



UNIVERSITY OF OXFORD
CENTRE FOR DOCTORAL TRAINING IN
CYBER SECURITY

DPHIL THESIS

**The Security of Human-Computer
Interaction by Speech**

Mary K. Bispham

supervised by
Prof. Michael Goldsmith and Dr. Ioannis Agraftotis

June 10, 2020

Abstract

This thesis investigates the security issues associated with human-computer interaction by speech, focussing on the context of voice-controlled digital assistants. The security of human-computer interaction by speech has become increasingly important as use of voice control has become more widespread. The research questions addressed in the thesis are whether the speech interface presents particular vulnerabilities that are not relevant to other types of interfaces, and, if so, what these vulnerabilities are and how attacks exploiting them can be defended. Based on a critical review of prior work, it is argued that the speech interface does represent a new attack surface with specific security vulnerabilities that have not as yet been comprehensively studied. These vulnerabilities arise both in relation to the inherently open nature of the speech interface, as well in relation to unintended functionality in the technologies implemented in voice-controlled systems to imitate human speech and language processing.

The thesis makes three main contributions towards closing the gaps in knowledge on the security of human-computer interaction by speech identified in the review of prior work. The first contribution of the thesis is a novel taxonomy of the types of attacks that might be executed via a speech interface, representing a systemisation of knowledge in this area. The second contribution of the thesis is experimental work demonstrating new types of attacks via the speech interface that are foreshadowed in prior work, but have not been validated in practice. The experimental work develops systematic methodologies for executing attacks that hide malicious voice commands in nonsensical word sounds and in apparently unrelated utterances. The methodologies applied in these experiments involve testing both machine and human responses to such input to assess the potential for exploiting differences in machine and human perceptions to execute covert attacks. The third contribution of the thesis is proposals for the development of new defence mechanisms to counter attacks via the speech interface for which no effective defence mechanisms are currently available. These proposals include feasibility tests on the application of two existing technologies for security purposes in voice-controlled systems. The proposals for new defence mechanisms are grounded in a novel attack and defence modelling approach for analysing the security of human-computer interaction by speech that enables conceptualisation of the security of the speech interface in an inclusive framework, and facilitates a review of currently available defence mechanisms.

Declaration of Authorship

I declare that this thesis is entirely my own work, and except where otherwise stated, describes my own research.

Funding

This work was funded by a doctoral training grant from the UK Engineering and Physical Sciences Research Council (EPSRC grant number EP/P0081X/1).

Publications

M. K. Bispham, I. Agrafiotis, and M. Goldsmith, “A taxonomy of attacks via the speech interface”, *Proceedings of Third International Conference on Cyber-Technologies and Cyber-Systems*, 2018.

M. K. Bispham, I. Agrafiotis, and M. Goldsmith, “Nonsense attacks on Google Assistant and missense attacks on Amazon Alexa”, *Proceedings of International Conference on Information Systems Security and Privacy*, 2019.

M. K. Bispham, I. Agrafiotis, and M. Goldsmith, “Attack and defence modelling for attacks via the speech interface”, *Proceedings of International Conference on Information Systems Security and Privacy*, 2019.

M. K. Bispham, I. Agrafiotis, and M. Goldsmith, “The speech interface as an attack surface: An overview”, *International Journal On Advances in Security*, v 12 n 1 & 2, 2019.

(forthcoming) M. K. Bispham, A. Janse van Rensburg, I. Agrafiotis, and M. Goldsmith, “Black-box attacks via the speech interface using linguistically crafted input”, *Information Systems Security and Privacy, ICISSP 2019, Prague, Czech Republic, February 23-25, 2019, Revised Selected Papers*, 2019.

(forthcoming) M. K. Bispham, I. Agrafiotis, and M. Goldsmith, “The security of the speech interface: A modelling framework and proposals for new defence mechanisms”, *Information Systems Security and Privacy, ICISSP 2019, Prague, Czech Republic, February 23-25, 2019, Revised Selected Papers*, 2019.

