 Speech-to-Text	22
5.3.2 POS Tagging	23
5.3.3 Slot-filling Matching	23
5.3.4 Semantic Similarity	25
5.3.5 Context-Action Mapping	26
5.3.6 Intent Execution	27

5.4	System UML Class Diagram	27
5.4.1	<code>MainActivity</code> Class	28
5.4.2	<code>VoiceControl</code> Class	29
5.4.3	<code>VoiceProcess</code> Class	29
5.4.4	<code>GlobalState</code> Class	29
5.4.5	<code>ContextActionMap</code> Class	29
5.4.6	<code>SemanticSimilarity</code> Class	29
5.5	System Features	30
5.5.1	Synonym Mapping	30
5.5.2	Ignoring Incorrect Matches	30
5.5.3	Confirmation and Suggestions	30
5.5.4	Multiple Commands	30
5.5.5	Multiple Targets	30
5.5.6	Sentence Mapping	31
5.6	Game Design Overview	32
5.6.1	Game Flow	32
5.7	Game Design UML Class Diagram	33
5.8	Room Generation	35
6	Implementation	36
6.1	System Architecture	36
6.2	Voice Recognition Interface	37
6.3	Graphical User Interface	37
6.4	User Flow Diagram	39
6.5	Voice Recognition System Implementation	40
6.5.1	Data Flow Diagram	40
6.5.2	<code>GlobalState</code>	40
6.5.3	<code>VoiceProcess</code>	41
6.5.4	<code>AmbiguousHandler</code>	41
6.5.5	<code>SemanticSimilarity</code>	42
6.5.6	<code>CustomLexicalDatabase</code>	42
6.5.7	<code>ContextActionMap</code>	44
6.5.8	<code>Entity</code>	45
6.5.9	<code>Action</code>	45
6.5.10	<code>MultipleCommandProcess</code>	46
6.5.11	<code>SentenceMapper</code>	46
6.5.12	Synonym Mapping Implementation	47
6.6	Role-Playing Game Implementation	47
6.6.1	<code>GameState</code>	47
6.6.2	Battle Mode	48
6.6.3	Overworld Mode	48
6.7	Object Properties with WordNet	49
6.8	Room Generation Implementation	49
6.8.1	Rooms	49
6.8.2	Manual Room Generation	50

6.8.3	Using BooleanSuppliers	50
6.8.4	Room Generation Program	51
6.9	Video Conferencing Demo	52
6.10	Settings Activity	52
7	Evaluation	54
7.1	Mock Testing	54
7.2	Evaluation of Semantic Similarity	55
7.2.1	Combining Similarity Methods	57
7.3	Performance of Semantic Similarity Methods	59
7.4	Evaluation of System Features	60
7.5	Evaluation of Slot-Filling Grammar	62
7.6	Application Performance	63
7.7	Usability of the Voice Recognition System	63
7.8	Voice Recognition System Limitations	65
7.9	Room Generation Evaluation	66
8	Conclusion	67
8.1	Deliverables	67
8.2	Contributions	67
8.3	Future Work	68
8.3.1	Language Support	68
8.3.2	Deep Learning	68
8.3.3	Multi-threading	69
8.3.4	Hands-Free UI	69
Appendix A	Android Application Screenshots	70
A.1	Battle Mode Example	70
A.2	Overworld Mode Example	72
A.3	Multiple Commands Example	74
A.4	Ambiguous Utensil Example	76
A.5	Duplicate Contacts	78
A.6	Settings User Interface	79
Appendix B	Diagrams and Examples	81
B.1	System UML Class Diagram	81
B.2	Game Design UML Class Diagram	82
B.3	Battle Mode Example Commands	83
B.4	Overworld Mode Example Commands	84
B.5	Video Conferencing Example Commands	85
B.6	Room Generation Example	86

Appendix C Testing and Evaluation	87
C.1 List of Android Mock Tests	87
C.1.1 BattleTest	87
C.1.2 OverworldTest	89
C.1.3 CallTest	90
C.2 Actions and Contexts used in Evaluation	91
C.3 Evaluation Results for WUP	92
C.4 Performance Measurement Commands	94
C.5 Comparison of Usability	95
C.6 Other Survey Results	97
C.7 Room Generation Evaluation Results	99
Appendix D Supplementary Materials	100
D.1 Skyrim Kinect Command List	100
D.2 Penn Treebank Tagset	102
D.3 Pseudocode for ProcessInput()	103
References	104

1 Abstract

When adding voice commands to a video game, a developer may find themselves adding each variation of a phrase that a player can say, and mapping them to intents within the game. It becomes more work to include varying phrases of the same intent as the developer would have to include every permutation of the phrase (for which there could be infinitely many).

This project presents a text-based role-playing game (RPG) on Android which allows the player to issue actions within the game using their voice. The game uses a new voice recognition system that aims to make it easier for the developer to add voice commands without hard-coding the phrases, and works completely offline with no cloud processing required. The system uses the WordNet database created by Princeton University to process voice commands.

Various techniques for improving the voice recognition system are explored and evaluated, with the aim of decreasing the developer workload. These include using semantic similarity methods, learning mechanisms, and more. The voice recognition system is also applied to other domains such as video conferencing solutions, and cooking applications.

This project also explores other areas for easing the development of the RPG using WordNet, such as automatically assigning properties to objects, and room generation from text descriptions.

2 Acknowledgements

With thanks to Dr. Mark Wheelhouse of the Department of Computing at Imperial College London, for his guidance and supervision over this project.

Also many thanks to my family and friends for all their support.

3 Introduction

In recent years, video games have explored various forms of input that diverge from the traditional control schemes of a keyboard and mouse, or an analogue stick and buttons on a controller. Touch-screen controls and motion controls have had varying degrees of success over the years, and provide new forms of interactions with games. A recent form of input used in video games has been voice controls; a user utters a sentence or phrase and this would execute an action or intent within the game.

Voice recognition and control schemes have been used in various ways in the medium, whether in tandem with traditional control schemes, or as the only form of input. Most games which feature a voice control scheme are typically programmed to work with a specific set of keywords or phrases, hard-coded by the developer. This is usually the case with games where the voice control interface is entirely optional to the user, as it requires little to no input from the programmer; if the speech processing is handled by a separate library, the developer just has to map a text string (e.g. “open the door”) to the function that executes the intent (e.g. the function that opens a nearby door in the game).

3.1 Role-Playing Games

Role-playing games (RPG) are a genre of video games that involve the player controlling a character in a world featuring exploration and/or battle mechanics, along with a progression system, where the character gains new abilities or items as the game progresses. These include turn-based games such as Pok  mon and Final Fantasy, or more action-based games such as Dark Souls and The Elder Scrolls.

In RPGs, the player usually possesses several items that they acquire during the game, and accesses them by navigating menus. If the player has many different items, they may need to spend a lot of time searching the menus for a specific item, which can be quite tedious. It can also become quite repetitive if the user is constantly only pushing the “attack” button over and over again with no variation in the attack used. Using voice commands, a player could just simply say a phrase such as, “use a potion” or, “attack with an axe”, and the item will be used without having the player navigate several menu screens.

3.2 Problem Statement

It becomes increasingly difficult for a developer to add more phrases that can be accepted in the game and map to the same intent (e.g. when the user says, “push the door open” instead of, “open the door”) as they would have to hard-code each possible input that the user could possibly say: having too few phrases would mean that the user could become frustrated when their preferred way of stating an intent is not accepted by the game, while adding too many acceptable phrases would increase development time.

Natural Language Processing (NLP) has been used to attempt to understand what users

have said, and to extract meaning from the user’s phrases and sentences. For instance, if the user says, “push the door open”, NLP can be used to infer that the user’s intent is to open a door. Using NLP, several user phrases can be mapped to the same intent without having to hard-code each possible phrase.

This project will apply NLP to a voice control scheme used within a role-playing game, in order to reduce the amount of work required by a developer to add a flexible voice control interface, and to give more freedom of expression to the player while they play the game.

3.3 Motivation

NLP is already used to improve voice recognition in popular personal assistants such as Siri or Alexa, but these systems offload the NLP workload to the server in order to reduce the amount of local processing required, as modern NLP techniques require a sizeable amount of machine learning and pattern recognition [1].

Most video games do not require an internet connection to play, and therefore can be played anywhere, so it may not be possible to use these cloud services for adding voice recognition systems to games. In other applications such as embedded systems or remote robotics, no internet connection can be established, so the only option is to hard-code the acceptable voice commands. There may also be privacy concerns when sending audio recordings to an external server.

The voice recognition system proposed in this project works completely offline, but also supports near-instantaneous processing on modern mobile devices, and can be applied to other applications as well as video games.

Another motivation for this project is that very few popular games support voice recognition as a form of control, whether it is on consoles, PC or mobile. Voice controls allow games to become more accessible to people with certain disabilities that prevent them from using traditional controllers or a keyboard and mouse.

3.4 Objectives

A text-based RPG will be created that will use a flexible voice control interface. It will feature a simple exploration mechanic, with the user manipulating objects and environments using voice commands, and a simple battle mechanic. These will be designed to fully demonstrate the power of the voice-control system to make the game easier to develop, as well as easier to play (as opposed to using traditional mouse clicks and button presses).

Note that this project is not concerned with the real time digital processing required to convert speech to text; it is assumed that the speech from the microphone is already converted to a text string using an existing tool such as Google’s Speech-to-Text API. The focus is on the integration and usability of systems to process user commands.

As well as being designed to easily add voice commands, the RPG will be designed such that other areas of development are made easier, such as assigning object properties and easily creating new rooms, although the main focus will be on the voice recognition system.

The proposed design for the RPG and system aims to achieve the following goals:

- Allow a developer to add voice commands with little effort, using a system which works completely offline with near-instantaneous processing on modern devices, and applying the system to different domains such as games and video conferencing solutions.
- Evaluate the effectiveness of using different semantic similarity methods with Word-Net for mapping voice commands to actions, as well as evaluating the accuracy and performance of each method.
- Apply techniques used in the voice recognition system to other areas for improving the development of the RPG, including generating rooms from text descriptions, and automatically assigning properties to objects.

4 Background

4.1 Voice Recognition in Games

Below is a list of notable and recent games that have incorporated voice control functionality as a major component of the interactive experience.

4.1.1 In Verbis Virtus

In Verbis Virtus is an independent 3D fantasy adventure game developed by Indomitus Games for Windows [3]. It requires the player to solve puzzles in a 3D environment and battle enemies using magic.

When a microphone is connected, the player is able to cast spells using their voice by saying specific phrases defined by the game. For example, to cast a spell that produces a floating light source to brighten a room, a user must say, “let there be light”, or if the user wishes to shoot an energy beam from their hand, they must say, “Beam of Light”, and so. These phrases appear to be hard-coded into the game.

There are several limitations with this implementation of voice control. Firstly, since the phrases are hard-coded, there is no lenience in variations of the phrases, e.g. “create a light” or “Laser of Light”, etc. This can seem mundane and no different to a user pressing the same button repeatedly to cast a spell.

Another limitation is that the user is still required to use a keyboard and mouse to perform other actions in the game, such as moving around, looking around, and navigating menus and message boxes. The voice control interface cannot be used to perform these actions. The voice control interface itself does not seem integral to the experience, as the user could just simply press a key that performs the spell instead of saying the same phrase over and over again.

This voice recognition interface could be improved by allowing the player to perform any action within the game using their voice, as well as allowing for more variation in what the user can say for each action.



Figure 1: A screenshot of *In Verbis Virtus*, where the player is casting a spell. Image taken from the Indiegogo web page [2].

4.1.2 Skyrim with Kinect

The Elder Scrolls V: Skyrim is a first-person action role-playing game developed by Bethesda [4]. The Xbox 360 version of the game supports the Microsoft Kinect peripheral which allows the player to execute voice commands [5]. The user can use over 200 hundred pre-configured voice commands, and these can be used simultaneously with the traditional control scheme.

There are voice commands available that allow the user to navigate the game's menu screens. For instance, the player would say, “quick items” to open their inventory menu. Afterwards, the user can open menus for specific categories of items such as “potions”, “books”, etc. However, it doesn't seem to be possible to use any of the items selected from the menu within the game using only a voice command; this is only possible if you assign the item to a ‘hotkey’ (e.g. by saying, “assign health potion” once you have highlighted a potion item), and then saying “equip health potion” during the game. See Appendix D.1 for a full list of available voice commands.



Figure 2: A screenshot from *The Elder Scrolls V: Skyrim*. Image taken from the game's official web page [4].



Figure 3: Voice commands that are available within the items menu in Skyrim.

This system, despite being optional to the player, helps to reduce the time spent searching through the menu screens. However, it appears that this system also hard-codes the phrases that can be said by the player with no flexibility in how the player can say the command.

4.1.3 Star Trek: Bridge Crew

Star Trek: Bridge Crew is a virtual reality game developed by Ubisoft. This game allows the player to issue voice commands to AI crew members on a spaceship [7]. In order to develop a more interactive and realistic experience, the development team used IBM Watson's interactive speech capabilities.

The API allows commands to be delivered in a variety of ways, as the speech is parsed for its meaning using Watson's Conversation service. For example, a player could say, "show me the ship's status report", or they could say it in a different way such as, "damage report" to execute the same command [8]. This provides the user with more freedom in how they convey their request to the AI, giving the player a "new level of sense of presence" [9]. See Section 4.3.4 for more details on IBM's Watson API.



Figure 4: A screenshot of *Star Trek: Bridge Crew*. Image taken from the Steam web page [6].

The voice recognition design used in this game is similar to the proposed design outlined in this report. The user's intent should be extracted from the phrase that they speak, so similar phrases should map to the same intent. Unfortunately, Watson's API is a paid cloud service; so the game requires a constant internet connection to play.

4.1.4 Classic Zork on Alexa

Zork is a trilogy of classic text-based role-playing games that support a free-form style of input (e.g. in a terminal). The player is provided with an input scenario (e.g. "You are standing in an open field west of a white house, with a boarded front door. There is a small mailbox here."), and the user types commands that they want to execute (e.g. "open mailbox" or "attack troll with sword") [10]. The exploration part of the proposed game in this report will be loosely based off Zork's exploration mechanics of interacting with objects in the surroundings.

To add voice control to Zork, these commands could be delivered using a speech-to-text

interface. A group of developers created a port of Zork running on Amazon’s Alexa Skills Kit [11]. This allows the player to deliver commands using their voice as if they were typing the phrases into a terminal. However, this implementation still suffers from the limitations of Zork’s free-form input.

4.2 Limitations of Zork

The exploration mechanics of the proposed design in this report will be based on Zork’s general gameplay of interacting with objects in the environment to solve puzzles, and the limitations of Zork’s free-form input are discussed here.

Zork supports a number of text commands, which usually take the form of a verb-noun structure, such as “use <noun>” or “take <noun>”, but it also supports some more sophisticated commands with complex structures, such as “give all but the pencil to the nymph” or “drop all except the dart gun” [12].

Zork’s free-form input still has some limitations. For example, it still cannot accept all variations of specific intents. For instance, below is a list of accepted and rejected commands if the user wishes to open a mailbox in front of them:

- “open mailbox” — Accepted
- “open the small mailbox” — Accepted
- “can i open the mailbox” — Rejected
- “check the mailbox” — Rejected
- “find out what’s in the mailbox” — Rejected

Below is a list of results for closing the mailbox:

- “close mailbox” — Accepted
- “shut mailbox” — Rejected
- “close the letterbox” — Rejected
- “close the red mailbox” — Rejected
- “close the postbox” — Rejected

The rejected text commands could be valid phrases that express an intent to open some sort of mailbox. Zork has a pre-defined list of verbs and sentence structures that it can accept, and sticks to those without much flexibility.

Since most of the supported commands take the form of a verb-noun structure, this could be used in the proposed design as the general structure of commands expected by the user.

4.3 Voice Recognition Implementations

There are many tools and services available that allow developers and users to add voice recognition to video games, with some discussed here:

4.3.1 Tazti - Speech Recognition for PC Games

Tazti is a keyboard mapping tool available to players to be used with any type of PC game that uses the keyboard [13]. The user sets up a profile for each game separately by mapping speech commands to one or more keys. For instance, the user could say “fire” or “shoot fire” and this would map to the keystrokes required to cast a fire spell.

While this tool is not directly integrated into games by the developer, it has the advantage of being able to work with almost any PC game that uses a keyboard, such as role-playing games, first-person shooter games, and action games. However, the voice commands are limited to being hard-coded, so there is no way of varying the speech commands slightly, even if the meaning is the same. It is also not freely available and a license needs to be purchased.

4.3.2 PocketSphinx for Unreal Engine 4

Sphinx-UE4 is a speech recognition plugin for the game engine, Unreal Engine 4 [14], based on the PocketSphinx library which provides offline speech-to-text [15]. The plugin allows the developer to specify keywords and the commands that they map to (e.g. “turn right”, “enable sprint”, “kick the ball”, etc).

The plugin also supports grammar files. A grammar contains the grammatical structure that an accepted input can take, such as, <digit> <operation> <digit>. This will accept input such as, “three add five”, “six minus four”, “two times twelve”, and so on.

4.3.3 Houndify

Houndify is a paid cloud service developed by SoundHound that allows developers to add voice recognition and control to any application they wish [16]. The API takes as input either a text string or audio samples and returns a JSON string containing the response to the input. This response can range from anything: from home automation control, to weather updates, and more. The response would be in the form of an intent - a single word representing the action to be executed. For example, any attacking phrase would map to an intent called, “ATTACK”.

As well as having built-in voice commands already for specific domains (such as all weather commands or sports update commands), Houndify also supports the creation of custom voice commands by the developer. To create a custom command, the expression to be said by the user is specified using a syntax similar to regular expressions, but instead specifies the general phrase structure, and the expression is mapped to an intent [17].

Below is an example of a possible expression for attacking a troll:

```
"attack" . "troll"
```

Here, the player would have to say, “attack troll” to execute the intent. This can be expanded to have more variation in how the user can say this command:

```
(“attack” | “hit”) . [“the”] . “troll”
```

Here, the player can say phrases such as, “attack the troll” or “hit troll”. The expression can be further expanded. However, the expression will eventually become long and confusing:

```
(“attack” | “hit” | “damage”) . [“the”] . [“big” | “large”] . “troll” . [(“with” | “using”) .  
[“a”] . (“sword” | “hammer” | “axe”)]
```

Even though this expression allows for more variation in what the user can say, it still doesn’t cover all the possible variations of specifying an attack, and this long expression is only for one intent; many more expressions would have to be written for other intents, and this can become tedious for the developer.

While this API would be useful for building a personal assistant similar to Apple’s Siri or Amazon’s Alexa, it does not seem feasible for this project due to the time required to write the expressions for different commands. This is also a paid service, so users are charged for each query they make. It also requires a constant Internet connection to use.

4.3.4 Watson Conversation Engine

IBM’s Watson API features a Conversation service that allows developers to build voice recognition interfaces that understand natural language input [18]. The developer defines the training data to be used by the API (in order to train a natural language classifier) using ‘intents’ and ‘entities’ [19].

‘Intents’ are the goals that a user will have as they interact with the application. For example, in a voice-controlled car, an intent could be `wipers_on`, in order to activate the windscreen wipers. This intent would be paired with example utterances in the form of text strings, such as “turn on the wipers”, “switch on the windscreen wipers”, and so on. The more examples given, the more accurate the trained model will be.

#return

 Add a new user example...

exchange

I'd like to return this

I don't want this anymore

I need you to take this back

Figure 5: An example of an intent to return something, along with examples of phrases for the intent. Image re-used from an IBM blog post [20].

‘Entities’ are classes of objects that help to provide context to intents. For example, with the `wipers_on` commands above, it is not clear whether the user is referring to the wipers on the front of the car or the wipers on the back of the car (the API will assume the front by default). If the entity is included in the phrase spoken by the user, then the context of the intent becomes clear.

@returnItems

 Add a new value

<input type="checkbox"/> book	text	tome
<input type="checkbox"/> parrot	bird	macaw Norwegian Blue
<input type="checkbox"/> video cassette	movie	tape

Figure 6: An example of entities for items that can be returned. Entities that have been grouped together (usually synonyms) are placed on the same row [20].

Once the classifier is trained, it will be able to use the examples of intents provided to classify whether new examples of phrases have the same meaning as those. For example,

if the user says, “activate the windscreen wipers at the front of the car”, the classifier will return the `wipers_on` intent, even though this input is very different in structure to the examples provided above. A demo is available to test this in a web browser [21].

According to a research paper by McCord, Murdok and Boguraev on *Deep Parsing in Watson* [22], the engine consists of two deep parsing components: an English Slot Grammar (ESG) parser and a predicate-argument structure (PAS) builder. These allow Watson to produce parse trees of a sentence and extract pattern-based relations from it, such as question decomposition, hypothesis generation, and evidence scoring.

This API is very flexible and matches the first objective of this project; only a few examples of possible inputs need to be provided for each intent, and then the classifier can infer whether any new input strings will match the same intent. Unfortunately, despite being a very powerful system, this still remains a paid cloud service, so users are charged for each API request. It also requires a lot of processing power (which is offloaded to a cloud server) for the many machine learning algorithms being used here. However, a simpler system may possibly be implemented that uses a similar concept to the intents and entities framework, and the proposed design in this report uses a similar concept that runs locally.

There is a similar API called DialogFlow [23], which works very similarly to IBM’s Watson Conversation engine. However, it suffers from the same issues as above (e.g. cloud-based, pay-per-request, etc).

4.4 Natural Language Processing

Natural Language Processing (NLP) forms a major component of this project. This will give the developer the ability to take an arbitrary string of text and infer its meaning, and map it to an intent. Below are some explanations on NLP concepts related to this project.

4.4.1 Stages of Natural Language Processing

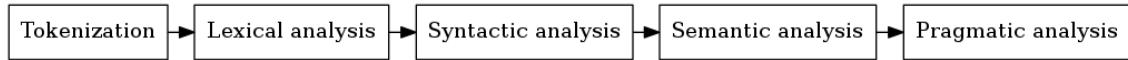


Figure 7: The general stages of processing a text string for its meaning using NLP.

NLP can be broken down into six stages [24]. The first stage is the tokenization, where a raw text string is broken down into words (usually separated by spaces).

The second stage is the lexical analysis, where the output of the tokenization process is improved upon by looking at words that can be taken apart even further (or rejoined

if needed) to uncover more information. These include words such as: clitic contractions (e.g. “what’re” would tokenise to “what” and “are”, while “we’re” would break down to “we” and “are”); end-of-line hyphens that split whole words into parts when using a justification alignment; and shortened titles or abbreviations (e.g. Dr., U.S.A., etc.) [25].

The third stage is syntactic parsing, where the grammatical structure of the sentence is determined. This is usually achieved by generating a syntax tree of each sentence, where each word is broken down into a terminal from a given grammar (e.g. verbs, nouns, adjectives, determiners, etc.) [26].

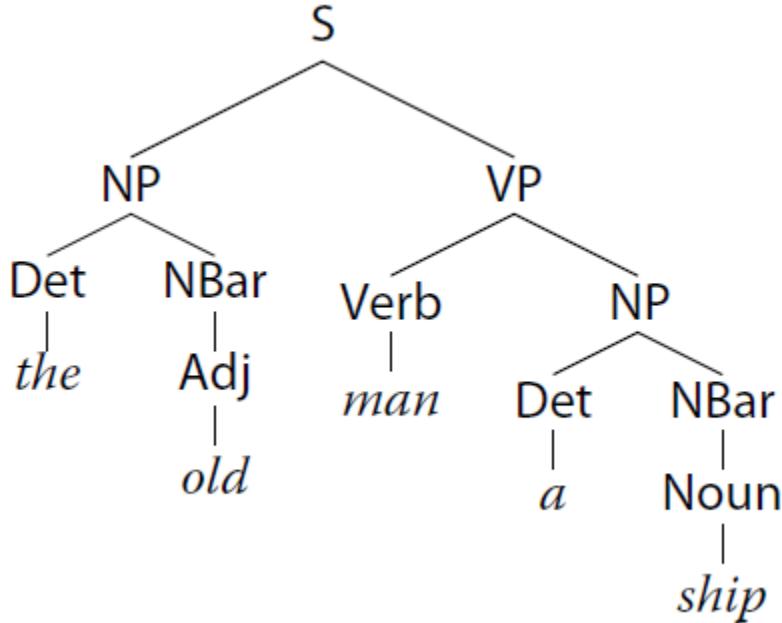


Figure 8: An example of a syntax tree for the sentence, “the old man a ship”. Diagram taken from the Handbook of Natural Language Processing [26].

In the example in Figure 8 for the sentence, “the old man a ship”, the word, “man” could be mistaken to be a noun, but in this case, the sentence is implying that the elderly control a ship, and therefore, “man” is actually a verb here. This meaning can be more easily inferred when the sentence is broken down using a syntax tree, as above.

The final two stages are very similar and is difficult to separate them into separate stages. Both stages involve determining the meaning of the sentence. The fourth stage is often concerned with semantic analysis, whereas the last stage is more concerned with discourse analysis [25].

4.4.2 Part-of-Speech Tagging

A technique often used in NLP is to try to label each word with its correct part of speech, known as part-of-speech (POS) tagging. POS-tagging systems generally use a tagset containing POS tags to assign a tag to each word. The most popular tagset is the Penn Treebank tagset. This consists of 48 tags, of which 12 are for punctuation and other symbols [27]. See Appendix D.2 for the full Penn Treebank tagset.

There are two main challenges regarding POS tagging. The first challenge is that words can sometimes be ambiguous. That is, the word can be assigned a number of possible POS tags depending on its context. For instance, the word, “can” can either be a verb (“I can...”), a noun (“a can of tuna”), or another verb with a completely different meaning (“can that noise!”) [28]. The many meanings that are attributed to a word are known as its *word senses*.

The second challenge involves words that are unknown to the tagger and cannot be tagged. A default tag is required for words which are unknown, but that could interfere with the rest of the POS-tagging process [27].

The proposed design for this project will use an open-source POS tagger in order to identify specific words in the player’s input, namely verbs and nouns.

4.4.3 Slot Filling

According to the authors of *Speech and Language Processing* [29], three tasks need to be done to understand a user’s utterance. The first is *domain classification*, to determine the topic that the user is talking about (e.g. weather, sport, etc), although this is unnecessary for single-domain applications such as the contents of a video game. The second task is *intent determination*, to determine the goal of the user, and finally, the last task is to do *slot filling*.

Slot filling involves determining the specific details of an intent by defining a grammar that the intent must adhere to. An example for arithmetic operations would be, <DIGIT1> <OPERATION> <DIGIT2>. The input would be parsed by a context-free grammar parsing algorithm to extract each bit of information in order to fill each slot in the grammar. Figure 9 shows an example of a slot-filling grammar.

A disadvantage to slot-filling using context-free grammar parsing is that not all valid expressions for an intent would match with the slots. A solution to this would be to define multiple grammars for different structures of inputs.

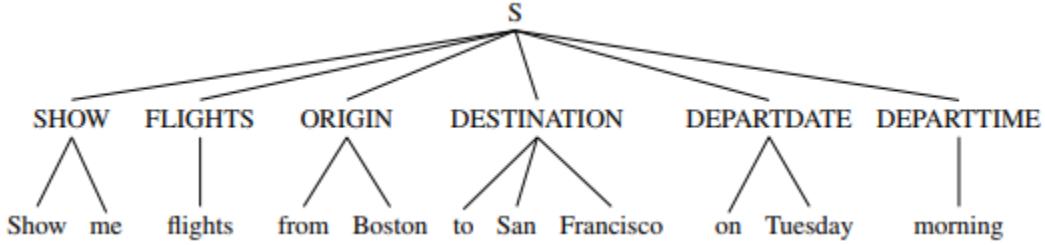


Figure 9: An example of a slot-filling grammar for flights, using slot names as the internal parse tree nodes. Diagram taken from Chapter 29 of the book, *Speech and Language Processing* [29].

4.4.4 Synonyms and Hypernyms

Synonyms are words or phrases which have a similar meaning to another word [30]. For example, the word “dog” has seven word senses (i.e. different meanings), and for one of the word senses with a meaning of, “a member of the genus Canis”, the synonyms would be, “domestic dog” and “Canis familiaris” [31]. Another word sense for “dog” would have a meaning of someone who is morally reprehensible (e.g. “you dirty dog”), and this sense would have its own synonyms.

A *hypernym*, on the other hand, is a word or phrase which generalises words for which it is a hypernym of. For instance, the word, “instrument” is a hypernym of “guitar”, because a guitar is a type of instrument. Similarly, the word “animal” is a hypernym of “cat” (or, at least, the sense which refers to the domestic animal, as “cat” can have several usages).

The synonyms and hypernyms of words and phrases can be represented as a hyper-tree structure, where the parent nodes are hypernyms of the child nodes. Each node is also connected to nodes that represent the other senses. For example, the word, ”dog” would have a node for each of its senses (and each one containing the synonyms for that sense).