Acknowledgements

My thanks are due to my supervisors, for their insights and guidance, to my fellow students, for their help and companionship, to my brother, for his invaluable support, and to my parents, for everything.

“The beginning of wisdom is to call things by their proper name.” (Confucius)

Contents

1	Introduction	9
2	Background on Human-Computer Interaction by Speech	15
2.1	Overview of Speech Dialogue Systems	15
2.2	Cloud-Based Voice-Controlled Digital Assistants	27
3	Prior Work on Attacks via the Speech Interface	30
3.1	Plain Speech	30
3.2	Inaudible Sound Injection	31
3.3	Adversarial Learning	32
3.4	Active Attacks	38
3.5	Human Factors	39
4	A Taxonomy of Attacks via the Speech Interface	41
5	New Types of Attacks via the Speech Interface	46
5.1	Summary	46
5.2	Nonsense Attacks on Google Assistant	47
5.2.1	Description and Context	47
5.2.2	Pilot Experiment	51
5.2.3	Main Experiment	60
5.2.4	Discussion	70
5.3	Missense Attacks on Amazon Alexa	73
5.3.1	Description and Context	73
5.3.2	Proof-of-Concept Study	77
5.3.3	Pilot Experiment	83
5.3.4	Main Experiment	85
5.3.5	Discussion	95
6	Attack and Defence Modelling in the Context of Human-Computer Interaction by Speech	98
6.1	Attack Scenarios	98
6.2	Modelling Technique	100
6.3	The OODA Loop Model in the Context of the Speech Interface	102
6.3.1	Attacks	102

6.3.2	Defences	103
7	Proposals for the Development of New Defence Mechanisms	115
7.1	Scope of Proposals	115
7.2	High-Level Defence Concept	117
7.3	Defences against Attacks on Speech Recognition	120
7.3.1	FlexSR Feasibility Tests - Noise Attacks	122
7.3.2	FlexSR Feasibility Tests - Nonsense Attacks	124
7.4	Defences against Attacks on Natural Language Understanding	129
7.4.1	Google Translate Feasibility Tests - Missense Attacks	130
7.5	Defences against Attacks on Dialogue Management	132
7.6	Recovery from False Positives	133
8	Conclusions and Future Work	135
8.1	Limitations and Future Work	137
	Bibliography	139
A	Nonsense Attacks on Google Assistant	157
A.1	Nonsense Words Sets	157
A.2	Pilot Experiment Results	167
A.2.1	Audio File Input Results (Successes)	167
A.2.2	Over-the-Air Results	175
A.2.3	Human Comprehensibility Results	178
A.2.4	Retest Results	184
A.3	Main Experiment Results	190
A.3.1	Audio File Input Results	190
A.3.2	Over-the-Air Results	200
A.3.3	Human Comprehensibility Results	201
B	Missense Attacks on Amazon Alexa	212
B.1	Pilot Experiment - Training Dataset	212
B.2	Experiment - Unsuccessful Adversarial Utterances	222

Chapter 1

Introduction

The use of speech interfaces for human-computer interaction is becoming more widespread, particularly in the form of voice-controlled digital assistants.¹ Such ‘assistants’ are intended to act as brokers between users and the vastly complex, often intimidating cyber world. Sarikaya [205] refers to digital assistants as a “metalayer of intelligence” between the user and various services and actions. Voice-controlled digital assistants are being used to perform an increasing range of tasks, including web searching, diary management, sending emails, and posting to social media. Their functionalities are being expanded from private to business use², as well as to personal banking.³ With the advent of assistants such as Amazon’s Alexa that can be used to control smart home devices, control of systems via a speech interface has extended beyond purely virtual environments to include cyber-physical systems. Pogue [188] describes voice control as a “breakthrough in convenience” for the Internet of Things. Kar and Haldane [110] describe the role of digital assistants as addressing the higher levels of the DIKW (Data-Information-Knowledge-Wisdom) pyramid as applied to the Internet of Things, defining the ‘Wisdom’ level of the pyramid in this context as ‘actionable intelligence’ embodied in such assistants.

The first voice-controlled digital assistant to be released commercially was Apple’s Siri in 2011. Siri was based on an earlier system named CALO that had been developed with US defence funding. Siri was followed by the release of Amazon’s Alexa in 2014, Microsoft’s Cortana in 2015, and most recently in 2016 by Google Assistant.⁴

Speech interfaces have particular potential benefits for elderly or physically disabled users (see Christensen et al. [39]). They may eventually also be used in time-sensitive and

¹A recent UK government survey, for example, reported that eight percent of adults in the UK now own a smart speaker, see <https://gds.blog.gov.uk/2018/08/23/hey-gov-uk-what-are-you-doing-about-voice/>

²See Inc., 30th November 2017, “Why Amazon’s Alexa may soon become your new colleague”, <https://www.inc.com/emily-canal/amazon-alexa-for-business.html>

³See for example The Asian Banker, 7th June 2018, <http://www.theasianbanker.com/updates-and-articles/a-voice-led-future-banking-with-alexa-and-similar-services-become-mainstream>

⁴Financial Times, 21st September 2016, “Google uses Assistant to square up to Siri in AI arms race”, <https://www.ft.com/content/f9423056-7efe-11e6-8e50-8ec15fb462f4>

even life-critical contexts, such as in medical or military applications.⁵ Another application that has been proposed for speech interfaces is the use of voice to control laboratory instruments (see Austerjost et al. [14]). There is some speculation that communication with computers via natural spoken language represents the next major development in computing technology.⁶

Notwithstanding its potential benefits, the introduction of a speech interface represents a potential expansion of a system’s attack surface. Researchers have demonstrated, for example, that it is possible to bypass the lock screen functionality on a Windows desktop using voice commands to Microsoft’s voice assistant Cortana.⁷ The speech interface is inherently difficult to secure, on account of the difficulty of controlling access to a system by sound. A speech interface potentially enables an attacker to gain access to a victim’s system without needing to obtain physical or network access to their device. Whereas physical access to a system may be controlled by physical measures such as locks, and network access can be controlled by technical measures such as firewalls, access to a system by sound is more difficult to control. This is implicitly recognised in a recommendation by Jackson and Orebaugh [103] to move Amazon Echo devices used for smart home control away from doors and windows to prevent them from being controlled by external voices. In addition to being difficult to prevent, attacks on a system via sound may also be difficult to attribute to a source, especially if they are not heard by the victim at the time of their execution. Attacks via the speech interface may be characterised as falling within the broader category of ‘bridgeware’ for gaining access to systems that are not connected to the public internet (see Guri and Elovici [77]).

The vulnerability of voice-controlled systems to attacks via the speech interface is compounded by the fact that the speech recognition and natural language understanding technologies incorporated in voice-controlled digital assistants are designed to respond to speech input in as flexible and natural as way as possible, so as to ensure ease of interaction with users. This design principle is at odds with the general cyber security principle of distrusting user input.⁸ Thus the human-like digital personas intended to give users a sense of familiarity and control in interactions with their systems may in reality be exposing users to additional risks. In the case of home assistive voice assistants for use by the elderly and physically disabled, the victims of attacks via the speech interface may be amongst society’s most vulnerable groups.

Internet security company AVG pointed out in 2014 the danger of the speech interface being exploited as a new attack surface, demonstrating how smart TVs and voice assistants might respond to synthesised speech commands crafted by an attacker as well as to their

⁵See Franzese and Coyne [62], Nathan et al. [173], Zinchenko et al. [258] and The Guardian, 12th September 2017, “British navy warships ‘to use Siri’ as technology transforms warfare”, <https://www.theguardian.com/uk-news/2017/sep/12/british-navy-warships-to-use-voice-controlled-system-like-siri>

⁶See for example James Vlahos, in his book *Talk to Me: Amazon, Google, Apple and the Race for Voice-Controlled AI*, Random House Books, 2019.