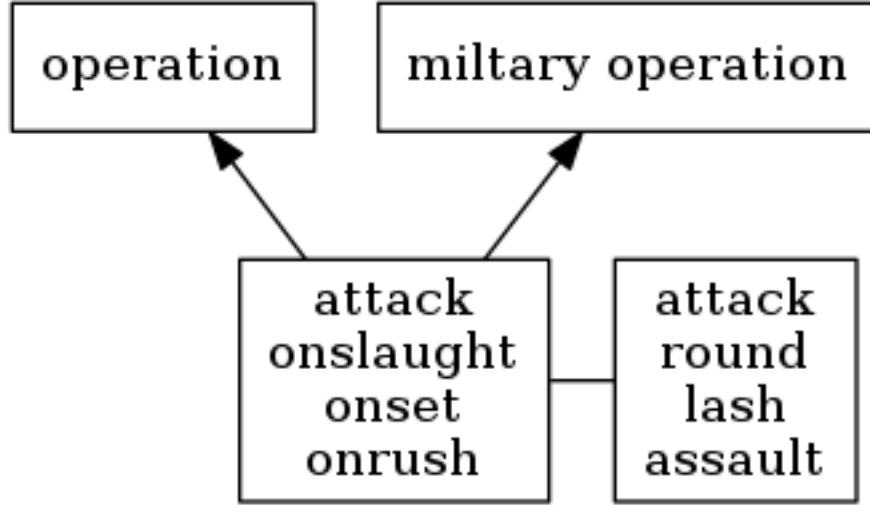


Figure 10: A partial hypernym tree showing the synonyms and hypernyms for one of the senses of the noun, “attack”. The “attack” verb sense (right) is a different node.

Synonyms and hypernyms could be used to determine the semantic similarity between words and infer whether the user’s intent matches a certain action.

Note that a *hyponym* is the inverse of a hypernym (e.g. the word, “cat” is a hyponym of “animal”).

4.4.5 Semantic Similarity

There are many ways of calculating the semantic similarity between two words or phrases. Many of these methods use the synonym and hypernym trees of the words to infer similarity, and assign a ‘similarity score’ to them. This score could be used to determine whether, for instance, a verb has semantic similarities to one of the actions in a game, and therefore can be mapped to that action.

The following are some examples of measures for semantic similarity [32]. The first is the Wu and Palmer (WUP) measure, which takes the depth in the hypernym tree of the Least Common Subsumer (LCS) (that is, the first hypernym that both words have in common) and compares it with the sum of the depths of each word. The formula for the similarity of words w_1 and w_2 is as follows:

$$sim_{wup} = \frac{2 * depth(LCS(w_1, w_2))}{depth(w_1) + depth(w_2)}$$

For example, the WUP score for the semantic similarity between “juice” and “water” is calculated as follows [33]:

```

1 T1 = HypernymTree( juice ) =
2     [Sense #1] *ROOT* < entity < physical_entity < matter
3         < substance < food < foodstuff < juice
  
```

```

5 T2 = HypernymTree( water ) =
    [Sense #3] *ROOT* < entity < physical_entity < matter
7     < substance < food < water
9 Lowest Common Subsumer(s) = argmax(depth(subsumer(T1,T2)))
    = { subsumer(T1[1], T2[3]) } = { food }
11 DepthLCS = depth( food ) = 6
13 Depth1 = min(depth( {tree in T1 | tree contains LCS } )) = 8
Depth2 = min(depth( {tree in T2 | tree contains LCS } )) = 7
15 Score = 2 * DepthLCS / ( Depth1 + Depth2 ) = 2 * 6 / ( 8 + 7 )
Score = 0.8

```

The WUP score of 0.8 out of 1.0 indicates that juice and water are semantically very similar [34], whereas the words “chair” and “biology” only have a WUP score of 0.17, suggesting they are not semantically similar.

Another measure of semantic similarity is Leacock and Chodorow’s method, which relies on the shortest distance between the two words in the hypernym tree, divided by the maximum depth of the whole tree; the shorter the distance, the more similar the words are [32]. There are many more different measures and it is possible to combine them to produce an overall score for semantic similarity.

4.5 WordNet

Princeton University created a large lexical database for the English language, called *WordNet* [35]. Words are grouped into synonyms (known as synonym sets, or *synsets* for short), and are connected to their respective hypernyms (more general words) and hyponyms (more specific words), creating the large hyper-tree mentioned in Section 4.4.4. The database is open-source and is used extensively in this project, using various APIs that interact with the WordNet database (which will be initially downloaded by the game for offline use).

4.6 Generating Objects from Descriptions

One feature for the proposed design of the RPG is to prevent the developer from manually placing objects in a room/environment within the game. A designer may want to just provide a description of the room without having to program anything. A simple solution to this would be to perform noun extraction on the description to extract the objects that should be in the room. However, it would also be useful to extract relations of objects to other objects as well. For instance, a lamp should only be on a table if the table exists in the room; if the table is removed/broken by the player, then the lamp should now be on the floor. This would allow for more information to be extracted about the objects in the room, as opposed to simply their existence.

Fader, Soderland and Etzioni from the University of Washington proposed a solution for extracting these types of binary classifications using an extraction algorithm which takes as input a POS-tagged input and returns a set of triples representing binary relationships.

They created a tool for this called *ReVerb* [36]. Taking the sentence, “A knife is on the floor”, it would extract the triple, (*a knife*, *be on*, *the floor*), representing the binary relationship (“be on”) between the knife and the floor. Using this, more interesting room generation may be possible just from a user-given description of the room.

4.7 Generating Descriptions from Objects

The opposite to the above is to be able to generate text-based descriptions of the current situation that the player is in based on objects around them. For instance, a description could be, “You are in a room lit by a candle on a table. There is a piece of broken glass on the floor and a door in front of you.” One could hard-code this sentence for each room the player visits, but if there are a lot of rooms in the game, this could become time-consuming. Developers would like to generate a description of the current situation based on sparse pieces of information (such as the items in the room, plus their locations, in a vector representation). This is the opposite of the above, where the object of a sentence/phrase is extracted; developers may want to go in the opposite direction.

This is a relatively new field in NLP as it is particularly difficult. One method of sentence generation is suggested by Iyyer et al [37], where recursive autoencoders (a type of neural network) are used to generate a decomposition model of the sentence that can be used to reconstruct the sentence again. However, they found that proper nouns were not reconstructed correctly, and the model would need to be built in the first place to reconstruct the sentence.

Another group of researchers proposed a method of sentence generation which involves taking a sum of word embedding vectors and modelling it as a mixed integer programming problem [38]. However, this also requires the original sentence to be de-constructed in the first place in order to know what the vector should contain.

One possible solution would be to only hard-code the general sentence structure of the description of rooms, and then insert the appropriate adjectives and nouns of the items, which is given as input in the form of a vector of strings.

4.8 Object Properties

One issue to note is when a user attempts to specify an item based on a description of it (e.g. the user says, “attack the enemy with a sharp object”); players would like to be able to select an item that matches the description (e.g. a sword is sharp, so it will be chosen if the user possesses one). One approach to this is to assign the items a set of properties such as whether it is sharp or blunt, whether it can be thrown, and so on.

A similar technique was used by the developers of Scribblenauts for the Nintendo DS. Using the company’s ObjectNaut engine, the game utilises a large database of a hierarchy of objects [39]. This allows them to design sprites for different objects (such as a cyborg, robot or android), but assign them the same interaction properties. However, this still

required the developers to go through encyclopaedias word-by-word to find which ones had similar properties. Using the semantic similarity scoring mechanism, this would not need to be done, as it can be inferred whether an object has a certain property based on how semantically similar it is to that property (e.g. an object may be sharp if it is semantically similar to the word, "sharp"), instead of hard-coding each property for every object in the game.

By assigning items different properties, the actions that the user can perform for each variation of items (e.g. sword, katana, knife, etc.) does not need to be hard-coded. Instead, there can be general objects which have certain properties, and then the actions for each general object can be coded. For example, there could be a general object with a 'sharpness' property, that would be able to cut things, and all swords and knives would be of this type of general object, and therefore also considered sharp.

5 Design

5.1 Target Platform

According to Ovum's Mobile Games Market forecast for 2017-2022 [40], mobile games continue to increase their market share as they grow in popularity, and within the next two years, almost half of the total video game revenue generated will come from mobile games. Therefore, it makes sense to design a system for a mobile platform such as iOS or Android. It also allows for hands-free input as a future extension.

The chosen platform to develop the voice recognition system and game will be Android mobile devices, as the availability of the open-source Google Speech-to-Text API means that less effort is required to transform the user's utterance into text.

5.2 Java Programming Language

The language used to develop the application is Java, which is an object-oriented programming (OOP) language. The use of an OOP language makes it easier to create objects in a video game using the inheritance hierarchy. For example, there would be an abstract class for physical objects called `PhysicalObject`, and another class such as `Glass` would be a derived class of a `PhysicalObject`, and so on.

A disadvantage to using Java is that it does not support multiple inheritance, unlike other languages with OOP constructs such as C++, due to methods such as `super()` having the same signature for all parent classes. For instance, a `Glass` class cannot inherit properties from both a `PhysicalObject` and a `BreakableObject` (i.e. they are both direct parents of the `Glass` class). Instead, this must be overcome by making the inheritance a single chain (e.g. `Glass` inherits from `BreakableObject` which inherits from `PhysicalObject`). While this can be considered a hindrance in some circumstances, it is not a major issue for the proposed design of the game.

While Java applications — and therefore Android — support the invocation of native code such as C++ using the Java Native Interface (JNI) for improved performance, this will not be used due to time constraints for the project.

5.3 System Overview

The following is a brief description of the voice recognition system, including the main components and the general flow of data.

The voice recognition system takes as input the raw microphone input containing the player's utterance, and outputs a response as a string, with the side-effect being the execution of a developer-defined action based on the player's utterance. Figure 11 shows a simple flow chart outlining the main stages of the voice recognition system from input to output.

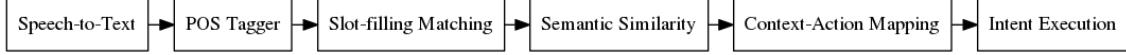


Figure 11: A flow chart showing the general stages of the voice recognition system.
Graph generated using Graphviz.

The steps are briefly described here and explained in more detail below.

- The audio from the microphone input is processed into a text string using a Speech-to-Text (STT) engine. This string is then split up into a list of words.
- The list of words is then used to generate a list of Part-of-Speech (POS) tags. Each word in the list is passed to the POS tagger and added to the tag list.
- The word list is parsed using a slot-filling grammar, e.g. "<ACTION> <TARGET> [with] <CONTEXT>", or "[use] <CONTEXT> [to] <ACTION> <TARGET>".
- When performing the matching, the current word being parsed is sent to a semantic similarity engine, and a similarity score is calculated.
- If the player has specified an action to be done under a certain context, then the context-action mapping is consulted to find the correct method to execute. Otherwise, the default action method is chosen.
- Finally, the in-game function is executed, and the text output is returned, with the game state updated.

5.3.1 Speech-to-Text

Android features a built-in speech recognition service, via its `SpeechRecognizer` class [41], and this is used for the speech-to-text engine. While the API's primary method of speech recognition is via remote server communication, there is an option with the Android OS settings to perform offline speech recognition instead, by downloading the language data (although this offline method has slightly less accuracy) [42].

Using this API, the player's utterance is transformed into a text string, which is then processed by the rest of the voice recognition system. Using an existing speech-to-text API means that less time is spent trying to get the player's utterance from the microphone input.

In terms of user interface, the application features an on-screen button containing the image of a microphone. When the user presses this button, a voice request is initiated and lasts until the user finishes uttering their input. If the user doesn't say anything, then the response times out and is cancelled.

5.3.2 POS Tagging

For the Part-of-Speech (POS) tagging of the input, the open-source log-linear POS Tagger provided by The Stanford Natural Language Processing Group will be used [43]. It is a Java archive file and can easily be imported and used in an Android Java project. The API handles all the POS tagging so no time needs to be spent implementing a POS-tagging system.

The library comes with two trained models. According to the software's README file, the first model uses a bidirectional architecture and includes training data using word shapes and distributional similarities. The second model uses a *left3words* architecture (which also uses word shape features), and the tags in the output are from the Penn Treebank tagset. According to the developers, the bidirectional model is "slightly more accurate" with a correctness score of 97.28%, but is much slower to tag with than the *left3words* model. Therefore, the *left3words* model was chosen to be used in the application. (It still features a performance correctness score of 96.97% which is not that different to the bidirectional model's performance.)

With the input string broken down into a list of words, each word is used as input to the POS tagger, and the word's tag (e.g. *VB*, *NN*, etc.) is added to a new list containing all the tags (i.e. the index of a tag is the index of its corresponding word in the other list). Both lists will be used simultaneously in other processes.

5.3.3 Slot-filling Matching

The next stage is to extract the important information from the input. In many applications, particularly in games, voice commands are usually given in the imperative form (such as, "do this", "attack that", etc.). It is assumed that the subject of the command (that is, the entity about whom the statement is made about [44]) is always the player (or the character that the player is controlling).

Generally, commands will have an imperative verb (e.g. "attack", "grab", "look"), optionally followed by the target entity for which the imperative verb should be applied towards, as well as an optional context entity for which the imperative verb should be applied with.

As an example, take the command, "hit the rock with a spoon". Here, the word 'hit' is the imperative action and it is applied to the 'rock'. The 'hit' action is applied using the 'spoon' as the context. The target and context help to provide a clearer intent; the user could have just said, "hit" without mentioning the target or context, and while the action to be performed is clear, it is ambiguous as to what should be hit and with what.

For this stage, slot-filling structures (SFS) are used to describe the grammar of the voice commands. The SFS grammar chosen is:

ACTION TARGET WITH CONTEXT

Here, parts of the input will be mapped to a *block* in the grammar. The WITH block is for words and phrases which signify that the following words represent the context to perform the action (e.g. “...with...”, “...using...”, “...with the power of the...”, and so on). Figure 12 shows how parts of an example command are matched to the grammar.

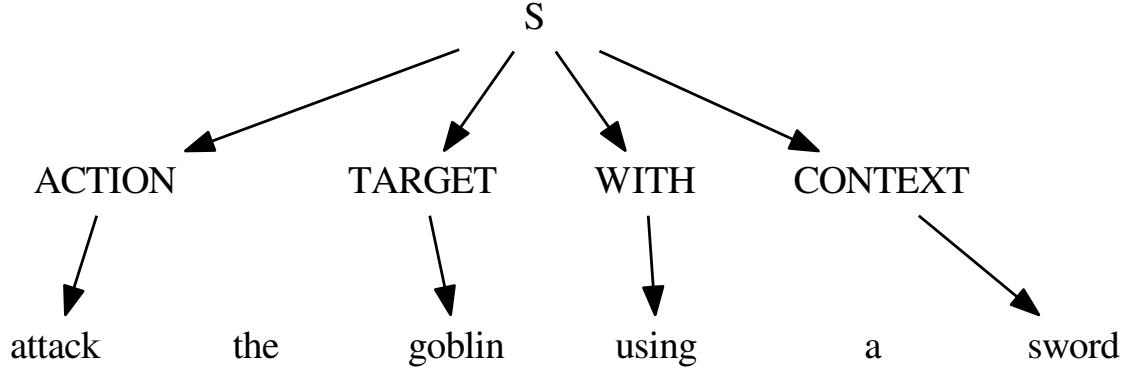


Figure 12: The slot-filling structure, {ACTION TARGET WITH CONTEXT} being applied to an example sentence, “attack the goblin using a sword”. Generated using Graphviz.

Not all the slots have to be filled; the TARGET, WITH and CONTEXT blocks are completely optional to fill, but they provide more information about the user’s intent. Should this information be left out, the system will fall back on the defaults for each one.

While this covers a large variety of commands that the user can give, such as, “attack the troll with a sword”, “pick up the knife”, or “look around the room using binoculars”, it does not cover commands that have a completely different structure. For example, the phrases, “use a potion to heal” and “with a key open the door” are perfectly valid commands that a user can give, but the matching blocks are in completely different order (i.e. the CONTEXT comes before the ACTION), so would not be matched with the above SFS.

An alternative SFS is applied when appropriate. The grammar is as follows:

WITH CONTEXT ACTION TARGET

Using this slot-filling grammar, the aforementioned example commands (“use a potion to heal” and “with a key open the door”) would be accepted. Figure 13 shows an example of the matching.

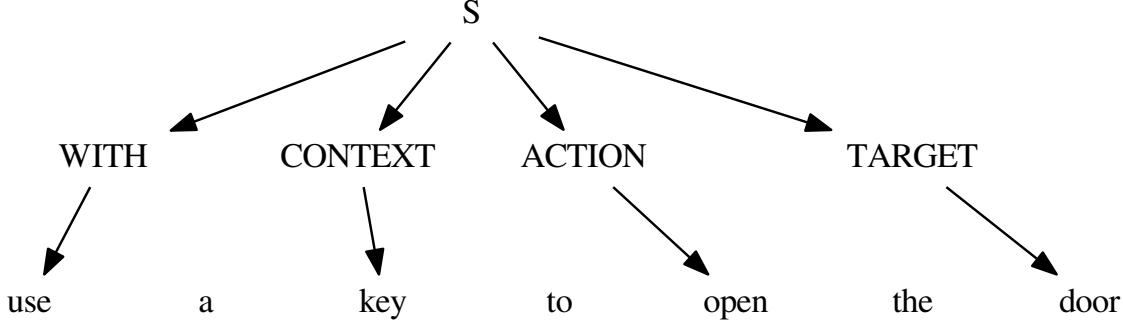


Figure 13: The slot-filling structure, {WITH CONTEXT ACTION TARGET} being applied to an example sentence, “use a key to open the door”.

During the processing, the correct SFS will be chosen based on the structure of the phrase (i.e. by determining the relative positions of each piece of information).

Using lists of possible actions, targets and contexts created by the developer (specified as strings), the slot-filling process will compare the words in the input to the words in the lists in order to find a match.

5.3.4 Semantic Similarity

During the slot-filling matching stage, if a word does not match any entry in the possible list of actions/targets/context character-by-character, then semantic similarity methods are used to determine whether the two words are similar. For example, if the possible list of actions are, ‘attack’, ‘heal’ or ‘run’, and the player’s input is, “hit”, then ‘hit’ is similar to the ‘attack’ action, so the ‘attack’ action would be chosen.

As mentioned in Section 4.4.5, there are many similarity methods that have been proposed. The first method to be used in the application is the *Wu and Palmer* method (WUP) described previously. Another method that is used is the *Lin* method (LIN), that calculates the similarity of two words based on the information content of their least common subsumer [32]. Information content (IC) is defined as the negative log of the probability of a concept. The probability of a concept is based on the sense frequency of a word within its synset.

Lin used IC to define a similarity measure similar to the Wu and Palmer method, but uses IC instead of the depths of the words in the word tree instead. The formula for the similarity of words w_1 and w_2 is as follows:

$$sim_{lin} = \frac{2 * IC(LCS(w_1, w_2))}{IC(w_1) + IC(w_2)}$$

Another method used is the *Lesk* method. This method determines the similarity of two words based on the overlaps in their definitions [45]. For instance, the definition for cer-

tain senses of the words, '*ash*' and '*coal*' (to describe a burnt substance) have the words, 'combustible', 'burn' and 'solid' in their definitions. Therefore, it can be deduced that the words are semantically similar. In another paper, Banerjee and Pederson extended the Lesk method by using all the definitions of all the senses for both words in the calculation, as only using one sense definition did not give enough information [46].

As well as using the definitions of the synonyms, hypernyms and hyponyms of the word, this method also uses the definitions of the word's meronyms and holonyms. Meronyms are words that represent part of something else [47]. For example, 'face' would be a meronym of 'person'. Holonyms are words that represent the whole of several parts, so a 'body' would be a holonym for words such as 'arm', 'leg' and 'torso' [47]. Using these definitions makes the Lesk method more accurate than other methods, but also adds more processing time for each calculation.

The WordNet database is accessed and traversed using the MIT Java WordNet Interface (JWI), which is an open-source library that interfaces with a WordNet database [48]. A full English WordNet database is downloaded from Princeton's website, and the JWI API is loaded with it.

Implementations for most of the similarity methods are provided by another open-source Java library, *WS4J* [49]. While this library provides the methods for doing the calculations (such as the WUP, LIN and LESK formulas), the methods for the navigation and analysis of the WordNet database (e.g. getting hypernyms, synsets, glosses, etc.) must be implemented from scratch (see Section 6.5.6 of the Implementation for more details).

5.3.5 Context-Action Mapping

Once matches have been found from the slot-filling, they are mapped to the methods that should execute that intent. If the player only specifies an action with no context, then a default action should be executed (the context itself defaults to a "default" context type). If the player specifies a context (e.g. "...with a sword"), then a mapping is consulted to find the correct method that should be called. The different targets are handled within the method itself.

Table 1 shows an example of what a context-action 'table' mapping would look like. The table is column-indexed using the action and row-indexed using the context type. Note the *null* entries which indicate that the context is not compatible with the corresponding action (for example, it does not make sense from a player to "attack with a potion"), and the context is therefore ignored (i.e. the 'default' method will execute instead).

Table 1: An example of a developer-specified context-action mapping

	actions		
context	<i>attack</i>	<i>heal</i>	<i>open</i>
<i>default</i>	AttackDefault()	HealDefault()	OpenDefault()
<i>weapon</i>	AttackWithWeapon()	null	null
<i>potion</i>	null	HealWithPotion()	null
<i>key</i>	AttackWithKey()	null	OpenWithKey()

5.3.6 Intent Execution

Each non-null field in the context-action map contains a class that extends from an abstract class, **Action**. This abstract class contains a method called `run()`, which executes the in-game intent. This method is overridden by a derived class (such as `AttackDefault` and `HealWithPotion`). These `run()` methods act as wrappers for the methods that would execute in the game (or execute some other actions unrelated to the game, such as displaying the available commands, etc.).

5.4 System UML Class Diagram

Figure 14 shows a simplified UML class diagram focusing on the classes used for the voice recognition system, including the inheritance and associations of the classes. Although the `GameState` class is related to the game system, it helps to illustrate the overall system design. Only the important member fields and methods are shown for each class.

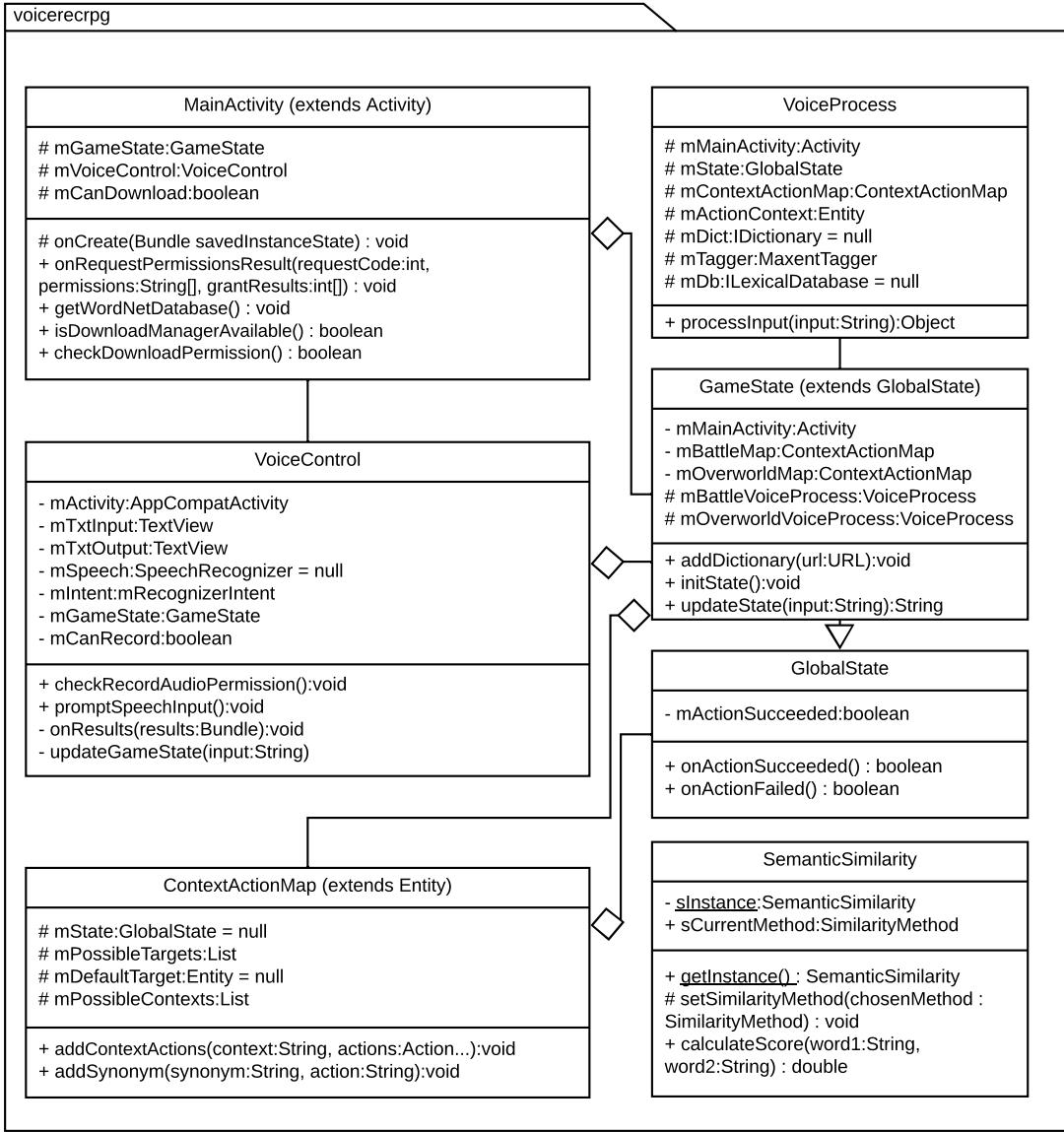


Figure 14: A simplified UML class diagram outlining the design of the voice recognition system. Diagram created using the online LucidChart tool. See Appendix B.1 for the full class diagram for the entire voice recognition system in the Android application.

5.4.1 MainActivity Class

The `MainActivity` class is the skeleton class which contains all the systems. It acts as the interface between the I/O and the rest of the application. It contains instances of the `VoiceControl` class, which handles the I/O, and the `GameState`, which contains the game environment.

When the application starts up, an instance of `MainActivity` is created, and searches the public directory of the phone's external storage (not necessarily an external device like an SD card, but the directory which the OS has defined the external storage to be). It searches for a WordNet database that is downloaded beforehand, and if it is not found, it will download/re-download the database again. It will also copy a model used for POS tagging to the external storage from the Android application package (APK) archive.

5.4.2 VoiceControl Class

The `VoiceControl` class is used to handle the input and output of the application, such as the microphone input and the text display. It contains an instance of the `SpeechRecognizer` class, which is provided by an Android SDK and uses Google's Speech-to-Text API, as well as two instances of the `TextView` class: one which is used to display the text output of the game, and the other to display the speech-to-text result at the top of the screen.

When the `SpeechRecognizer` instance performs a speech-to-text translation, the final result is forwarded to the `GameState` instance where it is processed by the voice recognition system using NLP, and used to update the state of the current game session in `GameState`.

5.4.3 VoiceProcess Class

The `VoiceProcess` class is responsible for processing the string input containing the player's utterance that is received from the `VoiceControl` class. It contains the majority of the NLP processing.

5.4.4 GlobalState Class

The `GlobalState` class is a generic abstract class that can be derived from in order to pass around environment objects between method calls easily. For example, if a method wishes to access an instance of an `Inventory` class (to access the game's items), then a reference to the inventory can be held in a class derived from `GlobalState` (e.g. `GameState`).

5.4.5 ContextActionMap Class

The `ContextActionMap` class contains the mapping of contexts to actions. It is an abstract class that is overloaded by the developer to add mappings of different ways of performing actions based on the given context (e.g. "attack with a sword" versus "attack with a hammer").

5.4.6 SemanticSimilarity Class

This class contains the semantic similarity engine used for calculating the semantic similarity of two words as a confidence value between 0.0 and 1.0, based on information in the WordNet database. It uses a singleton design pattern so the class does not have to be

instantiated every time a calculation needs to be performed with a different method, nor does a reference to an instance need to be passed around everywhere it is used. It also allows us to use the same static instance across different Android Activity contexts (and therefore allow the engine's parameters to be changed in a *Settings* screen).

5.5 System Features

In order to make the voice recognition system as flexible as possible and able to infer the meanings of a wide variety for utterances, several mechanisms have been added to the system. These mechanisms help to reduce the amount of work required by the developer, and some are present in order to compensate for the shortcomings of the system (see the Evaluation section).