⁷See BlackHat USA 2018 Briefings, “Open Sesame: Picking Locks with Cortana”, <https://www.blackhat.com/us-18/briefings/schedule/index.html>

⁸See for example in ENISA Info notes published 1st June 2016, “The Dangers of Trusting User Input”, <https://www.enisa.europa.eu/publications/info-notes/the-dangers-of-trusting-user-input>

users' voices.⁹ The reality of this possibility was first illustrated by a TV advertisement that contained spoken commands for activation of Google Home on listeners' phones for product promotion purposes. The advert was criticised as a potential violation of computer misuse legislation in gaining unauthorised access to listeners' systems.¹⁰ In another instance, it was shown to be possible to open a house door from the outside by shouting a command to digital assistant Siri (as discussed by Hoy [95]). Whilst there has been a considerable amount of debate on the threat to privacy posed by 'listening' devices, as highlighted in media reports on a request for speech data from Amazon's Alexa as a 'witness' in a murder inquiry¹¹ and discussed in research such as that of Hui and Leong [98], security concerns associated with human-computer interaction by speech have yet to be comprehensively addressed. The measures put in place to ensure privacy of speech data, such as the use of encryption to prevent 'sniffing' of speech data in transit between a local device and a cloud provider's server¹², do not solve the security problems associated with local access to a device by speech. This thesis makes a contribution to filling this gap.

The aim of this thesis is to investigate security issues that are specifically associated with human-computer interaction by speech. The research questions that the thesis seeks to answer are whether or not the speech interface presents new vulnerabilities not found in other types of interfaces, and, if so, what these vulnerabilities are, and what defences are required to prevent their exploitation by malicious actors. In answering these research questions, the thesis substantiates three key arguments. The first argument is that the speech interface represents a new type of attack surface, the vulnerabilities of which have not been comprehensively investigated in prior work. The second argument is that malicious input to a voice-controlled system is not limited to the space of speech sounds that are intelligible to humans as input to that system, but that malicious input may also come from the space of sounds to which humans allocate a different meaning to the system, or no meaning at all. The third argument is that current defence mechanisms are not sufficient to defend against all possible types of attacks via the speech interface, and that effective protection will require the development of new defence mechanisms. These key arguments are substantiated in the three main components of the thesis, representing its three main contributions, as follows:

⁹Forbes, 29th September 2014, "Voice Hackers Will Soon Be Talking Their Way Into Your Technology", <https://www.forbes.com/sites/jasperhamill/2014/09/29/voice-hackers-will-soon-be-talking-their-way-into-your-technology/>

¹⁰See Sophos Naked Security blog, 18th April 2017, "Burger King triggers Google Home devices with TV ad", <https://nakedsecurity.sophos.com/2017/04/18/burger-king-triggers-ok-google-devices-with-tv-ad/> The blog quotes privacy activist Lauren Weinstein: "... the federal CFAA (Computer Fraud and Abuse Act) broadly prohibits anyone from accessing a computer without authorization. There's no doubt that Google Home and its associated Google-based systems are computers, and I know that I didn't give Burger King permission to access and use my Google Home or my associated Google account. Nor did millions of other users. And it's obvious that Google didn't give that permission either. "

¹¹See Wired, 28th February 2017, "A Murder Case Tests Alexa's Devotion to Your Privacy", <https://www.wired.com/2017/02/murder-case-tests-alexa-s-devotion-privacy/>

¹²See for example <https://www.iot-tests.org/2017/02/testing-amazon-echo-dot-alexa-app/>

- A novel high-level taxonomy of the types of attacks that might be executed via a speech interface. In this taxonomy, attacks via the speech interface are categorised in terms of human perception, rather than in terms of attack methods or technical vulnerabilities, albeit that the taxonomy is aligned to technical vulnerabilities in the current generation of voice-controlled systems. Whilst the categories of attack based on human perception identified with this approach can be expected to remain stable over time, their alignment to specific technical vulnerabilities might be expected to shift as the state-of-the-art in voice control advances. Thus the taxonomy has the potential to encompass both attacks that are currently possible and new types of attacks exploiting different vulnerabilities that may become possible in the future.
- Experimental work demonstrating two new types of attack via the speech interface that are envisaged in the high-level taxonomy but that have not been demonstrated in prior work or seen ‘in the wild’. The novel attacks both rely on mismatches between machine and human perceptions of input to a voice-controlled system, and develop systematic methodologies for exploiting such mismatches. Specifically, the first attack exploits unintended functionality in the speech recognition functionality of voice-controlled systems by hiding voice commands in nonsensical word strings, whereas the second attack exploits unintended functionality in the natural language understanding functionality of such systems by hiding voice commands in apparently unrelated utterances. Methodologies for executing these attacks are first developed in pilot experiments and then applied to larger scale main experiments. The experiments test both machine and human responses to validate the existence of mismatches between machine and human perceptions of input that enable such attacks to be executed covertly.
- Proposals for the development of new defence mechanisms to counter attacks via the speech interface that would enable voice-controlled systems to detect potential attacks as part of their dialogue management functionality, and to issue verbal security alerts to users via their speech synthesis functionality. These proposals are grounded in a new attack and defence modelling approach for conceptualising the security of the speech interface that represents a novel application of the Observe-Orient-Decide-Act (OODA) loop model to the context of voice-controlled systems. This modelling approach facilitates an assessment of the effectiveness of currently available defence mechanisms in countering various types of attacks via the speech interface, and enables the identification of potential new points of defence. Following conclusions from the assessment of current defence mechanisms, the proposals for new defence mechanisms focus on attacks that exploit discrepancies between machine and human capabilities in speech recognition and natural language understanding.

The aim of the thesis to identify security issues that are specific to human-computer interaction by speech implies a focus on attacks that gain unauthorised access to a system via the speech interface itself. It is acknowledged that there are many possibilities for attacking a voice-controlled system other than via the speech interface. In a security analysis of Amazon’s Echo, Haack et al. [80] identify three means of attack on such systems. In addition

to sound-based attacks via the speech interface, the paper identifies network attacks (eg. sniffing of speech data in transmission between an individual user’s device and a provider’s servers) and API-based attacks (which might involve hacking a voice-controlled assistant’s API eg. to change the default wake word). Cho et al. [37] also include spoofing and tampering of packets transmitted between a voice assistant client device and a voice assistant cloud server in an analysis of potential attacks on a voice-controlled system. Kennedy et al. [112] have shown that it is possible to sniff speech data in transmission even if the data is encrypted, using a technique called voice command fingerprinting to detect traffic patterns that match common voice commands. Security researchers have further demonstrated that it is possible to gain a root shell on an Amazon Echo device via physical tampering¹³. It has also been shown to be possible to develop malicious third-party voice applications with misleading functionalities that can be uploaded to cloud-based voice-controlled systems and used to compromise user privacy, extract personal information, or distribute disinformation (Zhang et al. [255], Mitev et al. [164]).¹⁴ Notwithstanding the significance of these attacks, the compromise of voice assistants by such non sound-based methods is outside the scope of this thesis. Correspondingly, the scope of the thesis is also limited to defence mechanisms that are specific to preventing unauthorised access to a system via the speech interface, excluding defences that may also be applied to prevent access to a system via other modalities. Thus the thesis does not consider authentication measures such as passwords or PINs that may be employed to prevent unauthorised access via non-speech interfaces as well as via speech interfaces. Such measures are subject to the same strengths and weaknesses in relation to the speech interface as they are in other contexts. Traditional authentication measures such as passwords and PINs are in any case difficult to adapt to the speech interface.