5.5.1 Synonym Mapping

The developer can map synonyms to actions if the words are not considered semantically similar by the system. For instance, if the word, “regenerate” has a low semantic similarity score when compared to the “heal” action, then the developer can create a map of “regenerate” → “heal” so the system performs the correct mapping.

5.5.2 Ignoring Incorrect Matches

Semantic similarity matches can be ignored if they are incorrectly mapped by the system, as a juxtaposition to the above feature.

5.5.3 Confirmation and Suggestions

If the user gives a command that has a confidence value just below the threshold, then the phrase is considered ambiguous and, using the possible candidate matches, suggestions are given to the user to confirm their intent.

For example, if the user’s intent is to pick up a knife and they say, “pick up the utensil”, where the similarity score between “utensil” and “knife” is just below the threshold, the system will output, “Did you mean, ‘pick up the knife? (yes/no)’”, and waits for the user to respond. If they say no, then the system may suggest another alternative using the next ambiguous candidate (e.g, “Did you mean, ‘pick up the plate?’”, etc.).

5.5.4 Multiple Commands

The user has the ability to execute a chain of several commands one after the other using just one utterance. For example, if the user says, “attack the enemy and then heal using a potion”, this will execute two separate commands.

5.5.5 Multiple Targets

A user may want to execute a command that involves several targets. In this case, the intent is divided into several actions instead, with each action focused on one of the targets.

For instance, if the user says, “attack the troll and the goblin”, this will be split into two actions: one to attack the troll, followed by an attack on the goblin.

5.5.6 Sentence Mapping

The core mechanism of the system is the slot-filling structure described previously. This covers most of the forms of possible imperative commands that a user may utter to the system to execute an intent. However, one disadvantage of this system is that it will not detect phrases that do not map to the slot-filling grammar given (e.g. {ACTION TARGET WITH CONTEXT}), despite being a valid phrase. Such utterances are usually phrased as a question, such as, “What actions can I do?”, or, “What items are in my bag?”, etc. (see Section 7.5 for a detailed evaluation of the chosen slot-filling grammar). The system should still be able to detect these types of commands but prevent the developer from hard-coding each possible variation of the intents.

The proposed solution is to perform sentence matching on these phrases. The developer specifies an intent followed by several example sentences that a user would say to execute the intent. For example, if the intent is to show the player the items they possess, example sentences would include: “What is in my bag?”, “What items do I have?”, “What are the contents of my inventory?”, and so on. When parsing a new utterance from the player (e.g. “What items are in my bag?”), if the phrase fails the slot-filling stage, then the similarity of the sentence compared to the example sentences is calculated.

The calculation involves using a slightly altered cosine similarity method. Based on the formula for the cosine of an angle between two vectors, the cosine similarity method measures uses the cosine of the angle as a measure of similarity, using a vector space model [50] (the vector of a sentence is represented by its word frequencies). The similarity between two vectors, a and b is:

$$\text{similarity} = \cos \theta = \frac{a \cdot b}{|a||b|}$$

Take two examples: “the apples are in the tree” and “the tree contains apples”. The vectors for these sentences would be:

$$\begin{aligned} a &= \{\text{the}:2, \text{apples}:1, \text{are}:1, \text{in}:1, \text{the}:2, \text{tree}:1, \text{contains}:0\} \\ b &= \{\text{the}:1, \text{apples}:1, \text{are}:0, \text{in}:0, \text{the}:1, \text{tree}:1, \text{contains}:1\} \end{aligned}$$

The cosine similarity of these sentences is therefore:

$$\frac{a \cdot b}{|a||b|} = \frac{6}{2\sqrt{3} \cdot \sqrt{5}} = 0.77$$

The definition for cosine similarity has been altered to form a “soft cosine measure” [51]. This uses the semantic similarity between the two words (when the dot product is calculated) as a weighting: if a word only appears in one sentence and not the other, then the scores of its semantic similarity to the words in the other sentence are added to the

numerator.

While the cosine similarity is not ideal for all situations (e.g. it ignores word order so may match two sentences which are not similar at all), it is easy to implement for the sentence matching mechanism, which is supplementary to the voice recognition system.

5.6 Game Design Overview

The voice recognition system will be applied to a text-based role-playing game that has several different modes of play, in order to demonstrate the flexibility of the system. The first mode is similar to the Zork-style text-adventure gameplay, mentioned previously. This involves the player interacting with objects in the environment to solve puzzles while collecting items, in order to progress. This mode is named the *Overworld* mode.

The second mode of gameplay is similar to classic turn-based RPGs such as Pokémon or Final Fantasy, where the player battles an enemy, and the opposing sides alternate between taking turns to either attack, defend, heal, or perform some other action. This mode is called the *Battle* mode.

5.6.1 Game Flow

The focus of the project is not on the game itself, but rather on the power and flexibility of the voice recognition system to add voice commands to the game. Therefore, only a simple game will be created with a few rooms, with different rooms focusing on different game modes (*Overworld* or *Battle*).

For the first room, the player is placed in a room where they are presented with a puzzle as part of the *Overworld* mode. The room features the following description:

- “You are in a room with a locked door in front of you.”
- “There is a glass table in the middle of the room.”
- “There is a knife on the table.”
- “A painting of a tree hangs by a string on the left wall.”

The player must obtain a key hidden behind the painting by cutting it down using the knife (or any other sharp item they may possess such as a sword), and then use the key to open the door.

During the *Overworld* gameplay, the user is capable of performing the following actions:

[look, show, grab, open, cut, break]

Example commands include, “show my inventory”, “look around the room”, “grab the knife”, “cut the painting”, and so on.

Each room in the game which features the *Overworld* gameplay is represented by a finite state automaton. When the player performs certain actions (e.g. cutting the painting down), the state of the room will progress until the player has performed the correct actions in the correct order and is able to progress to the next room. Figure 15 show the state diagram for the first room.



Figure 15: The state diagram for one of the rooms in the game.

As part of the *Battle* mode, the player will be presented with a battle against an enemy (e.g. a *troll*). The following actions are available to the player:

[attack, heal, show, look]

Examples include: “attack the troll” “heal with a potion”, “look around”, and so on.

5.7 Game Design UML Class Diagram

Figure 16 on page 34 shows a simplified UML class diagram for the game logic.

The **GameState** class (shown previously in the system UML class diagram, to aid explanations) contains two different instances of classes that override the voice recognition system’s **ContextActionMap** class (which contain the voice command mappings): one for the overworld gameplay and the other for the battling gameplay.

The **Entity** abstract class represents anything that can be considered a possible target or context for the slot-filling. Each entity is given a name and a context type (as a string) that is used to determine which context it maps to (e.g. “weapon”, “key”, etc.). See Table 1 above for examples of these contexts. Each entity also contains a description field which is an array of words that describe the entity (e.g. a knife would be described as “sharp”, “pointy”, etc.). These descriptions are optional but they aid the voice recognition system to help identify similar words (if the semantic similarity match fails).

Finally, **Room** is an abstract class for generating new rooms for the game, each with their own objects.

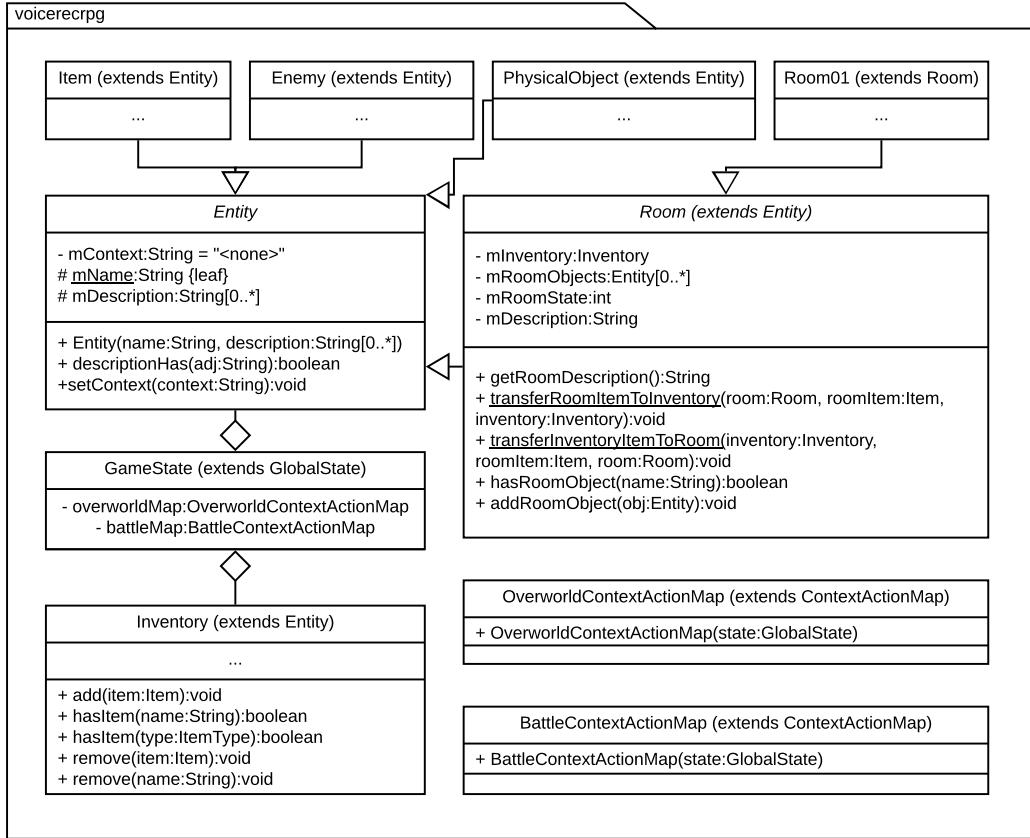


Figure 16: A simplified UML class diagram outlining the design of the gameplay mechanics. Diagram created using the online LucidChart tool. See Appendix B.2 for the full UML class diagram.

5.8 Room Generation

Instead of forcing the developer to manually instantiate objects within each room, a mechanism will be written that will allow a developer to generate a new room from just a description of it, using natural language processing techniques.

A command-line program has been created which takes as input a text file containing a description of a room and generates a source file for a `Room` class, containing the objects from the description, and also handles object relations (e.g. whether a lamp is on the table or on the floor). This allows the developer to easily create new rooms for the game without having to manually add the objects to a room. See Section 6.8 for the implementation.

6 Implementation

This section outlines the steps taken to implement the system and the game, mostly in chronological order where appropriate. For more details on the source code, please see the corresponding GitHub repository containing this project's work.¹

6.1 System Architecture

Figure 17 shows a UML component diagram for the system architecture of the Android application. The front-end consists of three **Activity** classes. An **Activity** is the environment where the user interacts with the application.

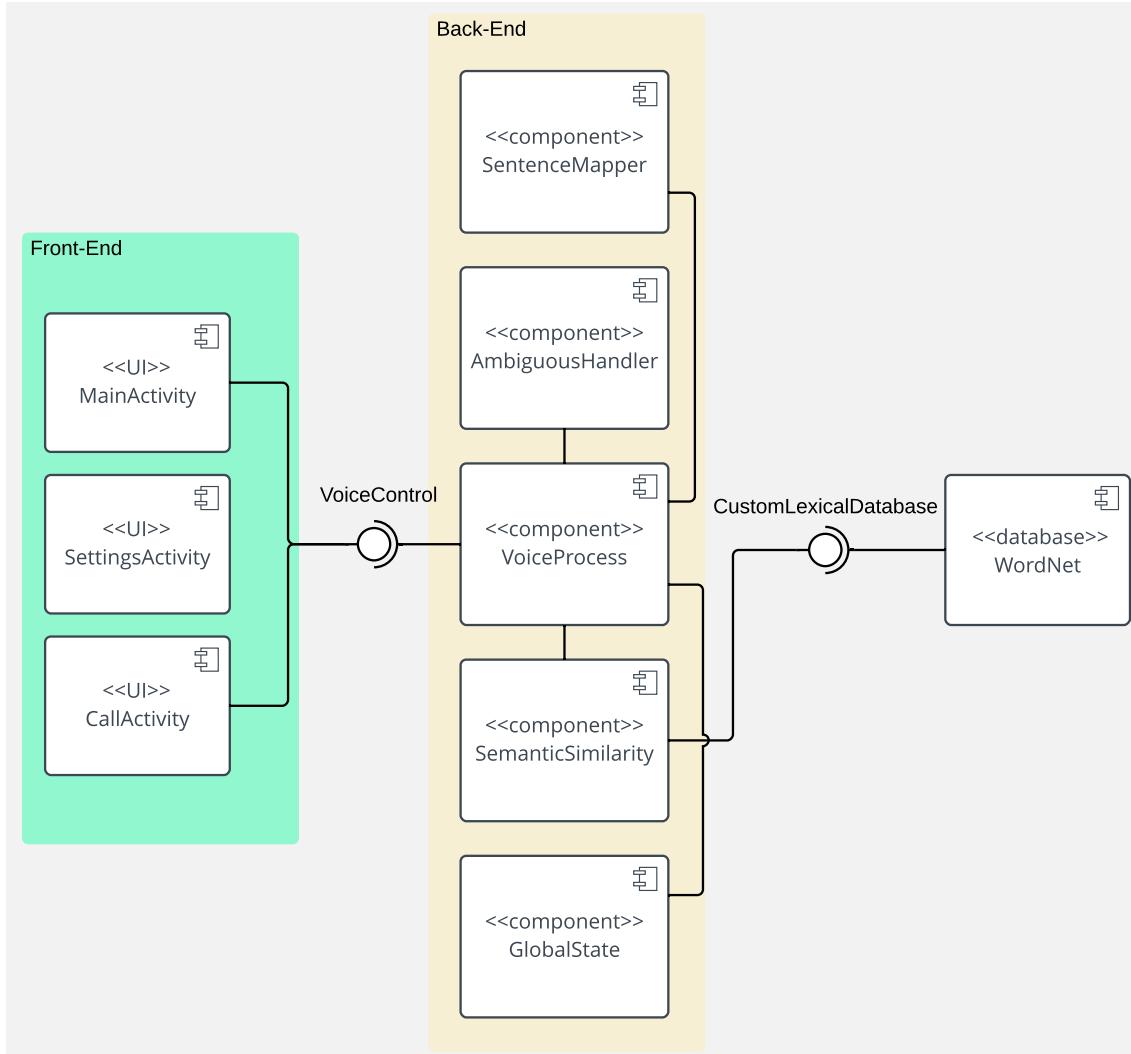


Figure 17: A UML component diagram showing the components of the application.

¹<https://github.com/BaronKhan/VoiceRecognitionRPG>

The `MainActivity` is where the role-playing game interactions will take place, while the `CallActivity` contains a demo of the system applied to a video conferencing example (see Section 6.9). There is also a `SettingsActivity` for setting the different parameters of the application (see Section 6.10).

The back-end accesses a WordNet database stored on the phone to perform similarity calculations through a `CustomLexicalDatabase` interface (see Section 6.5.6 for more).

6.2 Voice Recognition Interface

The first step in the implementation involved creating an Android project with a skeleton I/O to take as input the microphone audio, converts the speech to text, sends it to the voice recognition system for processing (initially just an empty class that pipes the input to output), before finally displaying the text on the screen.

The interface is within the `MainActivity` class. When the activity is created, it also downloads a tar ball containing a WordNet database and unzips it on the device's local storage. While it is possible to package the database with the model, the code was originally written to download the database and, due to time constraints, could not be changed.

The `MainActivity` class contains a `VoiceControl` class that handles the speech-to-text (STT) I/O, and implements a `RecognizerIntent`. The `VoiceControl` object contains an instance of the `SpeechRecognizer` class that has a `startListening()` method, which invokes an `onResults()` method when finished. This method has a `results` parameter of type `Bundle`, containing the string result of the speech recognition. The string is then processed by an instance of the `VoiceProcess` class (described below).

6.3 Graphical User Interface

Figure 18 shows a snapshot of the UI for the game logic with the following labels:

1. A `TextView` element containing the player's input for the game as a string (i.e. it is the output of the speech-to-text API).
2. A `TextView` containing the cumulative output of the game as each input is processed.
3. Pressing this `Button` element activates the microphone and starts a voice intent.
4. A timer used for debugging the time taken to process an intent by the system.
5. This drop-down list contains various options (such as an option to open the settings menu).

All UI elements (such as `TextViews`, `Buttons`, etc.) are accessed in the `MainActivity` class, and are separate from the voice recognition system.

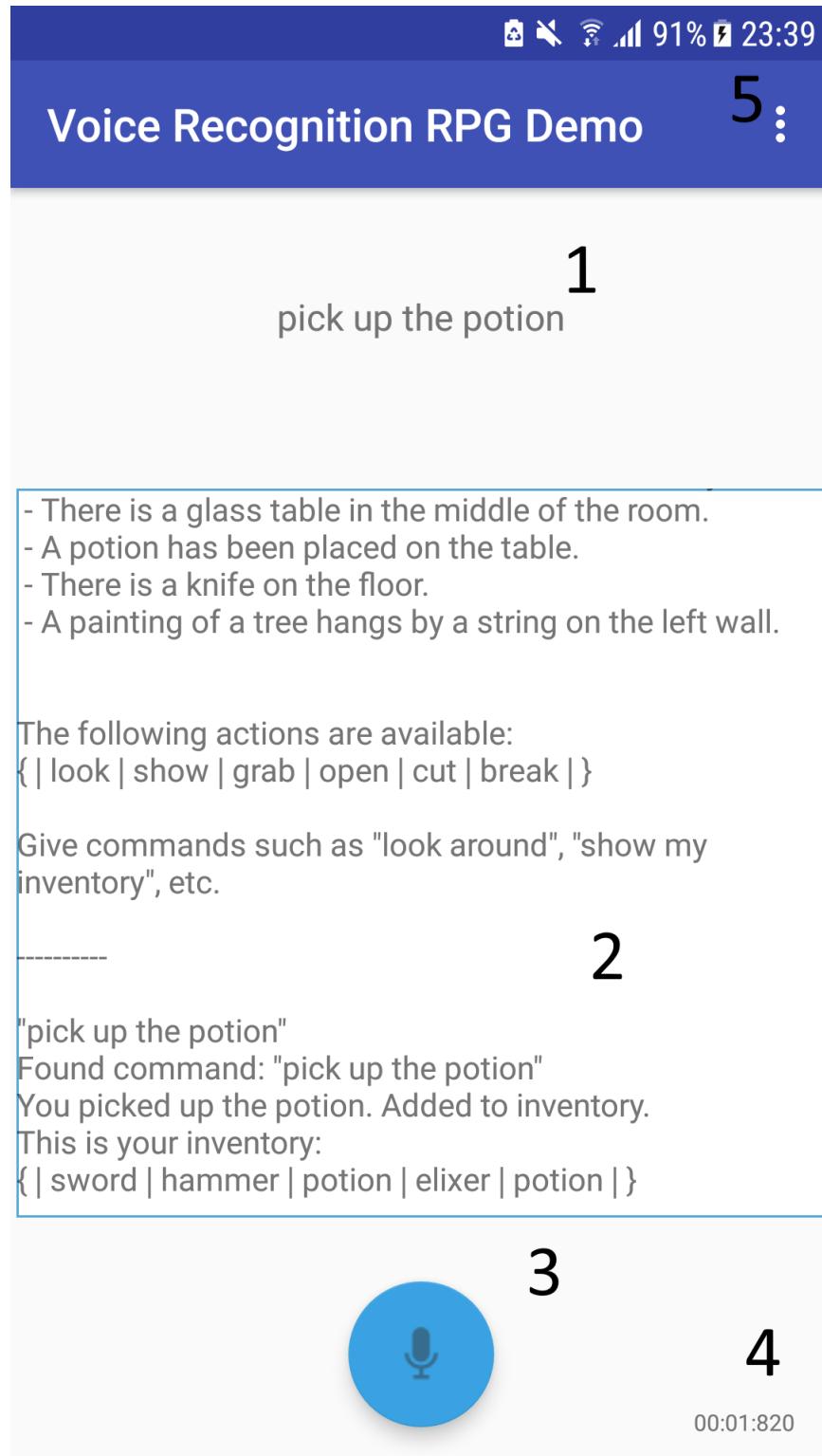


Figure 18: A snapshot of the user interface for the text-based role-playing game.

6.4 User Flow Diagram

Figure 19 shows a wireframe for the different screen transitions that occur as the user interacts with the application.



Figure 19: A wireframe showing the majority of transitions between activities as the user interacts with the application. The tails of the red arrows indicate where the user touches the screen (unless otherwise stated). Diagram created using MockFlow.

6.5 Voice Recognition System Implementation

This section outlines the general process of understanding the user's intent and mapping it to an action in the application. It also describes each component in the voice recognition system in more detail.

6.5.1 Data Flow Diagram

Figure 20 shows how data flows in the voice recognition system, starting from the user's audio input, and ending with the output sent to the UI. Note that not all data paths will be followed during a voice intent. The critical path is highlighted in red.

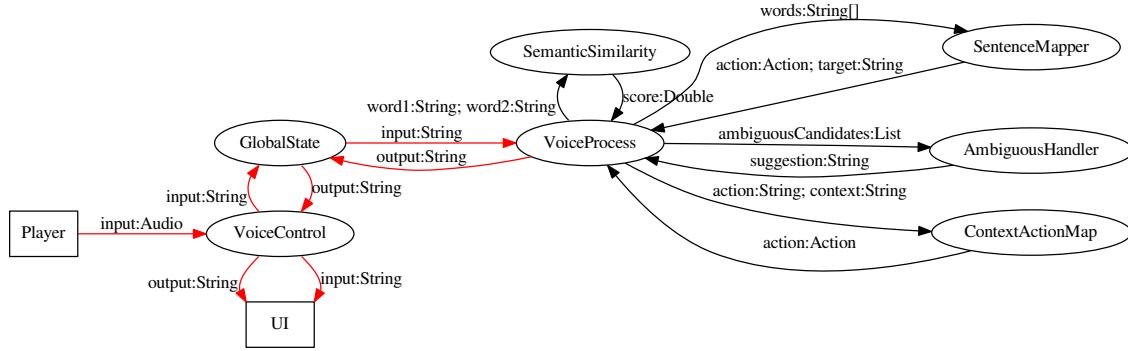


Figure 20: A data flow diagram showing how the player's input (starting from 'Player') is processed by the different components until a response is sent to the 'UI'. The critical path is highlighted in red. Diagram generated automatically using Graphviz.

6.5.2 GlobalState

The `GlobalState` class is an abstract class that sends the string input to a `VoiceProcess` instance via the `updateState` method. Derived classes (e.g. `GameState` for a game, `KitchenState` for a cooking application, etc.) contain the application-specific logic and environment that can be accessed by the action methods.

In some scenarios, a follow-up action may only take place if the user's action was successful. For instance, in a generic two-player game, the second player may only be able to make their move once the first player has made a successful move. A more relevant example is in a turn-based role-playing game, where the player is fighting an enemy. The enemy can make their move once the player has made their move.

To support this scenario, there is a Boolean field called `mActionSucceeded`, which is set within the action methods that are executed to tell the system whether the action has succeeded, and the developer can choose what happens next.

6.5.3 VoiceProcess

`VoiceProcess.java` is the main class that performs all the processing on the text string input (of the user's utterance). The main method is the `processInput()` method that takes the string input and returns the response. Figure 21 shows a flow chart for this method.

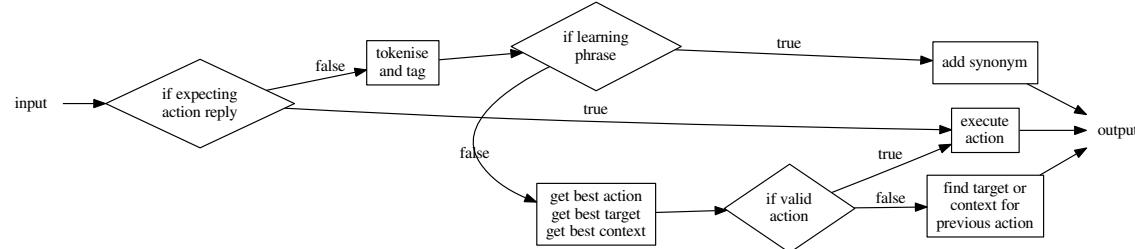


Figure 21: A flow chart outlining the `processInput()` method in the `VoiceProcess` class. Note that any minor steps are not shown here.

See Appendix D.3 for the pseudocode of the `processInput()` method. For the full Java method, refer to the source code in the project's repository.

When tokenising and tagging the input, a list of `words` and a list of `tags` are created. The index of a word in `words` corresponds to the index of its tag in `tags`.

For the `getBestAction()`, `getBestTarget()` and `getBestContext()` methods, the words in the input are checked sequentially to determine the best match. For instance, in the `getBestAction()` method, the candidate actions are first chosen from the input. For each candidate action, a score out of 1.0 is calculated against each possible action defined in the `ContextActionMap`: if the words match character-by-character or the candidate is a synonym for one of the actions, then the score is set to 1.0. Otherwise, the score is set to the semantic similarity score of the two words. The `getBestTarget()` and `getBestContext()` methods work similarly.

The candidates are chosen using POS tagging, to filter out words with no meaning such as “the”, “it”, etc. Words in the input are candidate actions if they are tagged as either a verb, adverb, or noun (occasionally, the POS tagger incorrectly tags some imperative verbs as nouns).

6.5.4 AmbiguousHandler

Candidate actions, targets and contexts may be marked as ambiguous/uncertain if their confidence score is just below a threshold value. If the best action/target/candidate has a confidence value lower than this threshold, then all the ambiguous candidates are queried to the user using the `AmbiguousHandler` class. This class will generate a suggestion and display it to the user, starting with the ambiguous candidate which gave the highest confi-

dence. If the user disagrees with the suggestion (i.e, they say “no” or something similar), then the next suggestion is given, until no more suggestions are available.

Figure 22 shows an example of a suggestion displayed to the user after an ambiguous intent. See the `AmbiguousHandler.java` file in the project’s repository for the full code.

6.5.5 SemanticSimilarity

The semantic similarity engine is implemented as a singleton design pattern with a private constructor and a private instance of itself within the class. The main method of the class is the `calculateScore()` method which returns the semantic similarity score for two words.

The engine uses up to two semantic similarity methods, and if two are chosen, then the score returned is the average score from both methods. As well as the Wu and Palmer (WUP), Lin and Lesk methods mentioned previously, below is a list of methods that can also be chosen:

- Fast Lesk (FASTLESK) - a faster implementation of the Lesk method, that removes string formatting of definitions and ignores meronym and holonym definitions of the words.
- Path-based (PATH) - a simpler version of the WUP method which uses the reciprocal of the shortest path between the words in the WordNet database as the similarity score [32].
- Leacock and Chodorow (LCH) - another path-based measure which uses the negative log of the shortest distance between the two words [32].
- Resnik (RES) - similar to the LIN method, but does not normalise the information content of the words’ least common subsumer [32].
- Cosine (COS) - similar to the Lesk method, but uses the cosine similarity of their definitions (as opposed to just overlaps in the definitions).

Most of the above similarity method implementations are provided by the WS4J library, except the FASTLESK and COS methods, which are unique to this project and are implemented manually in the code (though the FASTLESK method is altered from WS4J’s LESK method which was decompiled using Android Studio’s built-in Fernflower decompiler [52]). By default, only the WUP method is used by the system, though other methods can be chosen in the application’s settings.

6.5.6 CustomLexicalDatabase

Despite the WS4J library providing implementations of the algorithms for most of the semantic similarity methods used in this project, they require an interface to be implemented in order to access a WordNet database, called `ILexicalDatabase`.

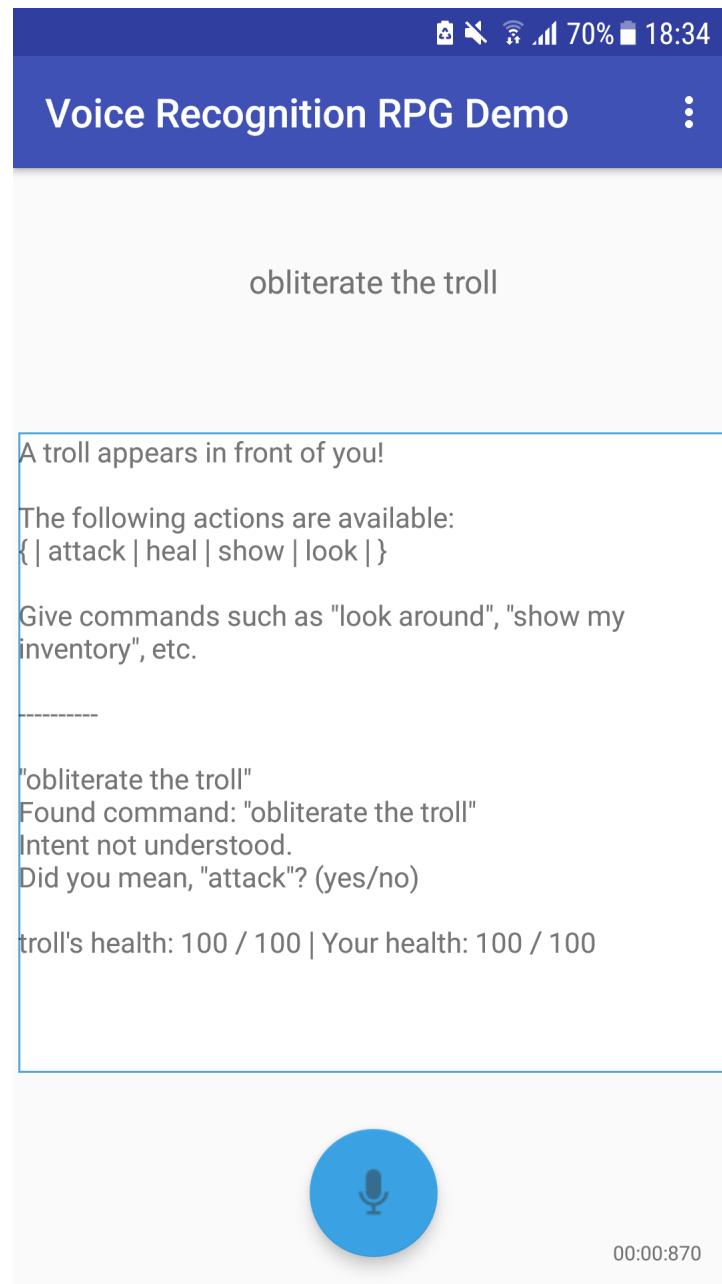


Figure 22: A snapshot of the ambiguous user intent, “obliterate the troll”. Here, the word, ‘obliterate’ is most semantically similar to the ‘attack’ action, but not enough to pass the threshold, so it is queried to the user.

The WS4J library is packaged with its own implementation of `ILexicalDatabase` interface, called `NictWordNet` (which accesses a Japanese WordNet database similar to the Princeton WordNet [53]). Unfortunately, due to its use of SQLite methods to access the database, the implementation is not compatible with the Android SDK, and could not be used.

Therefore, a new implementation of the `ILexicalDatabase` interface was created that works with the Android SDK, called `CustomLexicalDatabase.java`. It interfaces with a local Princeton WordNet database using the MIT JWI API mentioned previously. All methods were implemented using the JWI API except the `findSynsetBySynset()` and `conceptToString()` methods (these methods are not required for the above similarity methods).

Since this implementation using the JWI library may be useful for others in the future, a GitHub gist has been published online (renamed to `JwiLexicalDatabase.java`) that other developers may freely use ².

6.5.7 ContextActionMap

The `ContextActionMap` class contains the table of mappings of contexts to intents. The developer uses this abstract class to create a derived class that defines the mapping in a ‘table’ format (with contexts on the left and actions along the top) in the constructor. The table itself is implemented as a ‘map of maps’. The `setActionsList()` method specifies the possible actions, while the `addContextActions()` method describes which Java method should be invoked for each action based on the context. The entries in the ‘table’ are either `null` or classes that inherit from the `Action` class (see Section 6.5.9).

There is also an `addDefaultContextActions()` method, which defines the actions to be called if no context is given in the input. A `null` entry means that the action is invalid for the given context (e.g. you cannot attack with a potion). In this case, the default context action is called instead, and if the default action is `null`, then nothing happens.

Listing 1 shows an example of the table in a derived class, `GameContextActionMap`. Note that this example was only used for testing and evaluation of the system, and is not one of the tables used in the RPG demo.

Listing 1: Example of a `ContextActionMap` ‘table’

```

2 setActionList(      "attack",      "heal",      "move");
3 addDefaultContextActions( new Attack(), new Heal(), new Move());
4 addContextActions("weapon", new AtkWeapon(), null, null);
4 addContextActions("potion", null,      new HealPotion(), null);

```

In order to make it easier for the developer to create this ‘table’, a Python script was written called `generateTable.py` which takes as input a CSV file of a table, and gen-

²<https://gist.github.com/BaronKhan/5679157d1a8572debc8c75af1c875c16>

erates a Java source file similar to Listing 1. Table 2 shows the CSV table (where cells are separated by commas) that will generate the Java source file in Listing 1 using the command,

```
python generateTable.py game-map.csv GameContextActionMap
```

See the source code in the project’s repository for the full Python script.

Table 2: game-map.csv

	<i>attack</i>	<i>heal</i>	<i>move</i>
<i>default</i>	Attack	Heal	Move
<i>weapon</i>	AtkWeapon		
<i>potion</i>		HealPotion	

In the `ContextActionMap` class, there are two member fields that the developer must build in their application: the `mPossibleTargets` and `mPossibleContexts` lists. These are lists of objects that extend the abstract `Entity` class. Only entities in these lists are used in the target/context matching (when finding the best match).

6.5.8 Entity

`Entity` is an abstract class in the system which represents anything that could be a possible target or context. Each derived class of `Entity` must set the following fields: `mContext`, which is the context type of the entity (i.e. the left-hand side value in a `ContextActionMap` described in Section 6.5.7); `mName`, which is the name of the entity (e.g. “sword”, “potion”, etc.), and, optionally, `mDescription`, which is a string array of adjectives that are used to describe the entity. `mDescription` is optional, but it assists the `processIntent` method in the `VoiceProcess` class in identifying the best matches for targets and contexts.

Once a match for a context is found, its `mContext` string value is used to index the `ContextActionMap` table, along with the chosen action.

6.5.9 Action

The `Action` class is an abstract class that is used to wrap the methods that should be executed in the `ContextActionMap`. The `execute()` method is overridden by the derived class and executed when the action is chosen in the `VoiceProcess` class. The method takes as input the `GlobalState` object (so the method can update its state if needed) and the current target (if it exists).

Some methods may wish to ask the user for confirmation before executing (e.g. “Are you sure you want to do this?”). In this case, the `execute()` method would set the `mWantsReply` member field to `true`. On the next user intent, the input is processed by the `processReply()` method (i.e. whether the user said “yes”/“no”, etc.).

6.5.10 MultipleCommandProcess

Optionally, the developer can enable the ability to execute multiple commands with a single utterance. For instance, the player can say, “attack the troll with a sword and then heal with a potion”; this would execute the attacking and healing as separate intents. The developer can wrap their `VoiceProcess` class inside a `MultipleCommandProcess` class.

The `MultipleCommandProcess` class has a `splitInput()` method which splits the user’s input using the words, “and” and “then” as delimiters. It stores each result in a queue, and then processes each intent using the `executeCommand()` method, which takes the queue as input. If an intent is interrupted and more user input is required before continuing (e.g. if one of the intents is ambiguous, or an action requires a reply), then the remaining intents are stored in a second queue, and are processed once the current intent is resolved. See Appendix A.3 for an example of the aforementioned scenario. Also see the `executeCommand()` method in `MultipleCommandProcess.java`.

6.5.11 SentenceMapper

The `SentenceMapper` is used to match sentences provided by the developer which do not fit the default slot-filling grammars (e.g. question-based intents). Each instance of `ContextActionMap` holds a reference to an instance of a `SentenceMapper`. Invoking the `addSentenceMatch()` method in the `ContextActionMap` class will add a map of several example utterances to an action and target.

The sentence mapper is run at the very end of the `processInput()` method of the `VoiceControl` class, if the slot-filling matching failed.

The `checkSentenceMatch()` method of the `SentenceMapper` class goes through each group of sentences, and calculates the average similarity of each group of sentences compared to the user’s utterance, and matches with the highest value (if it is greater than a threshold value).

The similarity of each sentence with the utterance is based on the cosine similarity of the sentences (refer to Section 5.5.6). First, a vector of each sentence is created (which is a map containing the words as keys and the values being the word’s frequency) before the formula for cosine similarity is applied. Refer to the `calculateCosScore()` method in the class.

A disadvantage to running the sentence-matching at the very end is that the input may accidentally match to the slot-filling grammar. For example, if the user says, “What can I cut?”, and “cut” is an action in the game, then it may map to the cut action, and will not be processed by the `SentenceMapper`.

6.5.12 Synonym Mapping Implementation

Synonyms are used to aid the `VoiceProcess` system to match more words that the semantic similarity engine does not match automatically (due to a low similarity score). For example, the word, ‘obliterate’ has a similarity score with ‘attack’ just below the threshold; they are similar but not enough. Therefore, either the player or developer can add this mapping such that ‘obliterate’ maps to ‘attack’.

The following are the possible ways of adding a synonym mapping:

- Add the synonym in the code (e.g. `addSynonym("obliterate", "attack")`).
- The synonym is added when the same ambiguous intent is spoken twice (e.g. when the system asks the user, “Did you mean, ‘attack’?” twice, and the user replies with “yes” twice).
- The player defines the mapping as an intent (i.e. they say, “obliterate means attack”).

In some cases, a word could become a synonym for multiple words. Take the following scenario:

- The player is in a room with just a knife. The player utters, “pick up the utensil”, and creates a mapping from “utensil” to “knife”.
- The player progresses to a room with just a fork. They create a mapping from “utensil” to “fork”.
- The player is now in a room with both a knife and a fork, and says, “pick up the utensil”.

In this scenario, the system will use the `AmbiguousHandler` to ask the user whether they meant the knife or the fork. See Appendix A.4 for the room that demonstrates this scenario. There is also an example of a room with multiple utensils, and when several synonym mappings are made for the word, ”utensil”, all suggestions are displayed at the same time (if those objects are present).

6.6 Role-Playing Game Implementation

Simple rooms were created to demonstrate the system working with a typical text-based role-playing game. The first room of the game focuses on the *Overworld* gameplay, and has a door which leads to a battle with an enemy (transitioning to the *Battle* mode).

6.6.1 GameState

`GameState` is derived from `GlobalState` and acts as the ‘environment’ for the game; it contains all the game’s objects and handles the current state.

Since there are two distinct modes of gameplay (the *Battle* gameplay and the *Overworld* gameplay), there is a `ContextActionMap` for each one in the `GameState` class. It also keeps track of various data such as the current game mode; the current room/enemy; and the player's health.

The `GameState` class also has an instance of the `Inventory` class. This is a list of items currently in the player's possession. The player can transfer items from their inventory to the room and vice versa.

6.6.2 Battle Mode

In the *Battle* mode, the player is able to attack the enemy with a weapon in their inventory. If no weapon is specified as the context in the intent, then a random weapon is selected from the inventory, and if no weapon is available, the player attacks the enemy with their bare hands. The player can also heal using any healing items they have (e.g. potions, elixers, etc.), as well as look at their inventory or the available actions. See the `BattleContextActionMap` class in the project's source code for the 'table' of action mappings for the *Battle* mode (including the helper features such as synonym maps and sentence matches).

Weapons can be classified as either "sharp" or "blunt" based on whether they are instantiated with a description containing one of those words (or if the weapon's name is semantically similar to either of the properties) — each type of weapon has a different multiplier for the damage dealt to the enemy.

Enemies are entities with their own health that the player fights. The enemy in the demo is a troll that can attack the enemy once the player has performed a successful action.

See Appendix A.1 for screenshots of examples of actions performed during the *Battle* mode. Also see Appendix B.3 for examples of commands that can be used during this mode.

6.6.3 Overworld Mode

In the *Overworld* mode, the player interacts with objects in order to solve puzzles and progress to the next room. Refer to the `OverworldContextActionMap` class in the project's repository for the table of context-action mappings in this mode.

Objects in the room can inherit from the `PhysicalObject` class, and can have properties such as whether they are breakable, whether they can be cut or scratched, etc. For example, if the player breaks an object, the `onBroken()` method is invoked. By default, the object just breaks, but the user can override this method in the object to add special functions (e.g. if breaking an object reveals a key, etc.).

See Appendix A.2 for screenshots of the first room in the game, and the actions required to solve the puzzle in the room. Also see Appendix B.4 for examples of commands that can be used during the *Overworld* mode.

6.7 Object Properties with WordNet

Using the semantic similarity engine, a user can select an item to use/interact with based on its characteristics. For example, if the user possesses a sword in their inventory and says, “attack with something sharp”, the sword will be selected because the adjective, “sharp” has high similarity to “sword”. Similarly, the word, “heavy” has high similarity to “hammer”, so if the player says, “hit with something heavy”, the hammer will be used.

In order to aid this system, entities can also be given a list of words (`mDescriptionList`) which describe the object. This aids the voice recognition system in identifying objects with characteristics that the user has specified.

Note that these properties associated to objects are separate from the properties in `PhysicalObject`, which are related to interactions specific in the *Overworld* gameplay (e.g. whether an object is breakable), although these physical properties can be automatically assigned using the `autoAssignProperties()` method, which uses the semantic similarity engine to compare the name of the object to words associated with the property. This allows the user to create new objects without having to manually assign each property. See the `Chair` or `GlassTable` objects in the source code for examples of this.

If the automatic assignment does not correctly assign a property (there is a program in the repository that determines which properties are automatically assigned to an object³), the developer can manually assign them. For example, we do not want the player to break the painting in the room, even though the method will automatically assign the painting to be breakable, so it is set manually to not be broken by the player.

6.8 Room Generation Implementation

The use of room generation in this project aims to decrease the workload of the developer and designer as much as possible, using semantic similarity techniques discussed above.

6.8.1 Rooms

Rooms contain the objects that the player interacts with during the *Overworld* mode.

Each room is derived from the `Room` abstract class. Objects are added to the room when it is instantiated, and a description of the room is displayed when the player issues a command such as, “look around”. Listing 6.8.1 shows an example of the description that is output in the first room.

³<https://github.com/BaronKhan/VoiceRecognitionRPG/tree/master/archive/property-checker>

```

2   - You are in a room with a locked door in front of you.
   - There is a glass table in the middle of the room.
   - A potion has been placed on the table.
4   - There is a knife on the floor.
   - A painting of a tree hangs by a string on the left wall.

```

This description will change dynamically as the player adds, removes or interacts with objects in the room. For example, if the player breaks the glass table, then the third sentence involving the potion would change to, “A potion lays on the floor with the broken table.” Similarly, if the player cuts the painting down, then the last sentence would change to, “The painting that was on the wall now lies on the floor.”

6.8.2 Manual Room Generation

Initially, objects were added manually to the constructor of the room, and the description of the room was generated via a series of `if` and `else` statements in a `getRoomDescription()` method (based on the current state of the room when the method is called). However, this was too cumbersome for a developer/designer to use, as the `if` statements would become long and complicated, so the whole structure was changed. See the source code from a previous commit for an example of the first room in the game using this version of the room generation code (for comparison purposes)⁴.

6.8.3 Using BooleanSuppliers

For the second iteration of the room generation structure, objects are added to the room alongside the sentence in the description that mentions the object. In order to dynamically change these sentences as the state of the room changes, `BooleanSuppliers` are used. This is a functional interface that allows a developer to store a lambda expression (e.g. a Boolean expression) and evaluate it later on using the `getAsBoolean()` member method [54]. For example, instead of writing an `if` statement for whether a table exists in the room or not, we can just create a lambda expression in the constructor `((() -> getRoomObjectCount("table") > 0))`, and evaluate it later. This prevents the need to write a series of `if` and `else` statements when retrieving the room’s description.

The `Room` abstract class contains an `mDescriptionList` field. Each entry in this list represents a different sentence in the description and contains a triple of the following:

- `Pair<String, String>` — a pair of sentences: one for when a conditional is `true` and the other for when it is `false`.
- `BooleanSupplier` — a lambda expression that represents the conditional which chooses the sentence that should be displayed from the pair above.
- `BooleanSupplier` — a lambda expression that denotes the existence of the object tied to the sentence. If no object is tied to the sentence, then this is set to `null`.

⁴<https://github.com/BaronKhan/VoiceRecognitionRPG/blob/d18851f61a655457d9b7e9b24e0e79e43ba9657a/android/app/src/main/java/com/khan/baron/voicerecrpg/rooms/Room01.java>

In the `getDescription()` method (which is called to output the current description of the room to the player), we simply evaluate the `BooleanSuppliers` in all the entries in the `mDescriptionList` field, and adjust the description accordingly. See the project's source code for examples of rooms.

6.8.4 Room Generation Program

It is still quite difficult for a developer/designer to create new rooms using the system above. It would be ideal for a designer to just supply a text description of the room, with the room objects generated automatically.

A standalone Java program was created which takes as input a text file containing a description of a room, and outputs a Java source file for that room with the objects added to the constructor (as well as any conditionals for the description).

The program reads an *object file* containing the available objects that are in the game (and the corresponding constructor that should be used to create the object). It then takes the description in the input file and breaks it down into sentences and parses each one to determine which objects should be added to the room, as well as any conditionals. Object creation is denoted using an asterisk (*) just before the object word. The object must be in the object file for it to be added.

The semantic similarity engine used in the voice recognition system is used again here to create objects using synonyms in the room description (e.g. the sentence, “There is a utensil in the room” may add a spoon to the room since “spoon” is semantically similar to “utensil”, and a spoon is in the object file).

More interesting room generation is possible by finding binary relationships between two objects. For example, if a description in a room says, “There is a potion on the table”, then the potion being on the table should only be the case if and only if the table is still in the room; if the player destroys the table, then the description of the potion’s position should change (e.g. it should be on the floor now).

ReVerb is a Java library that takes as input a raw text file and extracts all the binary relationships in the text [36]. Each relationship is represented as a triple of `(argument1, relation, argument2)`. An example of a binary relationship is in the sentence, “The potion is on the table”, where the triple would be `(potion, be on, table)`.

These binary relationships can be used to add conditionals for descriptions. For example, when the room generation program detects the sentence, “A *potion has been placed on the table”, it detects the relationship between the potion and table, so a potion can only be on the table if the table is still in the room, and therefore adds the conditional as a `BooleanSupplier` `() -> getRoomObjectCount("table") > 0` when rendering the source file.

See the project’s repository for the full Java program source code (`GenerateRoom.java`). For convenience when building, the Java program is wrapped around a Bash script, `generateRoom.sh`, which builds the source code and executes it.

See Appendix B.6 for an example of using the room generation program to render a `Room` Java source file from a text description.

The developer can still edit this source file to add more complex functionality that was not added by the program. The room generation program was used to generate the tutorial rooms in the game (with some minor adjustments to determine which room a door leads to, as the program cannot handle that currently).

6.9 Video Conferencing Demo

In order to demonstrate the flexibility of the voice recognition system, it has been applied to a video conferencing domain. A simple demo of a video conferencing system has been created within the `CallActivity` class (separate from the RPG demo) that allows a user to call contacts using voice commands. Figure 23 on page 53 shows a screenshot of the activity.

The user is able to call a contact and add more contacts to the current call. All participants are displayed on the screen and the user can choose to end a call with a specific person, or end all of the calls. The user can also choose to mute the audio of a certain participant, or mute their own video/audio (with the video/audio icons changing accordingly).

See the `CallContextActionMap` class in the project’s repository for the actions available in the video conferencing demo. Also see Appendix B.5 for examples of commands that can be used as well as the list of possible contacts.

The system can also handle duplicate contacts if they exist (e.g. if there are multiple contacts with the first name, “Fred”). See Appendix A.5 for a screenshot of the scenario when an ambiguous contact is called.

6.10 Settings Activity

The application also has a settings activity, which allows the user to change various settings such as the current semantic similarity method used, as well as edit the user-defined synonym mappings. See Appendix A.6 for a screenshot of the settings activity.

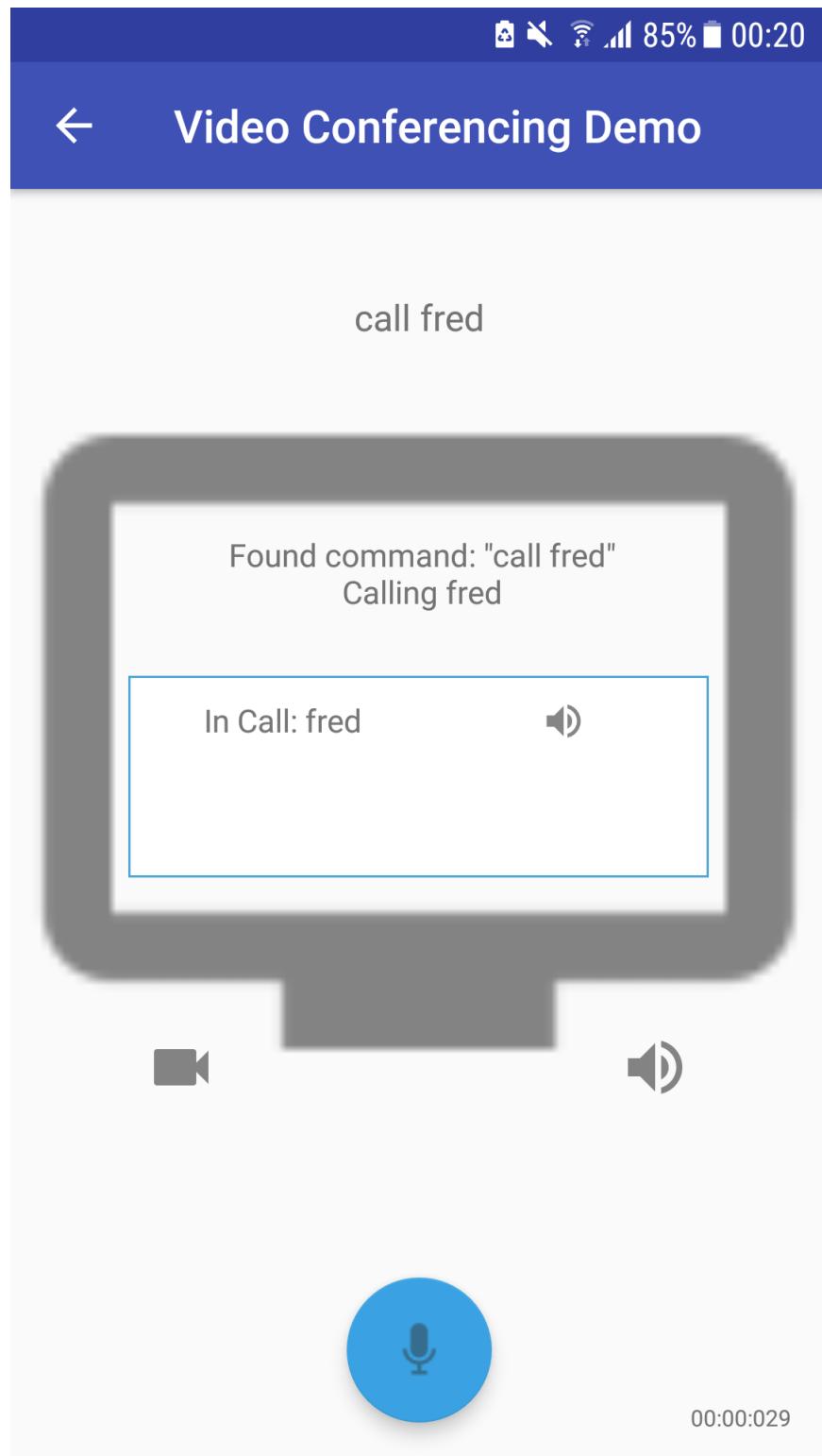


Figure 23: A screenshot of the user interface for the video conferencing demo.

7 Evaluation

This section outlines a quantitative evaluation of the project's correctness and performance. It also contains a qualitative evaluation of the practicality of the voice recognition system and whether it is easy to use via an anonymous survey⁵. It should be noted that the results use a small sample size of eight, and are only supplementary to this section.

7.1 Mock Testing

Different systems in the Android application are tested for correctness using Android Instrumentation tests and the *Android JUnit* framework. These tests run directly on the Android device without using the GUI; only the individual modules are tested in a closed environment.

There are four test classes in the Android project, each containing several test suites that test certain functionalities (a test suite can have between 3-20 different tests):

- `BattleTest.java` — tests the correctness of commands during the *Battle* mode.
- `OverworldTest.java` — tests the correctness of commands for the *Overworld* mode.
- `CallTest.java` — tests the correctness of video conferencing commands.
- `WordNetTest.java` — tests the interaction with the WordNet database; also tests the `SentenceMapper` class.

The listing below shows an example of a test suite for testing commands that show the inventory in `BattleTest.java`. Note that some tests actually test for *incorrectness* of a command (such as line 7 of the listing).

```
1  @Test
2  public void testInventorySuite() {
3      gameState.initBattleState(new Troll(9999999));
4      gameState.getInventory().add(new Weapon("hammer", "blunt", "heavy"));
5      gameState.getInventory().add(new Weapon("sword", "sharp", "metal", "pointy"));
6      testInventoryShown("show my inventory", true);
7      testInventoryShown("show my troll", false);
8      testInventoryShown("show my bag", true);
9      testInventoryShown("please show the contents of my possessions that I have",
10                      true);
11     testInventoryShown("look at my inventory", true);}
```

New tests were added frequently and were executed to ensure a new feature did not break another feature. New features could not be committed to the project's repository unless every test in every test suite of every test class passed with no errors.

See Appendix C.1 for a list of commands that are tested for correctness using the mock tests (note that not all tests are shown here). See the project's repository for the entire testing framework.

⁵<https://goo.gl/forms/8byr4MKo56bqJMvk2>

7.2 Evaluation of Semantic Similarity

The semantic similarity engine — on its own without any helper mechanisms (such as synonym-mapping, sentence-matching, match-ignoring, etc.) — has been evaluated based on the number of different commands that it can detect.

A separate Java project was created that runs the system on a set of commands across three entirely different domains: commands that may be found in a game; video conferencing commands; and commands that may potentially be used in a cooking system/application. On each run, a different semantic similarity method is used.

The number of tests that pass for each domain is calculated as a percentage of the total tests, and a final score is calculated for that semantic similarity method (as an average of all three domain scores), with the results saved as a CSV file.

In the evaluation project, the action methods do not perform any function; they simply return a string containing the output (e.g. “ATTACK_WEAPON” for attacking with a sword; “SERVE_SOUP” for serving soup, etc.). Each test case is a pair containing a string input (i.e. the user’s utterance) and the expected output. For each test case, the string input is used to update the game/calling/cooking state. The output from the object is compared to the expected output, and the results are written to the CSV file. The time taken for a test suite is also recorded, as well as the average time taken for each test.

The `ContextActionMaps` used in this evaluation project are not the same as the ones used in the Android application (they do not use any helper mechanisms). See Appendix C.2 for the context-action mappings used in this evaluation projects.

Figure 24 shows pie charts for the percentage of tests in the evaluation that pass using the Wu and Palmer (WUP) method.

The final score for the WUP method is 71.08%. See Appendix C.3 for the full table of results (including test cases) when using the WUP method.

Note that any failed test cases in the evaluation can be resolved using helper mechanisms (e.g. synonym-mappings, match-ignores, sentence-matching, etc.), which have been removed from the evaluation project in order to only evaluate the similarity method.

Table 3 on page 57 shows the final scores for each semantic similarity method outlined in Section 6.5.5, while Figure 25 shows the corresponding bar chart for the results.

See the full results for each semantic similarity method (including the result for each test case) in the corresponding CSV results file in the project’s repository⁶.

⁶The CSV files can be found in the `evaluation/evaluation-results` directory of the repository.

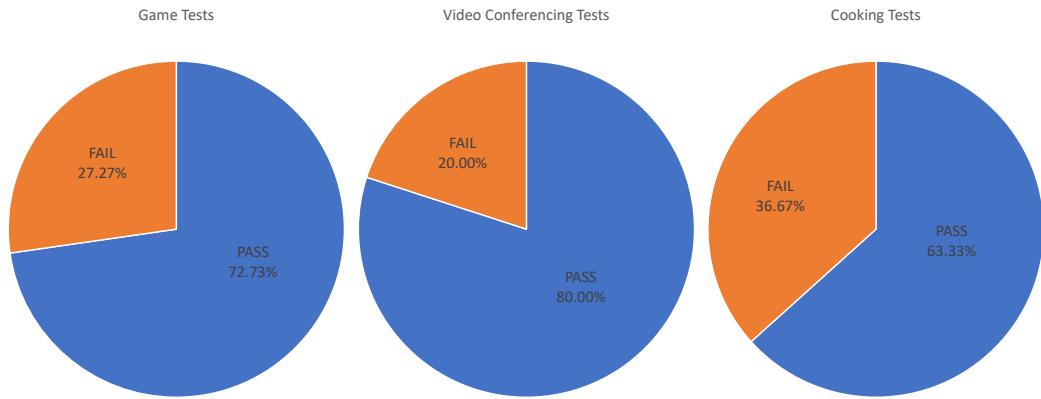


Figure 24: Pie charts showing the percentage of tests that pass in each domain using the WUP method.