The research presented in this thesis was conducted in a monolingual English context, whilst recognising an imperative to conduct similar work in other language contexts in future. Conversational interfaces in fact have particular advantages in language contexts with a writing system that is not very suited to text-based interfaces, such as Chinese.¹⁵

The remainder of the thesis is structured as follows. Chapter 2 contains an overview of the technologies behind voice-controlled systems and their implementation in cloud-based voice assistants. Chapter 3 presents a critical review of prior work relevant to attacks via the speech interface. Chapter 4 presents a novel high-level taxonomy of attacks via the speech interface, based on the review of prior work in Chapter 3. Chapter 5 presents the results of experimental work demonstrating two new types of attacks via the speech interface. Chapter 6 develops an approach to attack and defence modelling in the context of the speech interface based on OODA loop model. This includes a critical review of current defence mechanisms. Chapter 7 presents proposals for the development of new defence mechanisms that are grounded in the attack and defence modelling approach presented in Chapter 6. The proposals are supported by the results of feasibility tests on two potential defences

¹³See <https://labs.mwrinfosecurity.com/blog/alexa-are-you-listening>

¹⁴See also <https://www.wired.com/story/amazon-echo-alexa-skill-spying/>

¹⁵See MIT Technology Review, “Conversational Interfaces: Powerful speech technology from China’s leading Internet company makes it much easier to use a smartphone”, <https://www.technologyreview.com/s/600766/10-breakthrough-technologies-2016-conversational-interfaces/>

against attacks targeting vulnerabilities in speech recognition and natural language understanding, respectively. Chapter 8 concludes the thesis and describes the outlook for future work in this area.

Chapter 2

Background on Human-Computer Interaction by Speech

2.1 Overview of Speech Dialogue Systems

This section reviews the components of generic speech dialogue systems for execution of particular actions by voice command. Speech interfaces that facilitate the execution of particular actions in response to voice commands are referred to as ‘task-based’ or ‘goal-driven’ speech dialogue systems, as distinct from ‘chatbots’ or ‘open’ speech dialogue systems, whose purpose is simply to hold a conversation with the user without executing any actions (see for example Celikyilmaz et al. [34]). The current generation of voice-controlled digital assistants has some similarity with chatbots in that they are often anthropomorphised, with systems being given the persona of a friendly digital assistant in order to create a sense of communication with a human-like conversation partner that is capable of engaging in some casual chat as well as executing actions. However, this review focusses on the technologies implemented in such systems only in as far as they are relevant to the systems’ task-based functionalities.

Input to a speech dialogue system is provided by a microphone that captures speech sounds and converts these from analog to digital form. Following conversion to digital form, the system transcribes from the speech signal a sequence of written words in a process known as speech recognition, and then extracts from this sequence of words a representation of their meaning in a process known as natural language understanding. Bellegarda and Monz [16] describe the task of the speech recognition component as the task of extracting from a set of acoustic features the words that generated them, and the task of the natural language understanding component as the task of extracting from a string of words a computational representation of the user intent behind them. The paper by Bellegarda and Monz conceptualises the process of a user’s communication of intent to a speech dialogue system as information transmission across a noisy channel, whereby the user first formulates their intent in words and then vocalises these words as speech, and the dialogue system subsequently performs the same process in reverse. This process is illustrated in the diagram in Figure 2.1, as per Bellegarda and Monz’s paper.

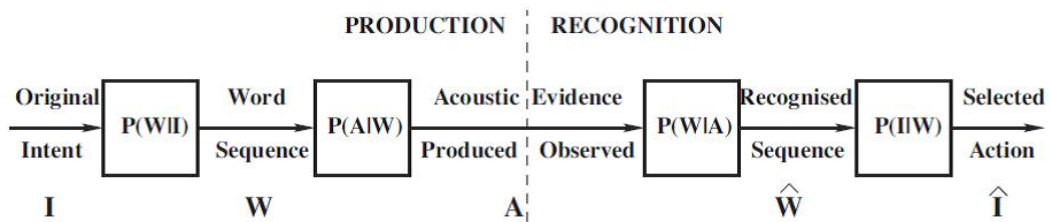


Figure 2.1: (from Bellegarda and Monz [16]) An example of integrated speech and language processing: personal assistance seen as information transmission across a noisy channel

The typical architecture of a generic speech dialogue system consists of components for speech recognition, natural language understanding, dialogue management, response generation and speech synthesis, as detailed by Lison and Meena [147]. The architecture of speech dialogue systems as presented by Lison and Meena is shown in Figure 2.2. An overview of the history and state-of-the-art of the technologies behind each of these components is provided in the sections below. As an attack on a speech dialogue system can be expected to target the input processing components of the system, i.e. the speech recognition and natural language understanding components, this overview focusses primarily on these components.