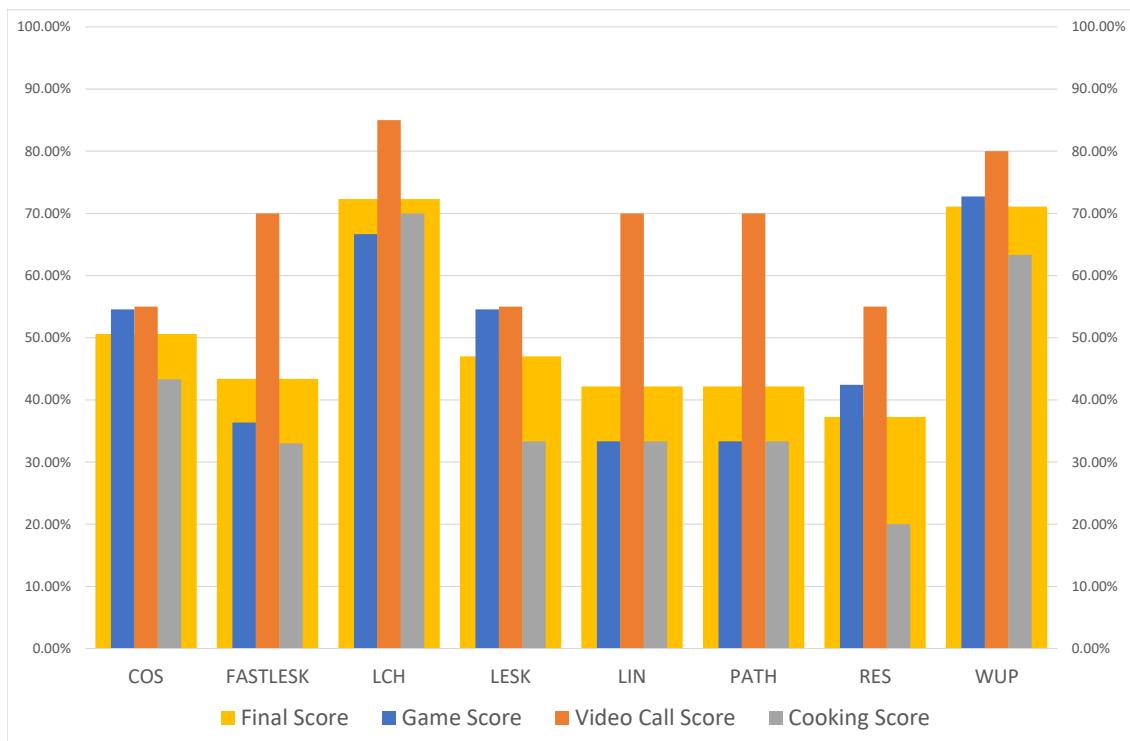


Figure 25: A bar chart showing the score for each domain and the final (average) score for each semantic similarity method in the evaluation project.

Table 3: Evaluation of Similarity Methods

Method	Game Score	Video Call Score	Cooking Score	Final Score
COS	54.54%	55.00%	43.33%	50.60%
FASTLESK	36.36%	70.00%	33.00%	43.37%
LCH	66.67%	85.00%	70.00%	72.29%
LESK	54.55%	55.00%	33.33%	46.99%
LIN	33.33%	70.00%	33.33%	42.17%
PATH	33.33%	70.00%	33.33%	42.17%
RES	42.42%	55.00%	20.00%	37.25%
WUP	72.73%	80.00%	63.33%	71.08%

The LCH method has the highest final score, closely followed by the WUP method. However, since the WUP method produced a higher score for the game domain (72.73% which is much higher than the LCH score of 66.67%), it was chosen to be the default method used in the role-playing game demo.

The results suggest that path-based methods perform better than methods which use information content or the definitions of the words. This could be because words which are close in proximity in the WordNet database tend to be very similar. The exception to this is the PATH method (i.e. using the reciprocal of the distance between the word nodes in the tree) which performed very poorly with a final score of 42.17%, perhaps because it is a very simplistic algorithm.

These results also suggest that different methods may perform better under different domains, and there is no optimal method that performs the best under all circumstances. However, it is clear that some methods such as the Resnik (RES) method perform poorly under all three domains, and are not suitable for the system.

It should be noted that the performance of the methods could also be dependent on the confidence thresholds that have been set in `VoiceProcess.java`; adjusting these thresholds could affect the accuracy of the methods either positively or negatively.

7.2.1 Combining Similarity Methods

It is possible to combine two different semantic similarity methods and take the average confidence value when comparing the similarity of two words. The evaluation project also calculates the overall scores for different combinations of methods.

Table 4 shows the final scores for the best combinations of methods. Note that not all combinations are shown; see the project’s repository for every combination’s final score. Figure 26 shows a bar chart of the results for all combinations.

Table 4: Final scores for the best-performing hybrid methods

Method 1	Method 2	Final Score
FASTLESK	LCH	74.70%
LESK	LCH	74.70%
WUP	COS	73.49%
LCH	PATH	72.29%
WUP	LCH	71.08%

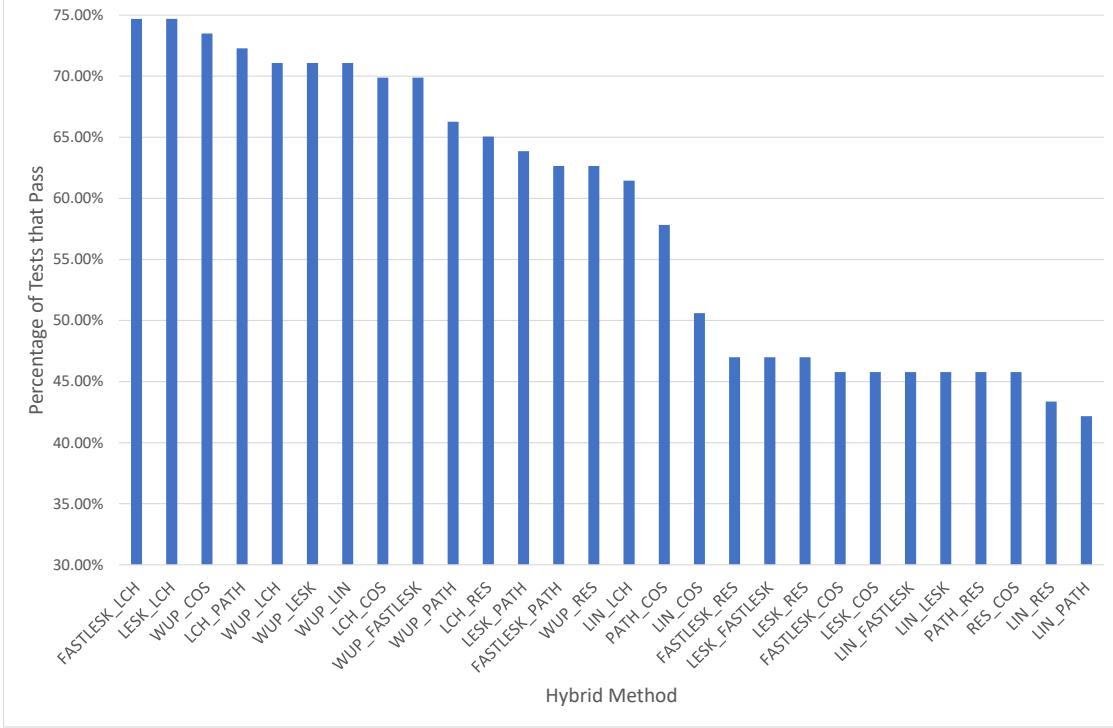


Figure 26: A bar chart showing the final (average) score for each combination of semantic similarity methods.

The Leacock and Chodorow (LCH) and Wu and Palmer (WUP) methods are involved in almost all of the top combinations (since they also have the highest individual final scores). The best combination is the LCH method with either the LESK method or FASTLESK method, and since the FASTLESK method can be up to 10 times faster than the LESK method (see Section 7.3 below for performance details), it is the preferred combination.

It should be noted that there is only a slight increase from 72.29% to 74.70% when combining the LCH method with the FASTLESK method (as opposed to just using the LCH method on its own), so the increase in accuracy may not be worth the increase in processing time when using two methods.

7.3 Performance of Semantic Similarity Methods

One user requirement for the system is for the voice recognition to be fast and responsive; the user does not want to wait a noticeable amount of time until their voice command has been processed. Therefore, a semantic similarity method which is fast is required.

According to Jakob Nielson of the Nielson Norman Group [55], there are three main response limits determined by human perceptual abilities that should be noted when optimising application performance:

- 0.1 seconds — the limit for which the user feels the system is instantaneous.
- 1.0 seconds — the user will notice the delay, but will feel they are not unduly waiting.
- 10 seconds — the limit for keeping the user's attention; anything longer requires a progress indicator (e.g. progress bar).

Ideally, command processing should be limited to 0.1 seconds, with 1 second being the limit for processing commands involving multiple intents, targets or contexts (e.g. “attack the troll with a sword and with a hammer and then heal using a potion”), as the intent would be broken down and processed separately, so may require more processing time.

Using the evaluation project described in Section 7.2 above, the average time taken to process one command was calculated for each semantic similarity method on three different machines. The following are the specifications of the machines. Note that the Java program is single-threaded, so multi-core processors are not fully utilised.

- Desktop PC
 - 4 x Intel Core i5-6600K quad-core CPU @ 3.5 GHz
 - 8 GB RAM
- Samsung Galaxy S7 Phone
 - 4 x Exynos M1 Mongoose octa-core CPU @ 1.6 GHz [56]
 - 4 GB RAM
- Raspberry Pi 3
 - 1 x ARMv8 quad-core CPU @ 1.2 GHz [57]
 - 1 GB RAM

Table 5 below shows the average time taken to process a command (in seconds) for different semantic similarity methods, across three machines.⁷ See the results in the project’s repository for the average times for each domain.

⁷The LESK method is taken over the average of six commands on the phone, and three commands on the Raspberry Pi, due to time constraints.

Table 5: Average Time Taken to Process a Command

Method	Average Time per Command / s		
	PC	Phone	Raspberry Pi
COS	0.023	0.85	6.40
FASTLESK	0.012	1.5	5.00
LCH	0.004	0.099	3.29
LESK	1.35	280.02	2068.81
LIN	0.009	0.29	5.19
PATH	0.004	0.088	2.96
RES	0.007	0.12	3.29
WUP	0.015	0.16	5.15

On the Desktop PC, all methods except the LESH method are below the 0.1 seconds limit, so would all be feasible to use. For the Android phone, only the LCH, PATH, RES and WUP methods are close to, or fall below, the 0.1 seconds limit. Unfortunately, all methods take longer than 0.1 seconds on the Raspberry Pi, but all methods fall below the 10 seconds limit (except LESH). The slow speeds on the Raspberry Pi suggest that the system is not suitable for low-power processors, such as embedded systems, but the system is suitable for modern Android devices given the speeds on the phone.

Clearly, the LESH method is not feasible on any machine, with the FASTLESK method being manyfold faster. This suggests that the use of definitions from the meronyms and holonyms of the words to calculate the similarity between two words (and any formatting to remove punctuation from the definitions) is not worthwhile in terms of speed (and given that the accuracy of both methods are very similar as seen in Section 7.2, the FASTLESK method is unanimously the preferred method).

Methods which produced high accuracy such as LCH and WUP also have fast performance, so they are ideal methods to use in the application. It should also be noted that multi-threading the system may be able to, at least, double the speed of processing.

7.4 Evaluation of System Features

As mentioned previously, the system supports mechanisms that the developer can use to improve the recognition of various commands. The following is an evaluation of the effectiveness of the helper mechanisms in increasing the number of commands that the system can process correctly.

The system supports the following features that can be toggled in the `VoiceProcess` class:

- Using the semantic similarity engine to detect similar words.
- An alternative slot-filling grammar, {WITH CONTEXT ACTION TARGET}, when needed.

- Creating a mapping of synonyms to their actions, targets or contexts.
- Ignoring matches for similarities between two words.
- Handling ambiguous intents with a suggestion and confirmation.
- Performing sentence matching.

Using the test cases from the Android mock testing (see Appendix C.1 for the list of commands that are tested), the number of tests that pass are calculated while gradually activating more features. It is not possible to test the features separately as some are dependent on others being enabled (e.g. the suggestions mechanism requires the semantic similarity engine to be enabled).

Figure 27 shows a cumulative line chart of the number of tests that pass as more features are enabled, starting with no features, followed by adding the semantic similarity engine, and so on. There are a total of 158 tests, and with all features enabled, all of the tests pass.

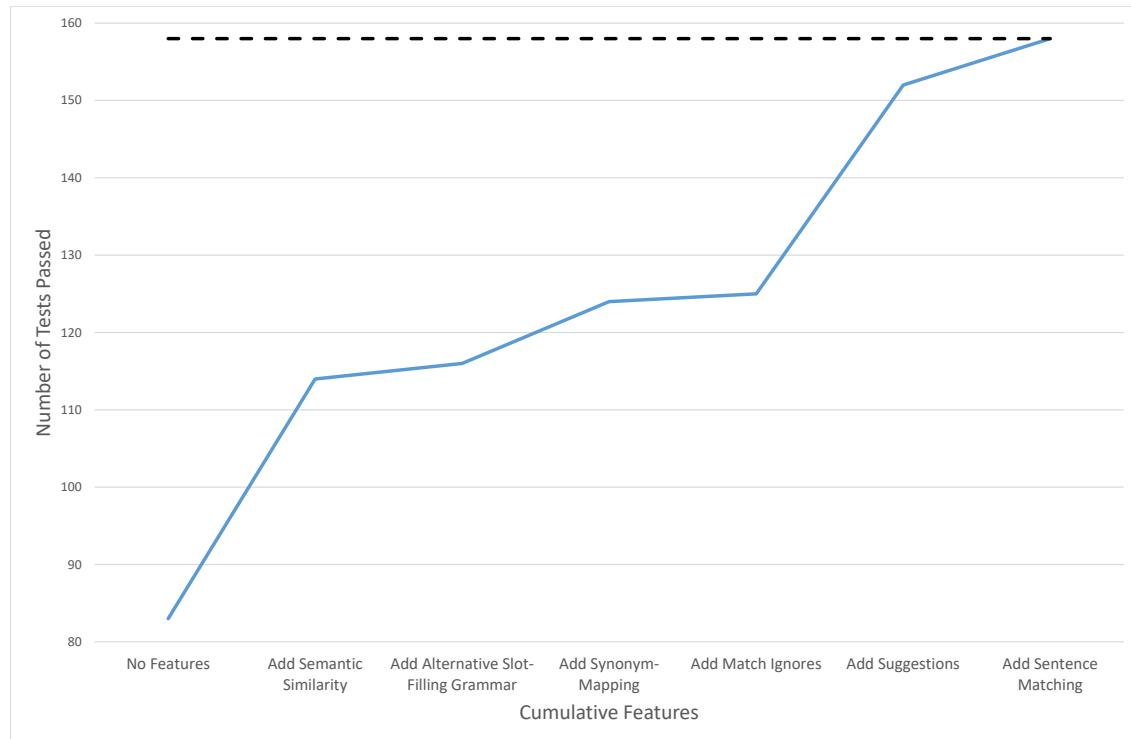


Figure 27: Cumulative line chart showing the number of tests that pass as more system features are activated.

The large increases in the number of tests that pass when adding the semantic similarity engine or the suggestions mechanism indicates that they have a more sizeable impact

compared to the other features (i.e. most variations of commands can only be detected if these features are enabled).

7.5 Evaluation of Slot-Filling Grammar

The sentence-matching mechanism was added to compensate for commands that do not have an imperative structure (which is what the current slot-filling grammars detect; see Section 5.3.3 on page 23). These include intents that are styled as questions, or greetings.

Using the anonymous survey, participants were asked what kind of phrases they would say in a role-playing game. Figure 28 shows the results for this from eight responses. The question posed was, “If you were playing a role-playing video game (e.g. Pokémon, Final Fantasy, Elder Scrolls) that supported voice commands, what type of commands are you likely to eventually say?”

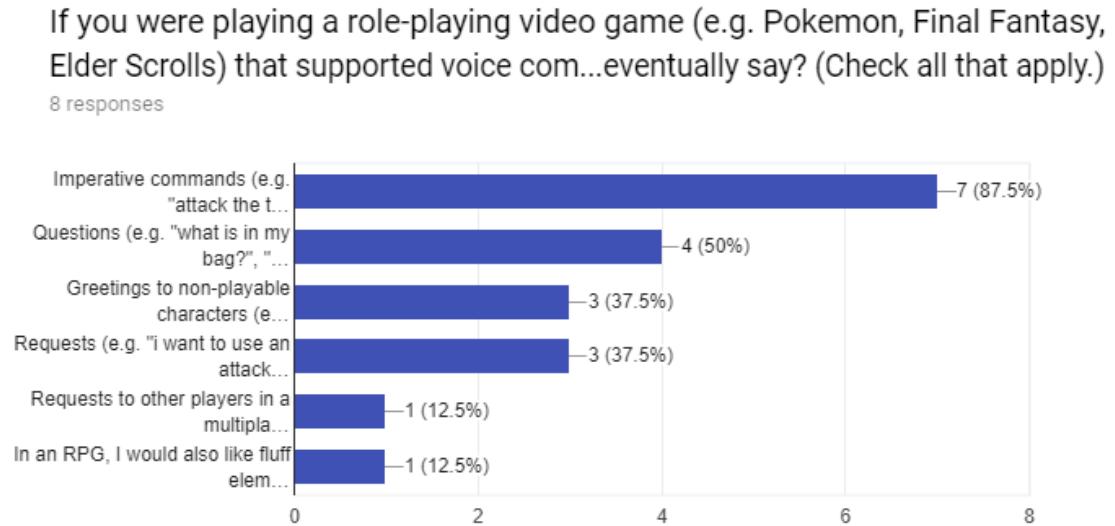


Figure 28: A bar chart showing the number of users that would eventually say certain types of commands in a role-playing game.

While the system covers imperative commands and requests, it does not cover questions and greetings (for which 50% and 37.5% of respondents would eventually say during the game, respectfully). Therefore, the sentence-matching mechanism would need to be used to satisfy these users. Unfortunately, since questions can vary greatly in structure (e.g. “What can I do?” versus “Where should I go next?”), as well as greetings, it is not possible to fit them to a single slot-filling grammar.

7.6 Application Performance

Figure 29 shows the CPU, memory and network usage of the Android application running on a Samsung Galaxy S7 for the first 60 seconds when the application is opened (according to the Android profiler included in Android Studio). See Appendix C.4 for the list of commands that were executed in the game during the profiling.

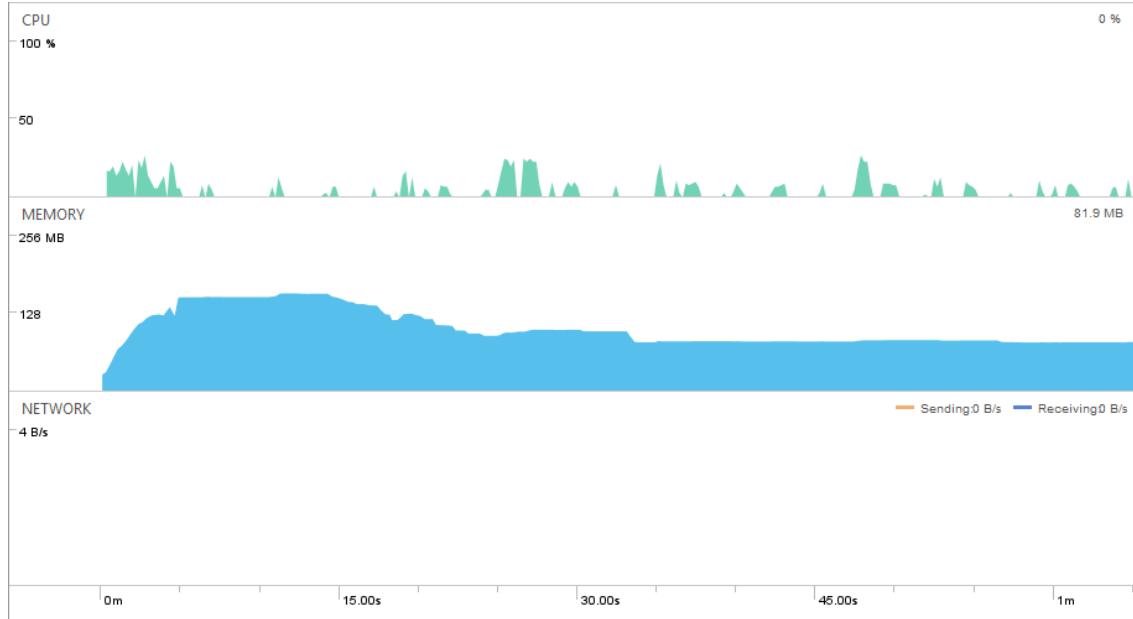


Figure 29: Profile of the Android application for the first 60 seconds when the application is opened. Graph generated by the Android profiler in Android Studio.

Figure 29 shows that no network data is sent or received by the application, as desired. However, the application has a large memory footprint of approximately 100 MB, due to the POS tagger model and the WordNet database being loaded. Since the application is single-threaded, it is expected that the maximum CPU usage is approximately 25%.

7.7 Usability of the Voice Recognition System

It is difficult to quantitatively measure how easy it is to add voice commands with this voice recognition system compared to more traditional methods. One method is to compare the number of lines of code required for a developer to write in order to implement the same number of voice commands in their application.

The voice recognition system implementation will be compared to two naïve/traditional systems for adding voice commands to an application.

Hard-coding Strings

Adding every possible command as a string and comparing the user's input to them

using a series of `if` and `else` blocks, or a `switch` statement.

Regex-style Matching

Write a word-based regular expression to describe the structure of acceptable commands (see Section 4.3.3 on page 10 for an example using Houndify).

When mapping a large number of commands, both of the above methods become increasingly complex and infeasible. With the system presented in this project, only a few lines of code need to be written to map the same number of commands, and more. See Appendix C.5 for a comparison of the methods with an example.

According to the anonymous survey, developers also preferred a system where they only have to write a few lines of code to implement countless voice commands. See Appendix C.6 for these results.

In order to qualitatively evaluate the usability of the system further, a standalone Java library was created, and the process of using the library in order to add voice commands to a Java application is outlined in a user guide for the library. After studying the guide, developers were asked about the ease of integrating the library into a project via the anonymous survey. Figure 30 shows the responses to this question.

Consider the following user guide (link below) for the above Java library.
How complex are the instructions?

8 responses

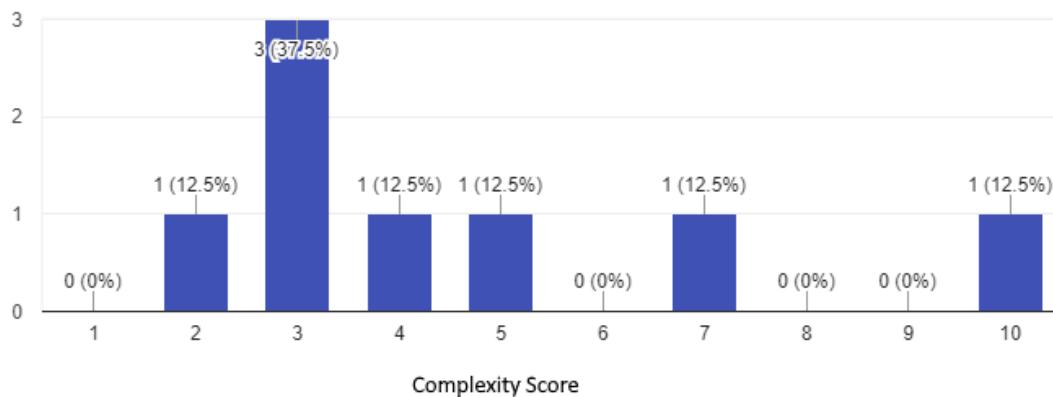


Figure 30: A bar chart showing the responses for how easy it is to use the voice recognition system in a project on a scale of 1-10, with 1 being “not complex at all” and 10 being “extremely complex”

Despite the small sample size, the majority of respondents found the system was not extremely complex to use (by giving it a score of 5 or less). However, a couple of respondents

found the system to be very complex to use. This could be because the system requires a lot of dependencies and models to be loaded in order to work, or because the initial steps to set-up the system (e.g. creating the `ContextActionMap`, `Entity` and `Action` classes for the first time).

7.8 Voice Recognition System Limitations

Several limitations of the voice recognition system were identified over the course of this project, and are described below:

Only works for Java

The system relies heavily on OOP inheritance, as well as Java-based libraries such as JWI and WS4J. While it may be possible to port the system to another language, it would require a lot of work and re-writing of the source code.

Cannot handle negatives in intents

The system currently cannot handle negatives in a user's utterance. For example, in the video conferencing example, if the user says, "call Dan with no audio", the "no" will be ignored and the system will create a call with Dan using audio. It may be possible to fix this by searching for any negatives and altering the intent as needed, but due to time constraints this could not be implemented and tested.

Only single-word actions, targets and contexts are allowed

The voice recognition system scans the input word-by-word, searching for the action, target and context (if they exist). Therefore, it currently cannot identify open compound words such as "ice cream" or "full moon", or hyphenated compound words such as "well-being", "merry-go-round" or "on-site" [58]. It may be possible to fix this by implementing a lookahead mechanism which looks at the next word in the input and decides whether it is part of the current word.

This also includes instances where two words are used to describe an action. For example, in the phrase, "take a look around", the word, "take" will map to the "grab" action first, so the system thinks the player wants to grab the surroundings which is not correct.

Not designed for many different targets per action

The use of a `ContextActionMap` is to handle many different contexts that the action can be performed with. For example, an attack using a blunt weapon may crush the enemy, but an attack from a sharp weapon may slice the enemy. However, in situations where an action is applied differently to many different types of targets, a large block of `if` and `else` statements or a `switch` statement is required to perform different things for different targets within an `Action` class (e.g. writing a series of `if (currentTarget instanceof ...) {...} else if ...`).

This may be resolved by having a `TargetActionMap` instead of a `ContextActionMap`, which indexes the 'table' using the target 'type' instead. However, this would mean

that the same issue will occur when having multiple contexts perform different actions (unless both a `TargetActionMap` and a `ContextActionMap` can be used together).

Actions cannot have the same string as targets

There is a bug where an action cannot have the same name as a target, as the system becomes confused when parsing and may think the target is actually the action (if the word for the target is a better match than the word of the action in the input).

For example, in the video conferencing example, there may be two actions, “call” and “stop”. If the user intent is, “finish the call” where “call” is the target entity, the system will think that “call” is a better match to the “call” action than “finish” is to the “stop” action, and will invoke the “call” action instead. This can easily be fixed by checking the ordering of the matches to see if the action precedes the target.

Cannot handle fully ambiguous inputs

This occurs when a user’s intent can be mapped to two completely different actions in the application and both would be valid due to the ambiguity of the intent. For instance, take the phrase, “throw the match”: it could refer to giving up in a competitive competition, or throwing a matchstick.

Many of these limitations are quite severe and hamper the usability of the system under certain circumstances, particularly the restriction to single-word actions, targets and contexts, which would mean actions such as, “show off” or “show up” cannot be disambiguated. However, many are relatively easy fixes, and almost all of them could be fixed over time, without breaking any existing functionality.

7.9 Room Generation Evaluation

Since the semantic similarity engine used in the room generation program has already been evaluated above, only the effectiveness of the object relationship mechanism needs to be evaluated. This can be done by creating a complex description of a room with all but the first sentence containing a different relation between two objects, and seeing if the system detects the relationship. See Appendix C.7 for the full results.

While the program can detect relations in sentences such as, “A letter sits upon the table” or, “A potion is underneath the table”, it cannot detect some simple relations such as, “There is a chair under the table” or “A fork lays on the table”. This could be because the *Reverb* library does not understand the relations, “lays on” and “under”. However, the majority of relations are detected by the program, and only a few conditionals have to be added manually by the developer.

The room generation program currently cannot handle multiple objects being created in the same sentence, so only one object creation can be associated to each sentence in the room description.

8 Conclusion

This section summarises the findings of this report, as well as the deliverables and contributions of this project. The project work can be found in the main repository ⁸.

8.1 Deliverables

A prototype for a text-based role-playing game has been created on Android which uses the new voice recognition system to process commands, and is available to download ⁹. The system makes it easier for developers to detect a wide variety of commands without the developer needing to hard-code any of the phrases (or specify the structure of acceptable phrases); they only have to create the table of context-action mappings, and use the system features, to process commands.

Unlike paid voice recognition services such as IBM’s Watson Conversation, Houndify or DialogFlow, this project’s voice recognition system is able to work completely offline (and therefore protects the user’s privacy) and the processing time per command on an Android device is within an acceptable limit (close to 0.1 seconds).

While there are some limitations to the system, using the WordNet database to calculate semantic similarities between words has produced a working foundation for the system that can be improved, and the system has been applied to other domains such as a video conferencing application.

A separate standalone Java library has been created for the voice recognition system that is open-source for developers to explore, called *Voice Commands with WordNet* ¹⁰. It includes a user guide explaining how to integrate the library into other Java projects.

8.2 Contributions

This project also explored and evaluated different semantic similarity methods in terms of accuracy and performance, with path-based methods such as the Leacock and Chodorow (LCH), and Wu an Palmer (WUP) methods performing the best.

A new interface for the WS4J library was created which uses MIT’s JWI API to access the WordNet database ¹¹. This interface has been made public and may be useful to other developers in the future.

Finally, the semantic similarity engine created for the voice recognition system has also been applied to other areas of development for the RPG (e.g. object property assignment and room generation). A room generation program has been created which takes as input

⁸<https://github.com/BaronKhan/VoiceRecognitionRPG>

⁹<https://play.google.com/store/apps/details?id=com.khan.baron.voicerecrpg>

¹⁰<https://github.com/BaronKhan/voice-commands-with-wordnet>

¹¹<https://gist.github.com/BaronKhan/5679157d1a8572debc8c75af1c875c16>

a text description of the room, and generates a source file for the room, making it easier to create more complex rooms for the game.

Ideas and techniques used in the room generation program, such as identifying binary relationships between objects and using the semantic similarity engine for object creation, may be useful for other systems and projects.

8.3 Future Work

This section outlines any future work that may be done as a result of the work already done in this project.

8.3.1 Language Support

The voice recognition system — from the slot-filling grammar to the WordNet database — only works for the English language, and needs to be changed for other languages.

While approximately 42% of the world's language use a *subject-verb-object* sentence structure (e.g. "He hits the goblin"), such as English, Mandarin and Spanish, approximately 45% of languages use a *subject-object-verb* sentence structure (e.g. "He the goblin hits"), such as Bengali, Greek and Japanese [59]. Therefore, the target would precede the object, so the voice recognition system's slot-filling grammar needs to be rearranged.

The Princeton University WordNet database only supports English words. While there are WordNet databases for most of the world's languages [60], each database requires a new interface to be implemented (although most of them are SQL databases which can be queried), and not all languages have a WordNet database currently.

8.3.2 Deep Learning

This project has used classical natural language processing techniques, but it may be possible to use deep learning to improve the accuracy and usability of the voice recognition system, such as using recurrent neural networks (RNN). While deep learning requires a large amount of processing and would not be feasible to train offline, models can be trained by the developer and included in the application, with updates to the model included with updates to the application or game.

One use for deep learning is to be able to choose the most appropriate semantic similarity methods to use for a specific domain. As shown in Section 7.3 on page 55, different methods perform better for different domains. Currently, only a combination of upto two methods can be used, with each method having a 50% weighting. Using deep learning, the best weightings for several methods can be trained to produce the most accurate and best performing combination.

Another area where deep learning (or classical machine learning) techniques can be used

is the suggestions mechanism. A global or local model can be trained whenever an ambiguous intent is spoken, so that the system learns that the user meant another intent instantly (instead of constantly asking them for confirmation of their intent).

8.3.3 Multi-threading

During this project, an attempt to multi-thread the application was made, but was unable to achieve a sizeable improvement in performance, and also occasionally experienced deadlocks in the mock tests. Multiple processors could be used to fork and join tasks instead of executing them sequentially. This can be achieved using Java 8's parallel streams [61].

For example, in the `SentenceMapper` class, instead of iterating through a `Collection` of `sentences` to calculate their cosine similarity score, and accumulating the output sequentially, as follows:

```
1 for (String sentence : sentences) {  
2     totalScore += calculateCosScore(words, sentence);  
3 }
```

A parallel stream can be used instead, where each iteration is executed in parallel, and the outputs reduced to the sum:

```
1 totalScore = sentences.parallelStream()  
2     .mapToDouble((sentence) -> calculateCosScore(words, sentence))  
3     .sum();
```

Most of the multi-threading can be performed in the `VoiceProcess` class, where words are currently processed sequentially. Instead, they could be processed in parallel, although a lot of re-factoring is required to make loop iterations independent within this class.

Another possible use of multi-threading is to run semantic similarity calculations in parallel if the application is using a hybrid method (the average of two methods).

8.3.4 Hands-Free UI

One of the novelties of voice controls is that the user doesn't have to touch any physical buttons or screens in order to execute commands in the application. The interface of the role-playing game in this project requires the player to push a blue button in order to start a voice intent.

It may be possible to replace the button with a wake word detection mechanism. This is a phrase that the user says to indicate the start of a voice intent (such as saying, "Hi Siri" or "Ok, Google"). This requires the microphone to always be listening and consume more power, but it will work offline and there are open-source APIs available that offer this, such as *Snowboy* [62].

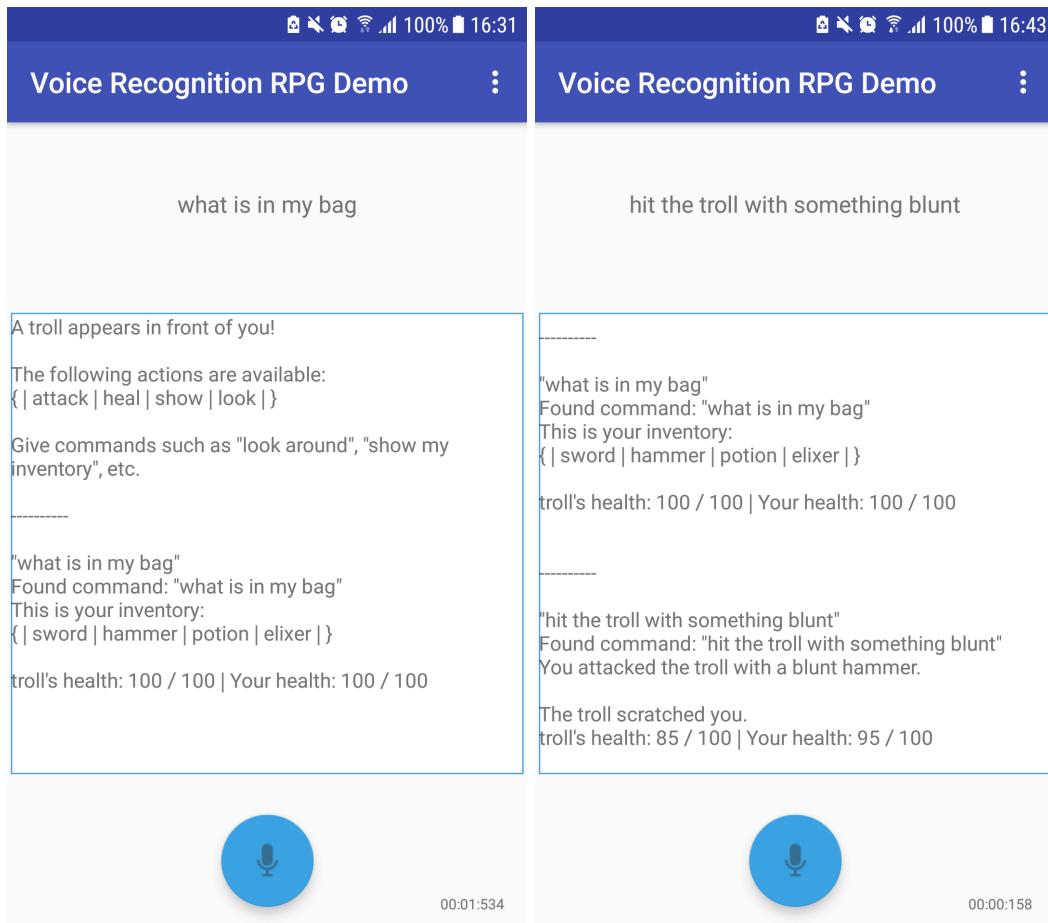
It is also possible to turn the output text into audio, using text-to-speech (TTS) APIs such as Google's TTS engine. Combining these two features would make the application completely hands-free, and could be beneficial to some players.

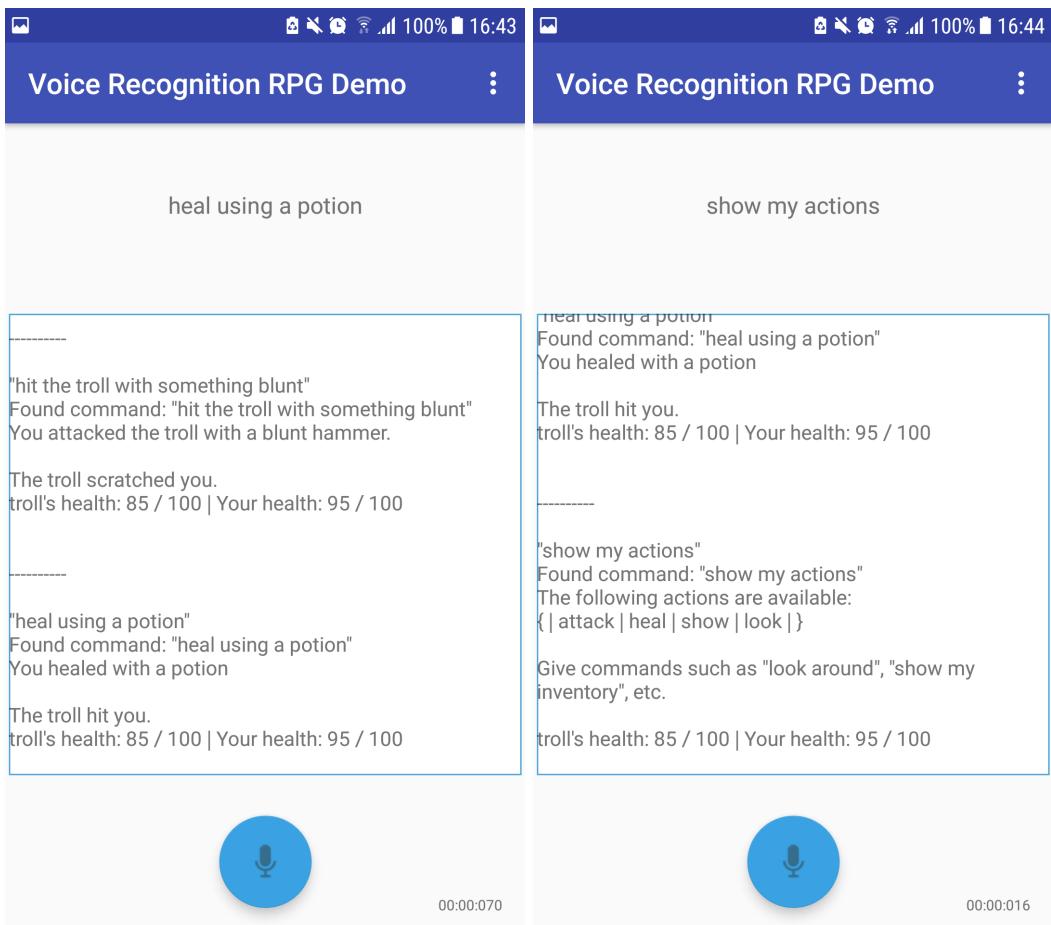
A Android Application Screenshots

This section contains snapshots of the Android application during different scenarios.

A.1 Battle Mode Example

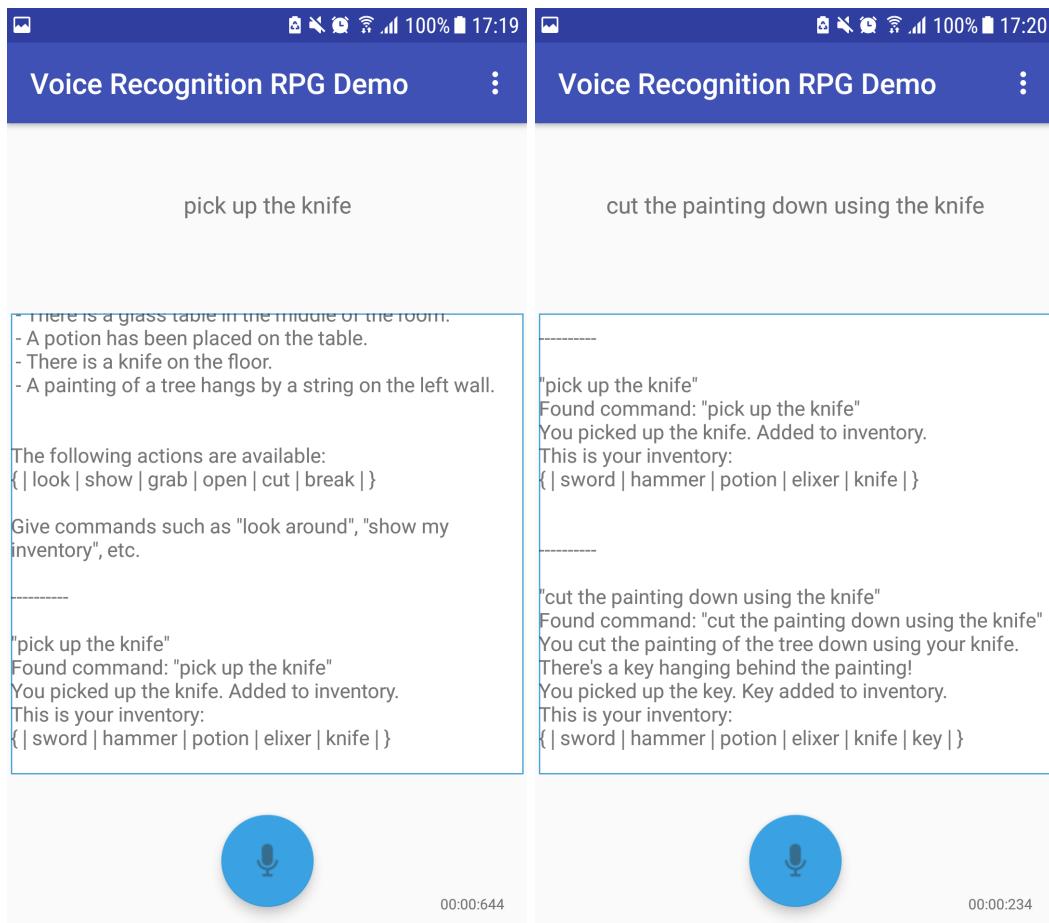
Below are screenshots of actions performed as the player battles a troll enemy.

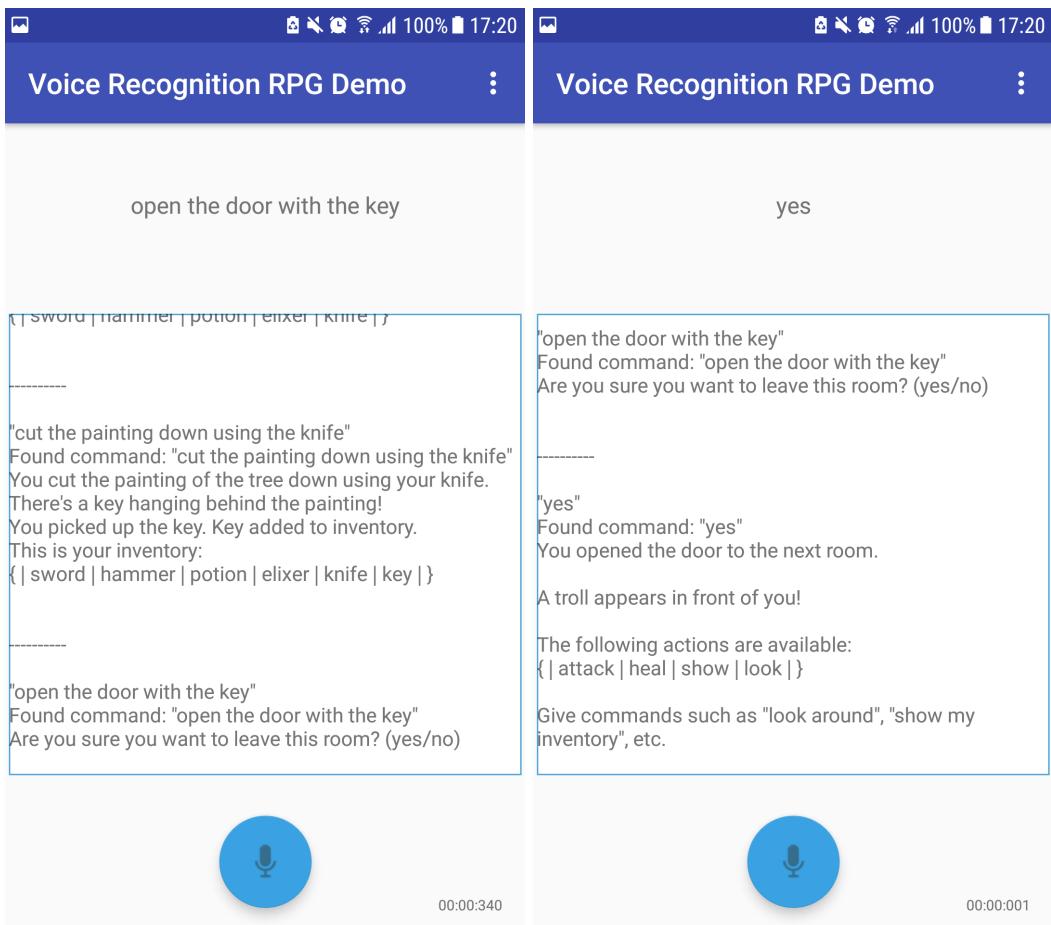




A.2 Overworld Mode Example

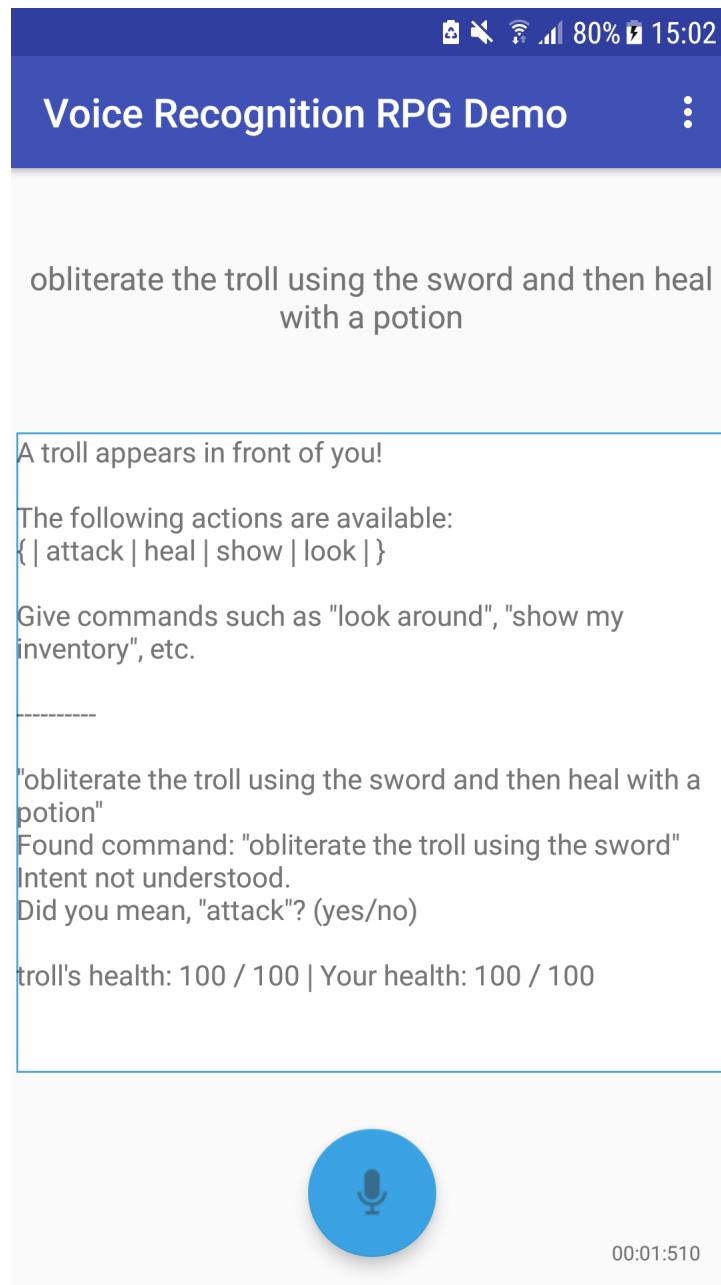
Below are screenshots of the first room in the game, and the actions required to solve the puzzle in the room.



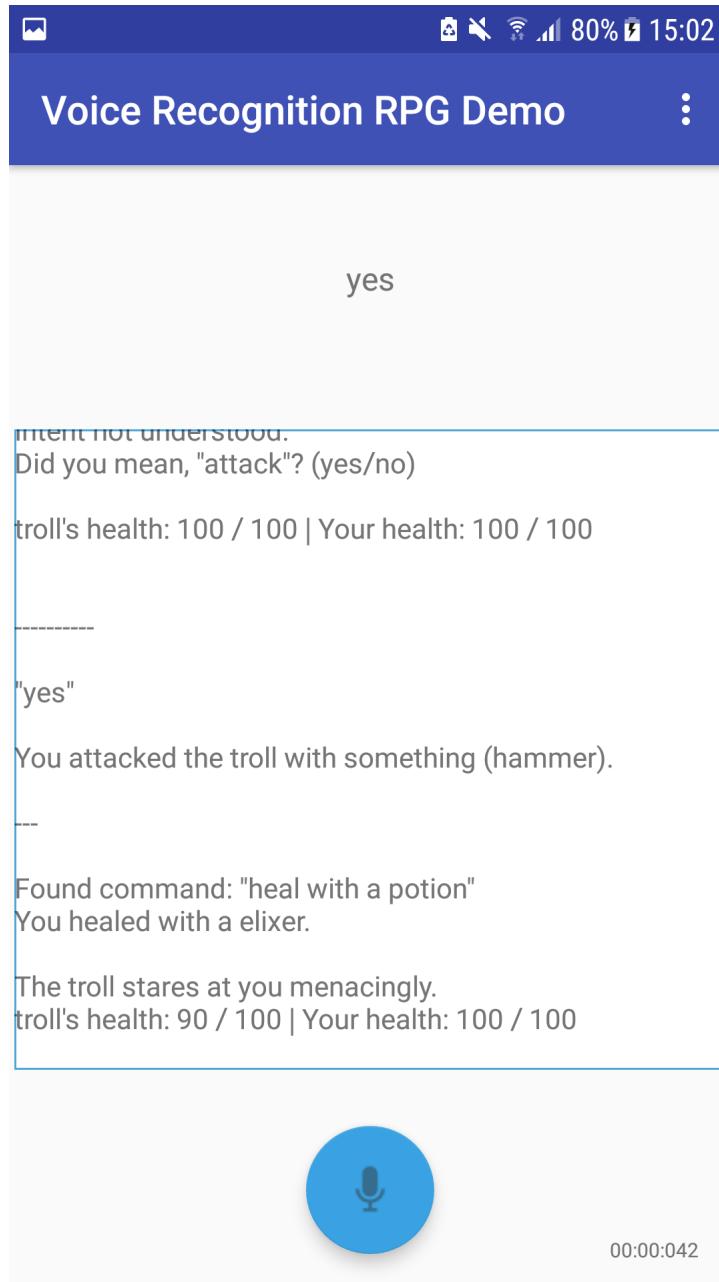


A.3 Multiple Commands Example

The following is an example of a multiple-command intent being interrupted and then continued. In the first snapshot, the system asks the user to resolve the first intent ("obliterate the troll using the sword"):

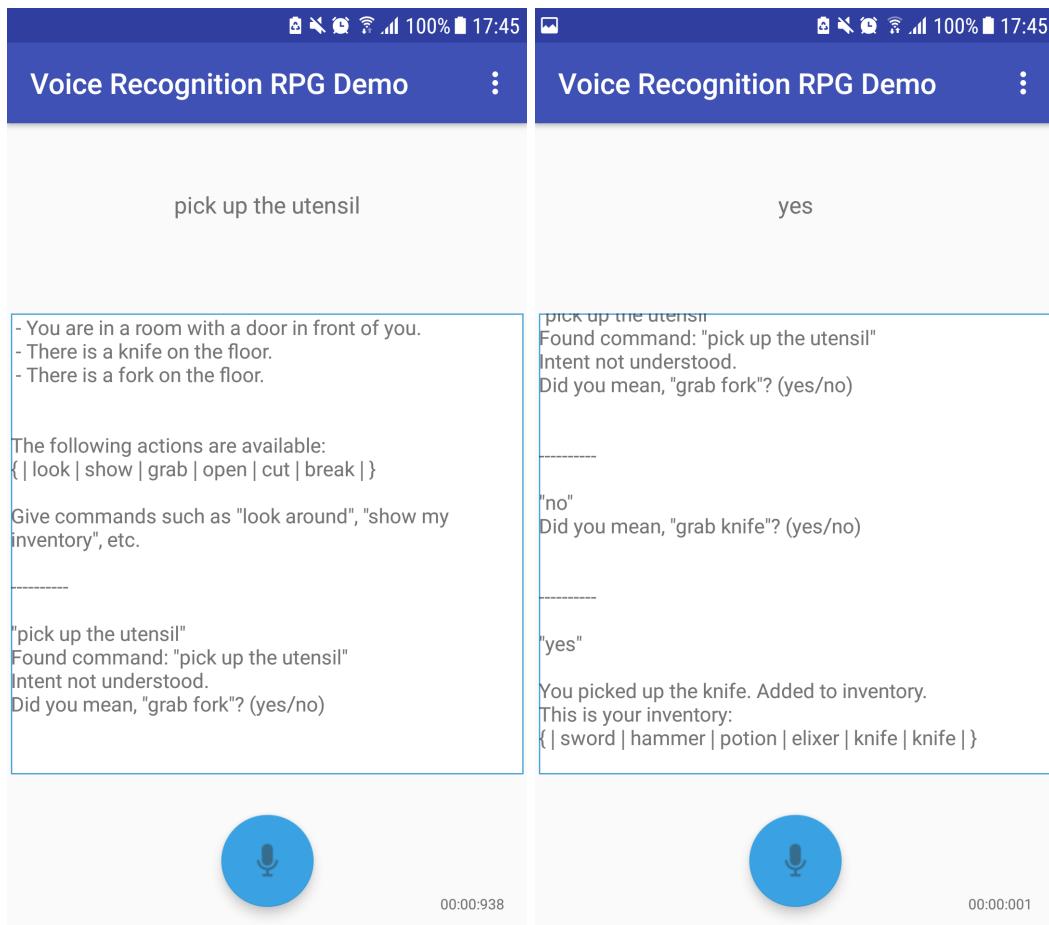


Once the user resolves the first intent, it executes it followed by the second intent ("heal with a potion"):

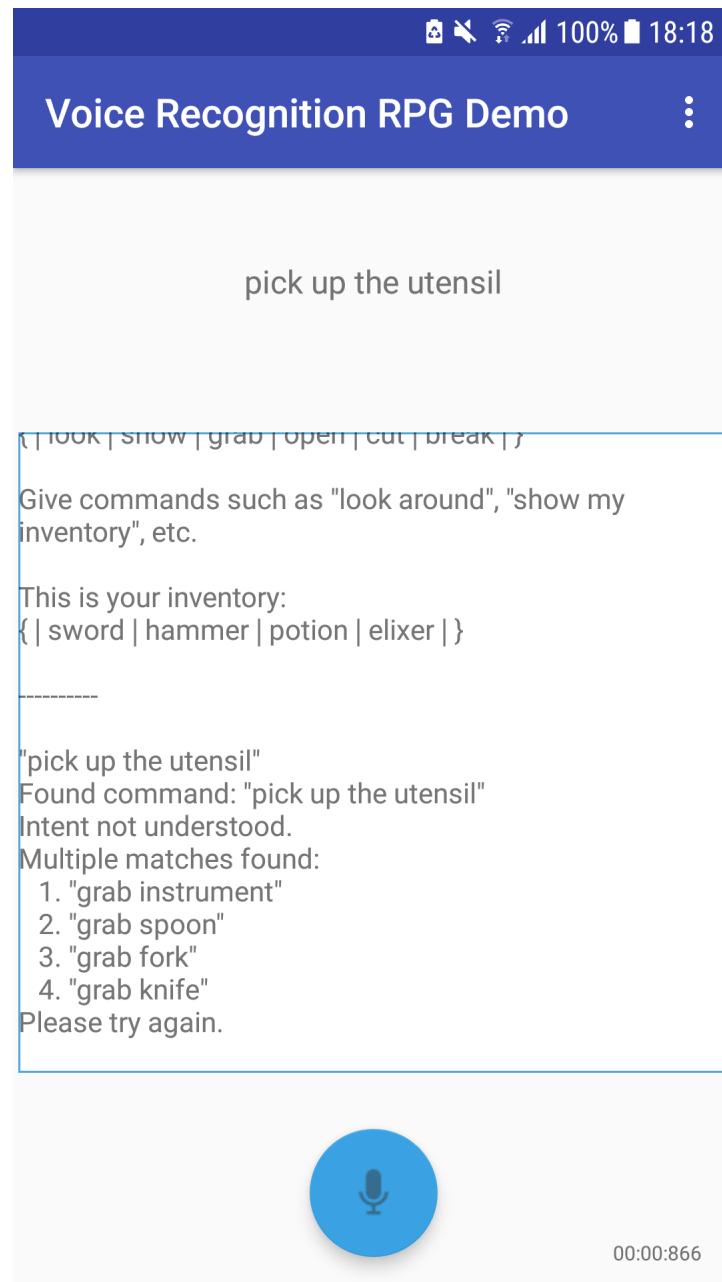


A.4 Ambiguous Utensil Example

Below are screenshots of the room containing both a fork and a knife, with the synonym mapping of `utensil --> [fork, knife]` defined.



If there are multiple mappings for utensil to various objects (e.g. `utensil --> [fork, knife, spoon, instrument]`), then all matches are displayed at the same time. Below is a screenshot of a room with these objects.



A.5 Duplicate Contacts

Below is a screenshot of the video conferencing demo when an ambiguous contact is called. The user would reply with the surname of the correct contact.

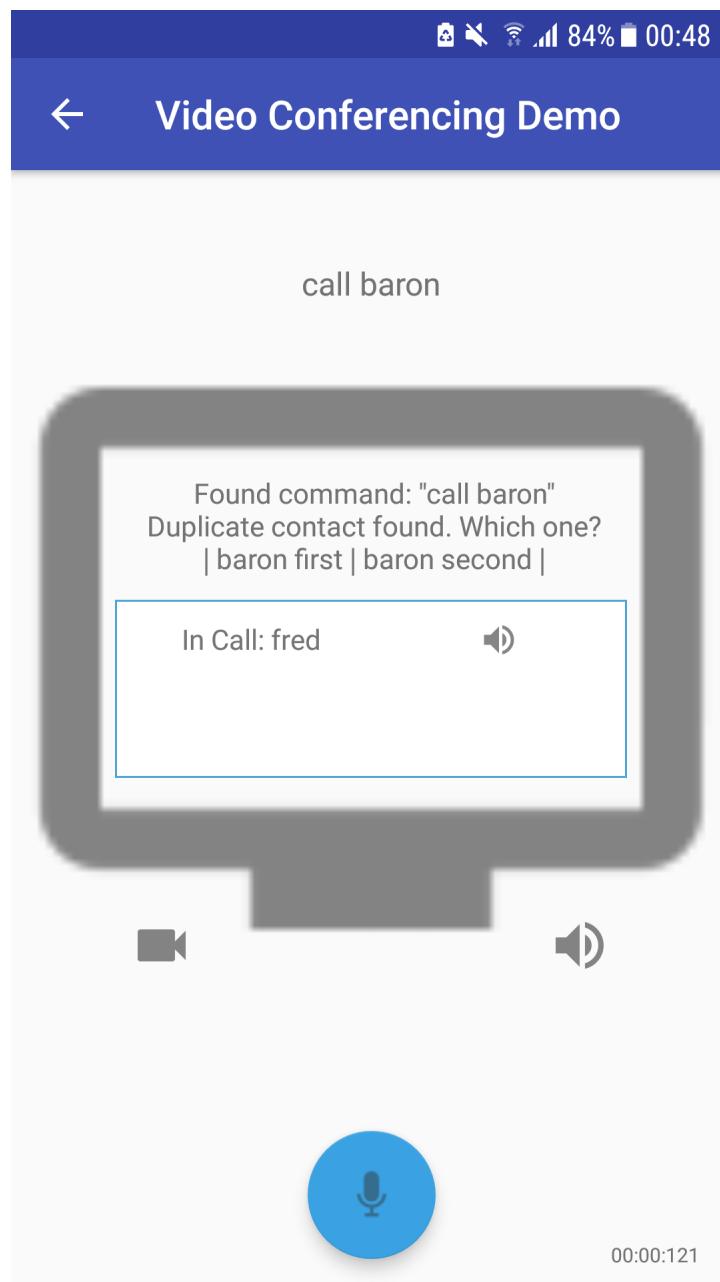


Figure 31: A snapshot of scenario involving duplicate contacts.

A.6 Settings User Interface

Below is a screenshot of the `MainActivity` with the drop-down menu open. It contains an option to access the settings activity:

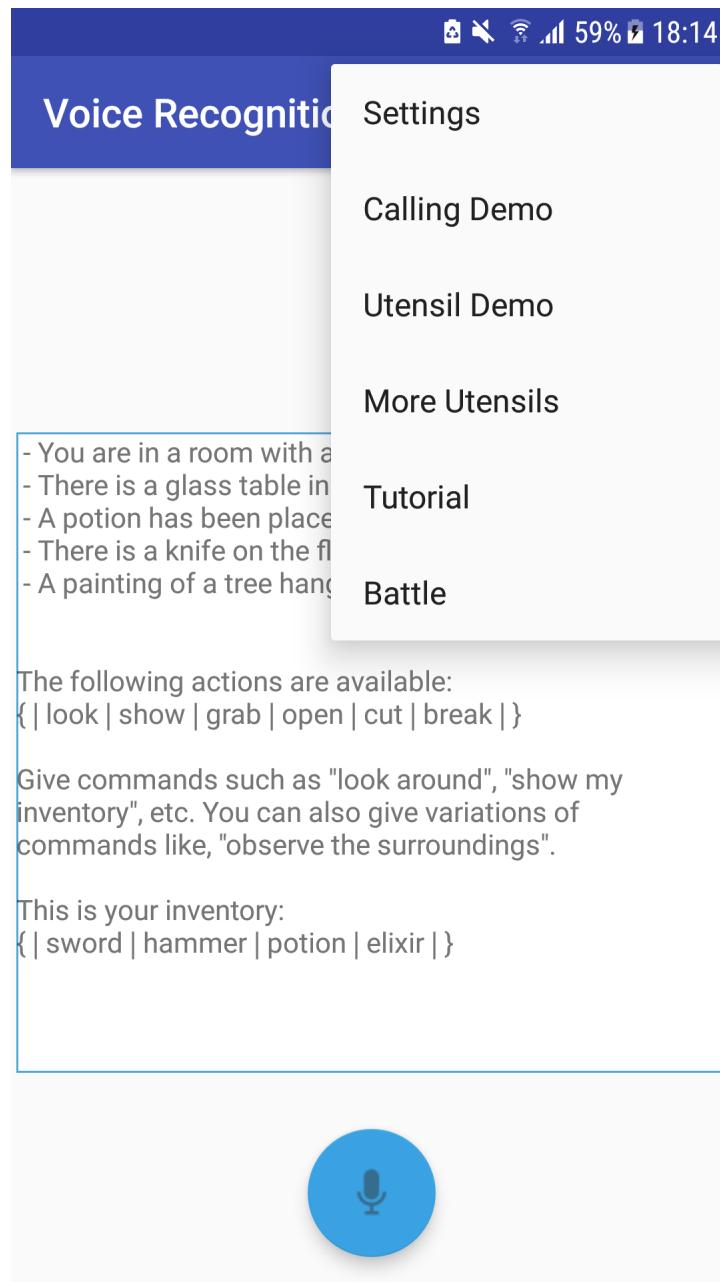


Figure 32: A snapshot of the `MainActivity` with the drop-down menu open.

Below is a screenshot of the graphical user interface for the settings.

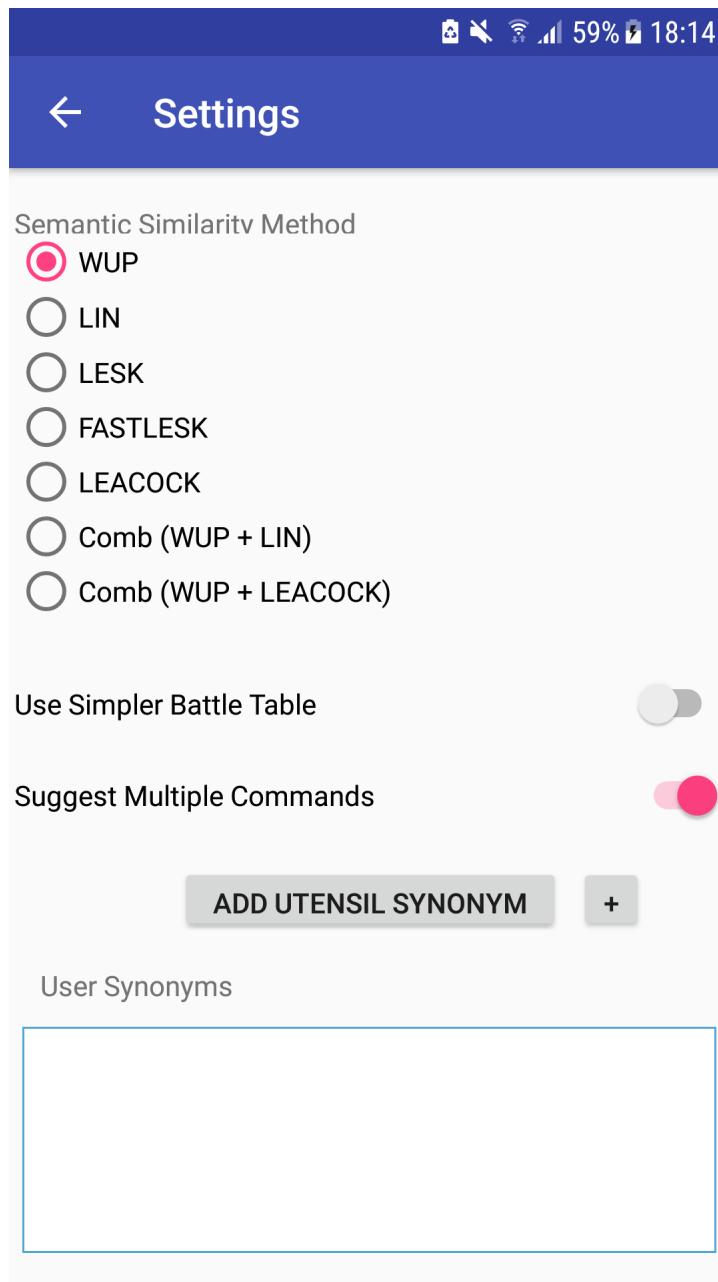
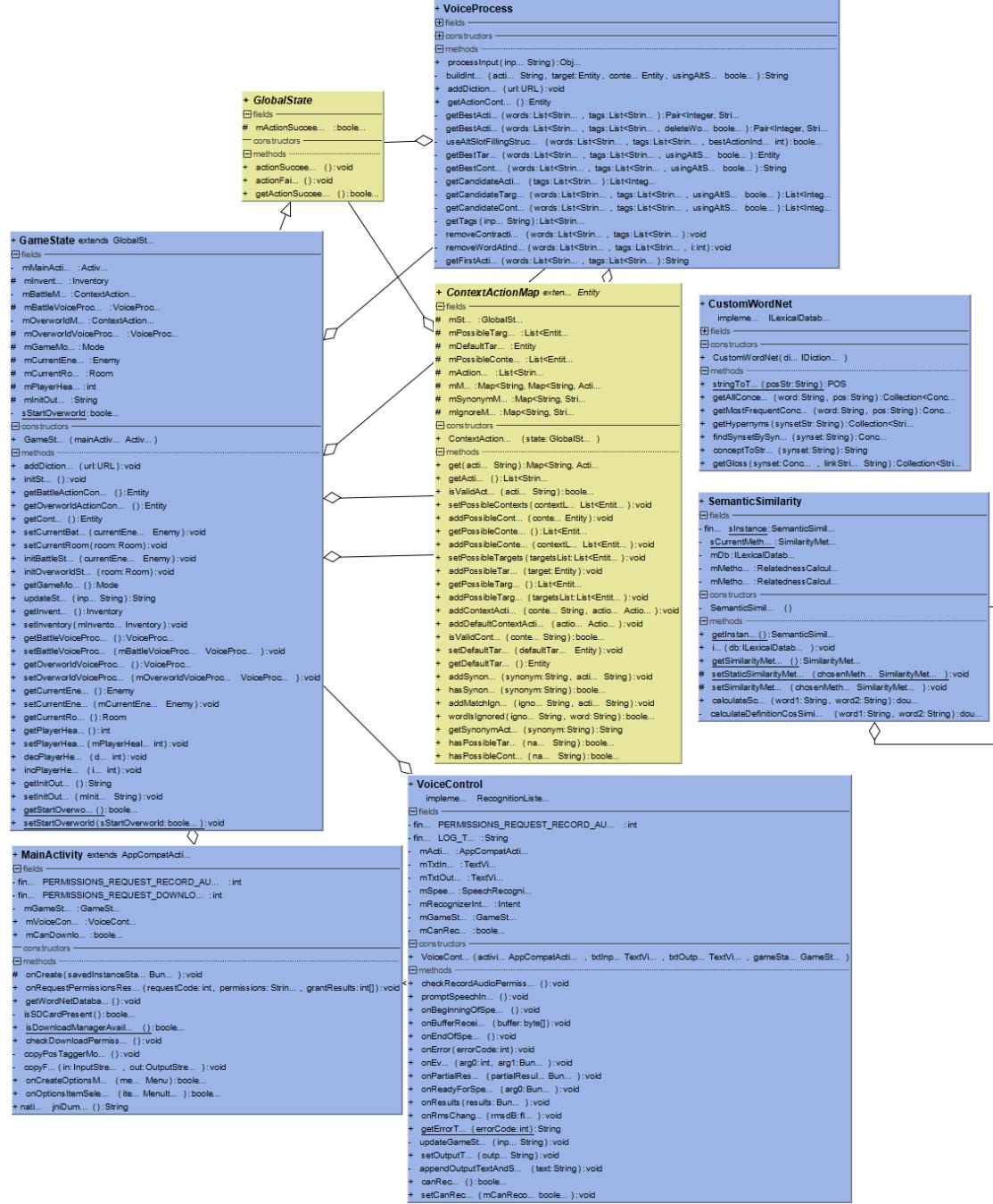


Figure 33: A snapshot of the settings activity.

B Diagrams and Examples

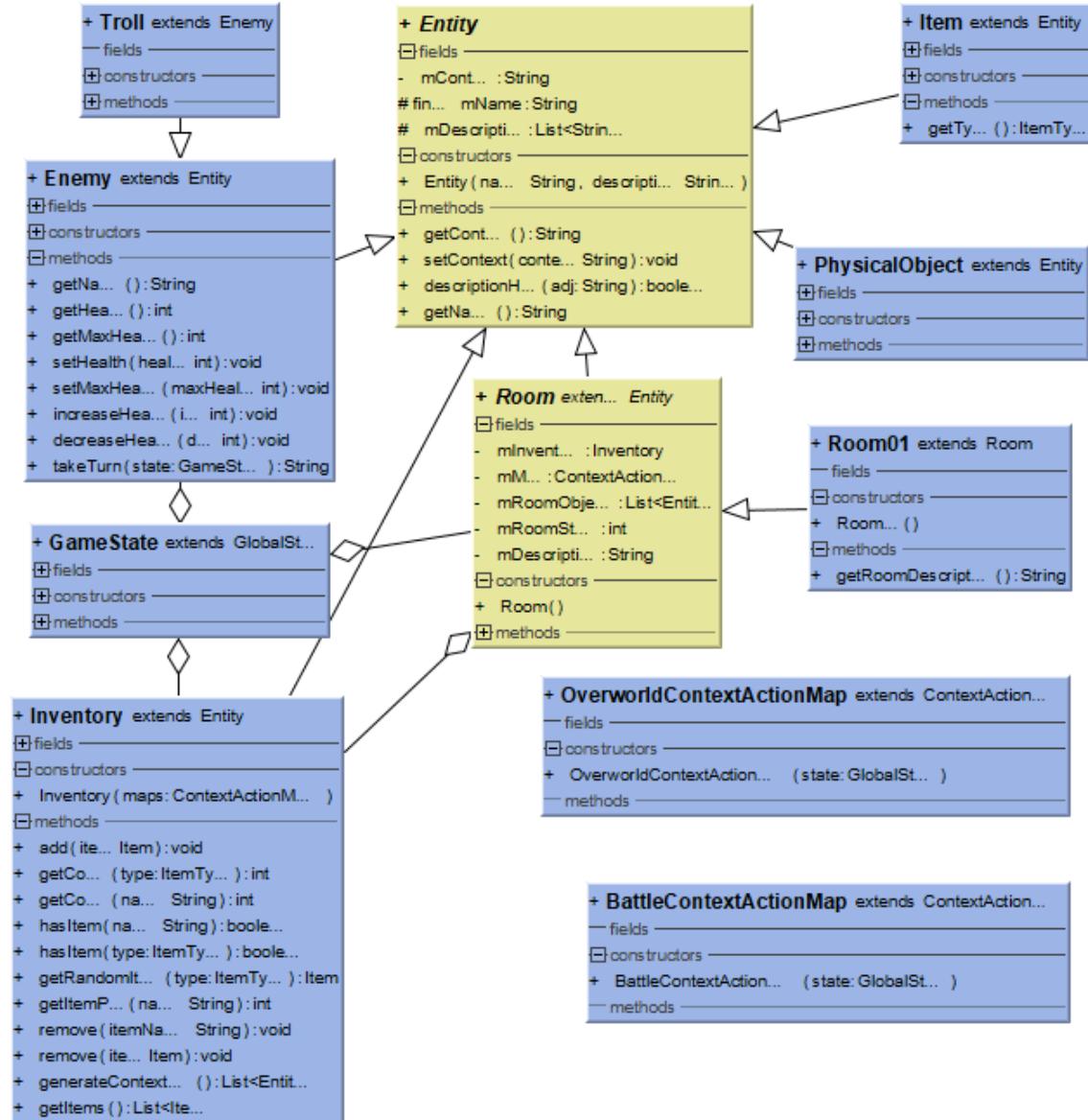
B.1 System UML Class Diagram

Below is the UML class diagram for the voice recognition system used in the Android app. This was automatically generated using the third-party *simpleUML* plugin for Android Studio [63].



B.2 Game Design UML Class Diagram

Below is the UML class diagram for the game mechanics - both the *Overworld* and *Battle* game modes. Generated using the third-party *simpleUML* Android plugin.



B.3 Battle Mode Example Commands

Below is a table of some of the example commands that can be used in the *Battle* mode of the role-playing game. Note that this is not an exhaustive list, and more commands are supported.

“attack with a hammer”
“punch the troll”
“launch an assault”
“heal with a potion”
“show my inventory”
“hit the troll with a bang”
“obliterate the troll with a sword and then recover using a potion”
“fight with a sword and hammer”
“use the sword to attack”
“what actions can i do”
“what is in my bag”
“use a heavy attack”
“launch an assault towards the troll using something sharp”

B.4 Overworld Mode Example Commands

Below is a table of some of the example commands that can be used in the *Overworld* mode of the role-playing game. Note that this is not an exhaustive list.

“look around the room”
“observe the surroundings”
“slash the door using a knife”
“break the table”
“cut the painting”
“pick up the potion”
“grab the knife”
“scratch the door with the knife and then with the sword”
“use the knife to cut the table and the door”
“use the key to open the door”
“unlock the door”
“smash the table with the hammer”
“look at the painting”

B.5 Video Conferencing Example Commands

Below is a table of some of the example commands that can be used in the video conferencing demo within the application. As above, this is not an exhaustive list.

“call fred”
“phone jane with audio only”
“call fred and jane”
“use video only to call jane”
“stop the call with fred”
“stop all of the calls”
“end all of the conversations”
“hang up”
“hang up with fred”
“mute my video”
“mute fred”
“unmute my audio”
“silence my audio”

Below is a table of all the available contacts in the demo:

First Name	Second Name
Fred	N/A
Jane	N/A
Baron	First
Baron	Second

B.6 Room Generation Example

Listing 2 shows an example of a text file for a room with various objects in it.

Listing 2: room-table.txt

```
1 You are in a large room. There is a single *door in front of you. There is a glass
2 *table in the middle of the room. A *knife is on the floor. A *fork is with
3 the knife. A *potion is on the table. Another *utensil is also on the table.
```

Using the following command:

```
./generateRoom.sh room-table.txt game-objects.obj RoomTable
```

The following Java source file will be generated:

```
1 /* AUTOMATICALLY GENERATED SOURCE FILE */
2 package com.khan.baron.voicerecrpg.game.rooms;
3 /* TODO: insert object imports */
4
5 public class RoomTable extends Room {
6     public RoomTable() {
7         super();
8         addDescription(
9             "You are in a large room.");
10        addDescriptionWithObject(
11            "There is a single door in front of you.",
12            new Door(new Troll(100), false));
13        addDescriptionWithObject(
14            "There is a glass table in the middle of the room.",
15            new GlassTable());
16        addDescriptionWithObject(
17            "A knife is on the floor.",
18            new Weapon("knife"));
19        addDescriptionWithObjectCond(
20            "A fork is with the knife.",
21            "A fork is in the room.",
22            new Weapon("fork"),
23            () -> getRoomObjectCount("knife") > 0);
24        addDescriptionWithObjectCond(
25            "A potion is on the table.",
26            "A potion is now on the floor.",
27            new Potion("potion"),
28            () -> getRoomObjectCount("table") > 0);
29        addDescriptionWithObjectCond(
30            "Another utensil is also on the table.",
31            "Another utensil is now on the floor.",
32            new Weapon("spoon"),
33            () -> getRoomObjectCount("table") > 0);
34    }
35}
```

C Testing and Evaluation

C.1 List of Android Mock Tests

Below is a list of most of the test suites containing commands that are tested for correctness in the Android application. Note that commands that are tested for incorrectness, or any tests that are not testing commands, are not shown here. Also note that some test suites in `BattleTest` are repeated in `OverworldTest` (as the modes use different `ContextActionMaps`), so are not shown again, and any commands that are tested multiple times in the same test class are only shown once. See the project's repository for the full test framework.

C.1.1 BattleTest

- `testAttackInputSuite`
 - “attack”
 - “attack with everything you have got”
 - “hit”
 - “hit the troll”
 - “punch the troll”
 - “launch an assault”
 - “charge at the troll”
 - “fight the troll”
 - “attack the troll”
 - “attack with the hammer”
 - “attack with the sledgehammer”
 - “attack the troll with the sledgehammer”
 - “use something to attack with”
 - “use a sword attack”
- `testCorrectContexts`
 - “attack the troll with a sword”
 - “attack the troll with a hammer”
- `testHealInputSuite`
 - “heal”
 - “use a potion to heal”
 - “heal with elixer”
 - “use an elixer right now before it is too late to heal”
 - “recover with an elixer”
 - “use something to heal with”
 - “heal myself with a potion”
- `testDescriptions`
 - “attack with something blunt”
 - “hit the troll with something pointy please”
 - “launch an assault towards the troll using something sharp”

- “use a heavy attack”
- testInventorySuite
 - “show my inventory”
 - “show my bag”
 - “please show the contents of my possessions that I have”
 - “look at my inventory”
- testUseSuite
 - “use the sword to attack”
 - “use something to attack with”
 - “use a potion”
 - “use something sharp”
- testLookAroundSuite
 - “look around”
 - “look at the troll”
- testShowActionsSuite
 - “show my commands”
 - “look at my actions”
- testLearningSuite
 - “strike means attack”
 - “strike the troll”
- testConfirmation
 - “obliterate the troll with the sword”
 - “kick the troll”
 - “regenerate using a potion”
 - “use something to regenerate with”
 - “attack with a bang”
- testMultipleCommandsSuite
 - “obliterate the troll with the sword and then regenerate using a potion”
 - “attack with a sword and hammer”
 - “attack with a hammer and with a sword”
- testSentenceMatching
 - “what are my actions”
 - “what actions can i do”
 - “what is in my bag”

C.1.2 OverworldTest

- testLookSuite
 - “look at my inventory”
 - “look inside the bag”
 - “look around”
 - “look at the surroundings”
 - “observe the room”
- testPickUpSuite
 - “pick up the knife”
 - “grab the potion”
- testRoom01CutSuite
 - “cut the painting”
 - “cut the table using the knife”
 - “slash the door using the knife”
- testBreakSuite
 - “break the table”
 - “break the door”
- testOpenSuite
 - “open the table”
 - “open the door”
- testUseSuite
 - “use your hand to pick up the knife”
 - “use the knife to slash the door”
 - “use the knife to cut the glass table”
 - “use the knife to cut the painting”
 - “use the key to open the door”
- testMultipleCommandsSuite
 - “grab the knife and cut the painting then unlock the door”
- testMultipleCommandsAndMultipleSuggestionsSuite
 - “grab the utensil and then cut the painting with the knife”
- testMultipleTargetsSuite
 - “pick up the knife and the potion”
 - “break the table and the door”
- testMultipleTargetsWithContextSuite
 - “use the knife to cut the table and the door”
 - “cut the table with the knife and the door with the sword”
 - “cut the door with the knife and then with the sword”
- testMultipleSynonymMapping
 - “grab the utensil”
 - “pick up the utensil”

C.1.3 CallTest

- testCallingSuite
 - “call fred”
 - “call jane”
 - “phone jane”
 - “contact fred”
 - “call fred with video”
 - “use audio to phone jane”
 - “call fred and jane”
- testStopCallsSuite
 - “stop all the calls”
 - “end all of the conversations”
 - “stop call with jane”
 - “finish call with jane”
 - “hang up”
 - “hang up with fred”
 - “stop calling fred”
- testMutedSuite
 - “mute my video”
 - “mute my audio”
 - “unmute my video”
 - “unmute my audio”
 - “silence my audio”
 - “increase my audio”
 - “mute fred”
 - “unmute fred”

Please see the project’s repository for the full test framework.

C.2 Actions and Contexts used in Evaluation

Below are the `ContextActionMap` ‘tables’ used in the evaluation Java project for the three domains. See the project’s repository for the full classes.

Table 6: GameContextActionMap

	attack	heal	move	defend
default	Attack	Heal	Move	Defend
weapon	AtkWeapon			
potion		HealPotion		

Table 7: CallContextActionMap

	phone	stop	mute
default	PhoneContact	StopCall	Mute
video	PhoneContactVideo		
audio	PhoneContactAudio		
contact		StopCallContact	

Table 8: CookingContextActionMap

	make	stir	boil	pour	serve
default	Make	Stir	Boil	Pour	Serve
spoon	Make	StirSpoon			
cooker	Make		BoilCooker		
food	MakeFood				

C.3 Evaluation Results for WUP

Tables 9, 10 and 11 below show the results of the evaluation project using the Wu and Palmer (WUP) method.

Table 9: Game Tests

Test	Expected	Result
attack	ATTACK	PASS
charge	ATTACK	PASS
hit	ATTACK	PASS
tackle	ATTACK	PASS
fight	ATTACK	FAIL
assault	ATTACK	PASS
battle	ATTACK	PASS
launch an assault	ATTACK	PASS
attack with a sword	ATTACK_WEAPON	PASS
attack with something sharp	ATTACK_WEAPON	PASS
attack with something pointy	ATTACK_WEAPON	FAIL
attack with something long	ATTACK_WEAPON	FAIL
attack with something metallic	ATTACK_WEAPON	PASS
heal	HEAL	PASS
recover	HEAL	PASS
regenerate	HEAL	PASS
rest	HEAL	FAIL
restore	HEAL	PASS
defend	DEFEND	PASS
guard	DEFEND	PASS
safeguard	DEFEND	FAIL
shield	DEFEND	PASS
heal with a potion	HEAL_POTION	PASS
heal with an elixer	HEAL_POTION	FAIL
heal with an healing drink	HEAL_POTION	PASS
move forwards	MOVE_FORWARD	PASS
move straight	MOVE_FORWARD	FAIL
move backwards	MOVE_BACKWARD	PASS
move in reverse	MOVE_BACKWARD	FAIL
move ahead	MOVE_FORWARD	FAIL
continue forwards	MOVE_FORWARD	PASS
run forwards	MOVE_FORWARD	PASS
dash forwards	MOVE_FORWARD	PASS
Score:	72.73%	

Table 10: Video Conferencing Tests

Test	Expected	Result
phone	PHONE	PASS
phone fred	PHONE_FRED	PASS
phone jane	PHONE_JANE	PASS
ring jane	PHONE_JANE	PASS
phone jane with video	PHONE_JANE_VIDEO	PASS
phone jane with webcam	PHONE_JANE_VIDEO	PASS
use webcam to call jane	PHONE_JANE_VIDEO	FAIL
call jane with audio	PHONE_JANE_AUDIO	PASS
phone jane with sound	PHONE_JANE_AUDIO	PASS
contact jane with video	PHONE_JANE_VIDEO	FAIL
stop call	STOP	PASS
stop call with fred	STOP_FRED	PASS
end call with fred	STOP_FRED	PASS
close	STOP	FAIL
finish call with fred	STOP_FRED	PASS
halt call with fred	STOP_FRED	PASS
mute video	MUTE_VIDEO	PASS
mute screen	MUTE_VIDEO	PASS
mute jane	MUTE_JANE	PASS
silence jane	MUTE_JANE	FAIL
Score:	80.00%	

Table 11: Cooking Tests

Test	Expected	Result
make egg	MAKE_EGG	PASS
make eggs	MAKE_EGG	PASS
make soup	MAKE_SOUP	PASS
create soup	MAKE_SOUP	PASS
produce soup	MAKE_SOUP	PASS
cook egg	MAKE_EGG	PASS
make omelette	MAKE_EGG	FAIL
make something liquid	MAKE_SOUP	FAIL
boil egg	BOIL_EGG	PASS
heat egg	BOIL_EGG	PASS
boil soup with cooker	BOIL_SOUP_COOKER	PASS
boil soup with boiler	BOIL_SOUP_COOKER	PASS
boil soup with stove	BOIL_SOUP_COOKER	FAIL
heat pottage with stove	BOIL_SOUP_COOKER	FAIL
stir soup	STIR_SOUP	PASS
mix soup	STIR_SOUP	FAIL
blend soup	STIR_SOUP	FAIL
stir chowder	STIR_SOUP	FAIL
stir pottage	STIR_SOUP	FAIL
stir soup with spoon	STIR_SOUP_SPOON	PASS
stir soup with cutlery	STIR_SOUP_SPOON	PASS
stir soup with utensil	STIR_SOUP_SPOON	PASS
stir soup with tablespoon	STIR_SOUP_SPOON	PASS
stir soup with teaspoon	STIR_SOUP_SPOON	PASS
stir soup with soupspoon	STIR_SOUP_SPOON	PASS
blend chowder using soupspoon	STIR_SOUP_SPOON	PASS
pour buttermilk	POOR_MILK	PASS
discharge milk	POOR_MILK	FAIL
present eggs	SERVE_EGG	FAIL
deliver soup	SERVE_SOUP	PASS
Score:	63.33%	

C.4 Performance Measurement Commands

Below is the list of commands that were executed during the profiling in Figure 29 on page 63 (in order of execution):

- “pick up the knife and potion”
- “slash the painting using the knife”
- “what do i have in my inventory”
- “open the door using the key”
- “obliterate the troll using the knife”
- “yes”
- “attack the troll with something sharp”
- “hit the troll with something heavy”

C.5 Comparison of Usability

This section outlines a comparison between the usability of the voice recognition system implemented in this project, and other naïve/traditional systems. The example below is the same one used in the `README` file of the library's repository¹²

If the user wished to map the following commands to an intent which attacked with a sword:

- “attack with a sword”
- “hit with something sharp”
- “use a blade to fight”
- “launch an assault using the sword”
- “obliterate the enemy with a pointy weapon”

Or commands which mapped to healing with a potion:

- “heal”
- “recover”
- “rest”
- “heal with a potion”
- “regenerate using a healing drink”

A naïve approach would be to use a switch statement¹³ and add every possible variation of the command.

```
1 switch(input) {
2     case "attack with a sword":
3     case "hit with something sharp":
4         //...
5     case "obliterate the enemy with a pointy weapon":
6         attackWithWeapon();
7         break;
8     case "heal":
9     case "recover":
10        //...
11    case "rest":
12        heal();
13        break;
14    case "heal with a potion":
15    case "recover with a potion":
16        //...
17    case "regenerate using an elixer":
18        healWithPotion();
19        break;
20    //...
21 }
```

¹²<https://github.com/BaronKhan/voice-commands-with-wordnet/blob/master/README.md>

¹³Strings can be evaluated with a `switch` statement since Java 7.

However, this misses out on other variations of the commands that would be acceptable (e.g. “fight using a sword”).

Another method is to create a regex-like expression to describe the structure of acceptable inputs:

```
1 ("attack" | "hit") . "with" . ["a"] . ("sword" | "blade")
```

However, this can also become quite complex, and still doesn't cover a wide variety of commands.

```
1 ("attack" | "hit" | "obliterate" | ("launch" . "an" . "assault")) . ("with" | "using") . ["a"] . ("sword" | "blade" | ("something" . ("pointy" | "sharp")))
```

With the system proposed in this project, the developer only has to create the following table of mappings in order to map all of the above commands and much more.

	attack	heal	move
default	Attack	Heal	Move
weapon	AtkWeapon		
potion		HealPotion	

Therefore, the number of lines needed to add the above voice commands and can be considered reduced using this method compared to the other two methods.

C.6 Other Survey Results

This section presents the other results from the anonymous survey given to developers. It should be noted that the same size for the survey is only eight participants, and the survey is intended to be supplementary information.

As a developer that wants to add voice commands to an application, which of the following are the most important to you? (Select up to two.)

8 responses

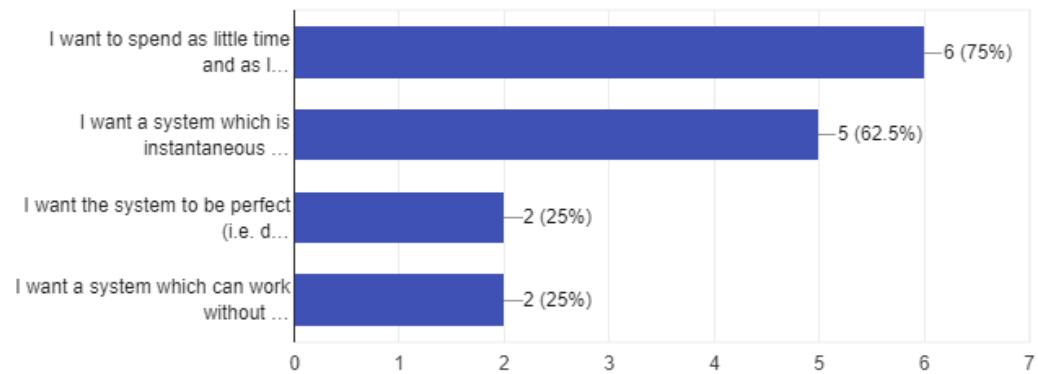


Figure 34: A bar chart showing the responses when asking developers what are the most important aspects of a system for adding voice commands to an application.

As a developer, when choosing a system for adding voice commands to an application (assuming you have the use...s it), which system would you prefer?

8 responses

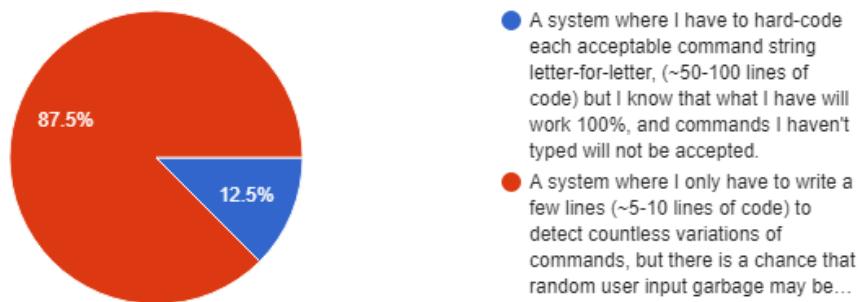


Figure 35: A pie chart showing the responses to a question asking a developer about which type of voice recognition system they would prefer.

Read the following README (link below) for a Java library that adds voice commands to a project. Is it clear what the library is aiming to do?

8 responses

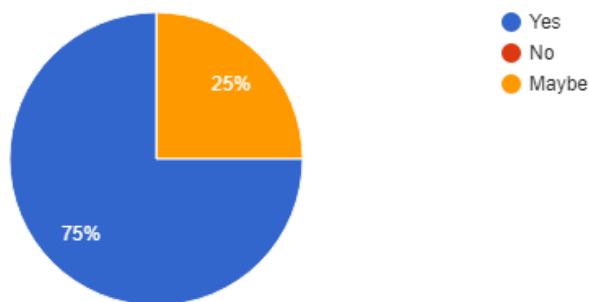


Figure 36: A pie chart showing the responses to a question asking a developer if they understood the purpose of the voice recognition system in making it easier to add voice commands.

C.7 Room Generation Evaluation Results

The following is a text file for a room description (with each sentence contain an object relation), followed by the methods in the constructor from the output of the program:

“There is a *table in the room. A *knife is on the table. There is a *spoon on top of the table. A *fork lays on the table. There is a *chair under the table. A *letter sits upon the table. A *key is placed above the letter. A *sword is next to the table. A *hammer is with the sword. There is a *potion is underneath the table. An *elixer is beside the potion.”

```

1 addDescriptionWithObject(
2     "There is a table in the room.",
3     new GlassTable());
4 addDescriptionWithObjectCond(
5     "A knife is on the table.",
6     "A knife is now on the floor.",
7     new Weapon("knife"),
8     () -> getRoomObjectCount("table") > 0);
9 addDescriptionWithObject(
10    "There is a spoon on top of the table.",
11    new Weapon("spoon"));
12 addDescriptionWithObject(
13    "A fork lays on the table.",
14    new Weapon("fork"));
15 addDescriptionWithObject(
16    "There is a chair under the table.",
17    new Chair());
18 addDescriptionWithObjectCond(
19     "A letter sits upon the table.",
20     "A letter is now on the floor.",
21     new Item("letter", Item.ItemType.ITEM_KEY, "document"),
22     () -> getRoomObjectCount("table") > 0);
23 addDescriptionWithObjectCond(
24     "A key is placed above the letter.",
25     "A key is in the room.",
26     new Item("key", Item.ItemType.ITEM_KEY, "access"),
27     () -> getRoomObjectCount("letter") > 0);
28 addDescriptionWithObjectCond(
29     "A sword is next to the table.",
30     "A sword is in the room.",
31     new Weapon("sword"),
32     () -> getRoomObjectCount("table") > 0);
33 addDescriptionWithObjectCond(
34     "A hammer is with the sword.",
35     "A hammer is in the room.",
36     new Weapon("hammer"),
37     () -> getRoomObjectCount("sword") > 0);
38 addDescriptionWithObjectCond(
39     "There is a potion is underneath the table.",
40     "A potion is in the room.",
41     new Potion("potion"),
42     () -> getRoomObjectCount("table") > 0);
43 addDescriptionWithObjectCond(
44     "An elixer is beside the potion.",
45     "An elixer is in the room.",
46     new Potion("elixer"),
47     () -> getRoomObjectCount("potion") > 0);

```

Out of the 10 relations in the description, 7 of them are extracted by the program.

D Supplementary Materials

D.1 Skyrim Kinect Command List

Below is a list of all 200 commands that can be used in the video game, The Elder Scrolls V: Skyrim, when a Microsoft Kinect peripheral is connected [5].



FAVORITES MENU & HOTKEY EQUIPPING

Only usable in the Favorites menu.

ASSIGN <hotkey command>	sets the selected item to the spoken hotkey. <hotkey command> can be any of the following:		
• HEALTH POTION	• DAGGER	• BATTLEAXE	• BOUND WEAPON
• MAGICKA POTION	• BOW	• WARHAMMER	• SUMMON SPELL
• STAMINA POTION	• SHIELD	• FIRE SPELL	• ARMOR SPELL
• POISON	• DUAL WIELD LEFT	• FROST SPELL	• CALM SPELL
• SWORD	• DUAL WIELD RIGHT	• LIGHTNING SPELL	• FRENZY SPELL
• MACE	• SOUL TRAP	• WARD SPELL	• HEALING SPELL
• AXE	• GREATSWORD	• RITUAL SPELL	• LIGHT

Only usable during main gameplay.

EQUIP <hotkey command>	EQUIP LEFT <hotkey command>	EQUIP RIGHT <hotkey command>
Equip item in the default hand	Equip item in the Left hand, if possible.	Equip item in the Right hand, if possible.
EQUIP DUAL <hotkey command>	EQUIP SWORD AND SHIELD / EQUIP MACE AND SHIELD / EQUIP AXE AND SHIELD / EQUIP DAGGER AND SHIELD	EQUIP DUAL WEAPONS
Equip the item in both hands, if possible.	Equip the item assigned to the Sword / Mace / Axe / Dagger hotkey in the right hand, and the item assigned to the Shield hotkey in the left hand.	

MAGIC MENU

After opening MAGIC, the following opens the menu to the corresponding category:

FAVORITES	ALTERATION	CONJURATION	POWERS
ALL	ILLUSION	RESTORATION	ACTIVE EFFECTS / EFFECTS
	DESTRUCTION	SHOUTS	

CHARACTER MENU

After opening CHARACTER, the following opens the menu to the corresponding category:

ITEMS	MAGIC	SKILLS	MAP
-------	-------	--------	-----

MAP MENU

After opening MAP, the following centers the camera on the corresponding locations:

WINDHELM	MORTHAL	RIFTEN
FALKREATH	SOLITUDE	WHITERUN
DAWNSTAR	MARKARTH	WINTERHOLD

PLAYER / WHERE AM I WAYPOINT QUEST MARKER / QUEST

Centers camera on you. Centers camera on marker set by you, if applicable. Centers camera on active quest target. Saying it multiple times will cycle through targets.

ITEMS MENU

After opening ITEMS, the following commands open the corresponding categories:

• FAVORITES	• POTIONS	• BOOKS	
• ALL	• SCROLLS	• KEYS	
• WEAPONS	• FOOD	• MISCELLANEOUS	
• APPAREL / ARMOR	• INGREDIENTS		

After opening any of the above sub-menu categories, you can use the following to sort your items.

SORT BY NAME	SORT BY WEIGHT	SORT BY VALUE
Sorts the item list by name, increasing. If the list is already sorted by name increasingly, sorts it decreasingly.	Sorts the item list by weight, decreasing. If the list is already sorted by weight decreasingly, sorts it increasingly.	Sorts the item list by value, decreasing. If the list is already sorted by value decreasingly, sorts it increasingly.

CLOSE MENU

Closes the menu

SKILLS MENU

After opening SKILLS, use the following commands to center the camera on that skill's perk tree:

ONE HANDED	HEAVY ARMOR	ALCHEMY	ILLUSION
TWO HANDED	LIGHT ARMOR	SPEECH	RESTORATION
ARCHERY	PICKPOCKET	ALTERATION	ENCHANTING
BLOCK	LOCKPICKING	CONJURATION	
SMITHING	SNEAK	DESTRUCTION	

D.2 Penn Treebank Tagset

Below is a table of the 32 tags from the Penn Treebank tagset, used for POS tagging. The 12 remaining tags for punctuation and other symbols have been left out [64].

CC	Coordinating conjunction	TO	to
CD	Cardinal number	UH	Interjection
DT	Determiner	VB	Verb, base form
EX	Existential there	VBD	Verb, past tense
FW	Foreign word	VBG	Verb, gerund or present participle
IN	Preposition or subordinating conjunction	VBN	Verb, past participle
PRP\$	Possessive pronoun	NNS	Noun, plural
RB	Adverb	NNP	Proper noun, singular
RBR	Adverb, comparative	NNPS	Proper noun, plural
RBS	Adverb, superlative	PDT	Predeterminer
RP	Particle	POS	Possessive ending
SYM	Symbol	PRP	Personal pronoun
JJ	Adjective	VBP	Verb, non-3rd person singular present
JJR	Adjective, comparative	VBZ	Verb, 3rd person singular present
JJS	Adjective, superlative	WDT	Wh-determiner
LS	List item marker	WP	Wh-pronoun
MD	Modal	WP\$	Possessive wh-pronoun
NN	Noun, singular or mass	WRB	Wh-adverb

D.3 Pseudocode for ProcessInput()

Listing 3 shows the pseudocode for the method, with most of the minor details removed. For the full Java method, refer to the source code in the project's repository.

Listing 3: Pseudocode for processInput()

```
1 def processInputPseudocode(input):
2     checkIfWordNetIsLoaded()
3
4     if expectingReply:
5         return processPendingIntent(input)
6
7     //Check for action replies
8     if currentAction.wantsReply():
9         return currentAction.processReply(input)
10
11    words, tags = tokeniseAndTag(input)
12
13    //Check for learning phrase ("--- means ---")
14    if words.size() == 3 and words.contains("means"):
15        addSynonym(words[0], words[2])
16
17    //Parse the input
18    actionIndex, chosenAction = getBestAction(words, tag)
19
20    if isValidAction(actionStr):
21        chosenTarget = getBestTarget(words, tags)
22        chosenContext = getBestContext(words, tags)
23
24        if isValidContext(bestContext):
25            return action.execute(state, chosenTarget)
26
27    else if foundAnotherTargetOrContext:
28        return previousAction.execute()
29    else:
30        return performSentenceMatching(input)
```

References

- [1] M. Tabini, “Inside Siri’s brain: The challenges of extending Apple’s virtual assistant,” 2013, April 8. [Online]. Available: <https://www.macworld.com/article/2033073/inside-siris-brain-the-challenges-of-extending-apples-virtual-assistant.html>
- [2] Indomitus Games, “In Verbis Virtus: voice controlled fantasy game,” Nov 2013. [Online]. Available: <https://www.indiegogo.com/projects/in-verbis-virtus-voice-controlled-fantasy-game>
- [3] ——, “In Verbis Virtus.” [Online]. Available: <http://www.indomitusgames.com/in-verbis-virtus>
- [4] Bethesda Softworks LLC, “The Elders Scrolls Official Site,” 2018, January 1. [Online]. Available: <https://elderscrolls.bethesda.net/en/skyrim>
- [5] NowGamer, “Skyrim Kinect: Full List Of 200 Voice Commands,” 2012, May 1. [Online]. Available: <https://www.nowgamer.com/skyrim-kinect-full-list-of-200-voice-commands/>
- [6] Ubisoft, “Star Trek™: Bridge Crew on Steam,” 2017, May 30. [Online]. Available: https://store.steampowered.com/app/527100/Star_Trek_Bridge_Crew/
- [7] ——, “Star Trek Bridge Crew,” 2017. [Online]. Available: <https://www.ubisoft.com/en-gb/game/star-trek-bridge-crew/>
- [8] C. Watters, “Star Trek: Bridge Crew Integrates IBM Watson For Voice Commands,” 2017, June 22. [Online]. Available: <https://news.ubisoft.com/article/star-trek-bridge-crew-integrates-ibm-watson-voice-commands>
- [9] developerWorks, “Star Trek: Bridge Crew now live with IBM Watson Speech Technology,” 2017, June 22. [Online]. Available: <https://developer.ibm.com/tv/dwnewsblast-watson-star-trek/>
- [10] R. Lammle, “Eaten by a Grue: A Brief History of Zork,” June 2014. [Online]. Available: <http://mentalfloss.com/article/29885/eaten-grue-brief-history-zork>
- [11] J. Foster, M. Thompson, and M. Loffler, “Classic Zork on Alexa,” Dec 2016. [Online]. Available: <https://www.hackster.io/devops-dungeoneers/classic-zork-494ff1>
- [12] D. Lebling and M. Blank, *Zork Trilogy Instruction Manual*. InfoCom, 1984. [Online]. Available: <http://infodoc.plover.net/manuals/zork1.pdf>
- [13] VoiceTechGroup, “tazti — Speech Recognition for PC Games,” 2018, January 1. [Online]. Available: <https://www.tazti.com/speech-recognition-software-for-pc-games.html>
- [14] Epic Wiki, “Speech Recognition Plugin,” 2018, March 17. [Online]. Available: https://wiki.unrealengine.com/Speech_Recognition_Plugin

- [15] CMU Sphinx, “PocketSphinx,” 2015. [Online]. Available: <https://github.com/cmusphinx/pocketsphinx>
- [16] SoundHound, “Houndify,” 2017. [Online]. Available: <https://www.houndify.com/>
- [17] ——, “The ClientMatch JSON Format,” 2017. [Online]. Available: <https://docs.houndify.com/reference/ClientMatch>
- [18] IBM, “Conversation - IBM Cloud,” 2017, December 12. [Online]. Available: <https://console.bluemix.net/catalog/services/conversation>
- [19] “IBM Cloud Docs — Conversation — About,” 2018, January 5. [Online]. Available: <https://console.bluemix.net/docs/services/conversation/index.html>
- [20] Zach Walchuk, “Build a chatbot in ten minutes with Watson,” Dec 2016. [Online]. Available: <https://www.ibm.com/blogs/watson/2016/12/build-chat-bot/>
- [21] IBM, “Conversation Demo,” 2017, December 18. [Online]. Available: <https://conversation-demo.ng.bluemix.net/>
- [22] M. McCord, J. Murdock, and B. Boguraev, “Deep parsing in Watson,” *IBM Journal of Research and Development*, vol. 56, no. 3.4, pp. 3:1–3:15, Apr 2012. [Online]. Available: <https://ieeexplore.ieee.org/document/6177729/>
- [23] “Dialogflow.” [Online]. Available: <https://dialogflow.com/>
- [24] N. Indurkha and F. Damerau, *Classical Approaches to Natural Language Processing*, 2 ed., ser. Handbook of Natural Language Processing. United States: Chapman and Hall, 2010, pp. 4–6.
- [25] C. Trim, “The Art of Tokenization,” Jan 2013. [Online]. Available: <https://ibm.com/developerworks/community/blogs/nlp/entry/tokenization>
- [26] N. Indurkha and F. Damerau, *Classical Approaches to Natural Language Processing*, 2 ed., ser. Handbook of Natural Language Processing. United States: Chapman and Hall, 2010, p. 61.
- [27] ——, *Part-of-Speech Tagging*, 2 ed., ser. Handbook of Natural Language Processing. United States: Chapman and Hall, 2010, pp. 205–206.
- [28] Dictionary.com, “The Dictionary of American Slang,” 2018, January 18. [Online]. Available: <http://www.dictionary.com/browse/can>
- [29] D. Jurafsky and J. Martin, *Dialog Systems and Chatbots*, ser. Speech and Language Processing, 2017. [Online]. Available: <https://web.stanford.edu/~jurafsky/slp3/29.pdf>
- [30] WikiDiff, “Synonym vs Hypernym - What’s the Difference?” [Online]. Available: <https://wikidiff.com/synonym/hypernym>

- [31] Princeton University, “Dog — WordNet Search - 3.1,” 2018, January 18. [Online]. Available: <http://wordnetweb.princeton.edu/>
- [32] B. McInnes and T. Pedersen, “Evaluating measures of semantic similarity and relatedness to disambiguate terms in biomedical text,” *Journal of Biomedical Informatics*, 2013, Sep 2013. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1532046413001238>
- [33] S. Gole, “Words similarity/relatedness using WuPalmer Algorithm - Sagar Gole’s Blog,” Jun 2015. [Online]. Available: <http://blog.thedigitalgroup.com/sagarg/2015/06/10/words-similarityrelatedness-using-wupalmer-algorithm/>
- [34] Hideki Shima, “Juice, Water — WS4J Demo,” 2018, January 18. [Online]. Available: <http://ws4jdemo.appspot.com>
- [35] Princeton University, “About WordNet,” 2013, January 9 2013. [Online]. Available: <http://wordnet.princeton.edu>
- [36] A. Fader, S. Soderland, and O. Etzioni, “Identifying relations for open information extraction,” in *Proceedings of the Conference of Empirical Methods in Natural Language Processing (EMNLP ’11)*, Edinburgh, Scotland, UK, July 27-31 2011.
- [37] M. Iyyer, J. Boyd-Graber, and H. Daume, “Generating sentences from semantic vector space representations,” *NIPS Workshop on Learning Semantics*, 2014. [Online]. Available: https://www.cs.umd.edu/~miyyer/pubs/2014_nips_generation.pdf
- [38] L. White, R. Togneri, W. Liu, and M. Bennamoun, “Modelling Sentence Generation from Sum of Word Embedding Vectors as a Mixed Integer Programming Problem,” in *2016 IEEE 16th International Conference on Data Mining Workshops (ICDMW)*, 2016, pp. 770–777, ID: 1.
- [39] M. Bozon, “World Debut: Scribblenauts,” 2008, December 5. [Online]. Available: <http://uk.ign.com/articles/2008/12/05/world-debut-scribblenauts>
- [40] G. Escofet, “Mobile games market forecast report: 2017–22,” Ovum, Tech. Rep., Dec 2017. [Online]. Available: <https://ovum.informa.com/resources/product-content/ces003-000052>
- [41] Android Developers, “SpeechRecognizer.” [Online]. Available: <https://developer.android.com/reference/android/speech/SpeechRecognizer.html>
- [42] V. Arputharaj, “How to enable offline Speech To Text in Android,” May 2015. [Online]. Available: <https://stackandroid.com/tutorial/how-to-enable-offline-speech-to-text-in-android/>
- [43] The Stanford Natural Language Processing Group, “Stanford Log-linear Part-Of-Speech Tagger.” [Online]. Available: <https://nlp.stanford.edu/software/tagger.html>
- [44] J. McNamara, “Sentence subjects.” [Online]. Available: <http://grammar.ccc.commnet.edu/grammar/subjects.htm>

- [45] M. Lesk, “Automatic sense disambiguation using machine readable dictionaries: How to tell a pine cone from an ice cream cone.” [Online]. Available: citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.178.2744&rep=rep1&type=pdf
- [46] P. T. and B. S, “Extended gloss overlaps as a measure of semantic relatedness,” in *Proceedings of the 18th international joint conference on AI*, 2003, pp. 805–810. [Online]. Available: <http://www.d.umn.edu/~tpederse/Pubs/ijcai03.pdf>
- [47] R. Nordquist, “Definition and examples of meronyms and holonyms,” Mar 2018. [Online]. Available: <https://www.thoughtco.com/what-is-a-meronym-1691308>
- [48] M. A. Finlayson, “Java Libraries for Accessing the Princeton Wordnet: Comparison and Evaluation,” in *Proceedings of the 7th International Global WordNet Conference*, 2014, pp. 78–85. [Online]. Available: <https://projects.csail.mit.edu/jwi/download.php?f=finlayson.2014.procgwc.7.78.pdf>
- [49] H. Shima, “WS4J,” Apr 2013. [Online]. Available: <https://code.google.com/archive/p/ws4j/>
- [50] C. S. Perone, “Machine learning: Cosine similarity for vector space,” March 2018. [Online]. Available: <http://blog.christianperone.com/2013/09/machine-learning-cosine-similarity-for-vector-space-models-part-iii/>
- [51] G. Sidorov, A. Gelbukh, H. Gomez-Adorno, and D. Pinto, “Soft similarity and soft cosine measure: Similarity of features in vector space model,” *Computacion y Sistemas*, vol. 18, no. 3, pp. 491–504, 2014. [Online]. Available: <http://www.cys.cic.ipn.mx/ojs/index.php/CyS/article/view/2043>
- [52] JetBrains, “Fernflower,” 2017. [Online]. Available: <https://github.com/JetBrains/intellij-community/tree/master/plugins/java-decompiler/engine>
- [53] National Institute of Information and Communications Technology, “Japanese Wordnet,” 2006. [Online]. Available: <http://compling.hss.ntu.edu.sg/wnja/index.en.html>
- [54] Oracle, “BooleanSupplier (Java Platform SE 8).” [Online]. Available: <https://docs.oracle.com/javase/8/docs/api/java/util/function/BooleanSupplier.html>
- [55] J. Nielson, “Response times: The three important limits,” Jan 1993. [Online]. Available: <https://www.nngroup.com/articles/response-times-3-important-limits/>
- [56] DeviceSpecifications, “Samsung Galaxy S7 Exynos - Specifications.” [Online]. Available: <https://www.devicespecifications.com/en/model/63063a47>
- [57] SocialCompare, “Raspberry Pi 3,” 2018, March 14. [Online]. Available: <http://socialcompare.com/en/review/raspberry-pi-3>
- [58] Grammarly, “Compound words: Open, closed or hyphenated?” [Online]. Available: <https://www.grammarly.com/blog/open-and-closed-compound-words/>

- [59] C. F. Meyer, *Typological Classification of Languages*, International Student ed., ser. Introducing English Linguistics. Cambridge, UK: Cambridge University Press, 2010, p. 36. [Online]. Available: <https://books.google.co.uk/books?id=MWbrvUiYzSkC&pg=PA36>
- [60] S. Chong, “A Complete Multilingual WordNet List by Language,” 2015, February 6. [Online]. Available: <https://www.certifiedchinesetranslation.com/openaccess/WordNet/>
- [61] E. Bamberg, “Parallel stream processing in Java 8 - performance of sequential vs parallel stream processing,” Jan 2016. [Online]. Available: <https://blog.io.de/2016/01/22/parallel-stream-processing-in-java-8-performance-of-sequential-vs-parallel-stream-processing/>
- [62] KITT.AI, “Snowboy Hotword Detection,” Feb 2018. [Online]. Available: <https://github.com/Kitt-AI/snowboy>
- [63] M. Schmid, “simpleUMLCE,” Apr 24 2010. [Online]. Available: <https://plugins.jetbrains.com/plugin/4946-simpleumlce>
- [64] University of Pennsylvania, “Penn Treebank P.O.S. Tags,” 2018. [Online]. Available: https://ling.upenn.edu/courses/Fall_2003/ling001/penn_treebank_pos.html