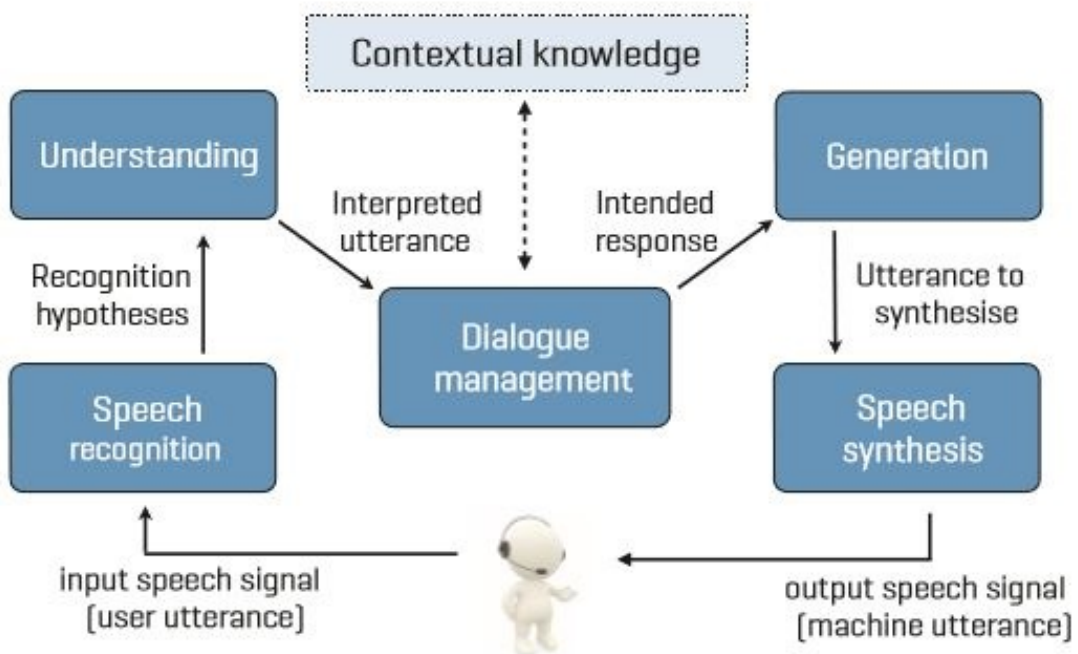


Figure 2.2: The Architecture of Speech Dialogue Systems (from Lison and Meena [147])

Speech Recognition Speech dialogue systems require in the first instance a capability for converting a series of speech sounds to a textual character string. This process, known as speech recognition, generally consists of combining an acoustic model for mapping speech sounds to word units with a language model for determining valid sequences of words. Juang and Rabiner [107] and Pieraccini and Rabiner [186] summarise the development of speech recognition technology. Speech recognisers for small sets of individual words, such as spoken digits, were first developed in the 1950s, and efforts to develop recognisers capable of transcribing continuous speech began in the 1960s, using dynamic programming techniques to match utterances of varying lengths to speech sound templates. The simplistic nature of such early systems led Pierce [187] to question the value of pursuing continuous speech recognition: “There are strong reasons for believing that spoken English is, in general, simply not recognizable phoneme by phoneme or word by word, and that people recognise utterances, not because they hear the phonetic features or the words distinctly, but because they have a general sense of what a conversation is about and are able to guess what has been said. – These considerations lead us to believe that a general phonetic typewriter is simply impossible unless the typewriter has an intelligence and a knowledge of language comparable to those of a native speaker of English.” However, as described by Juang and Rabiner, speech recognition continued to be pursued in the 1970s with the support of US government funding, as seen for example in the development of the Harpy system at Carnegie Mellon University, in which template matching was combined with grammatical rules that constrained the search space needing to be covered by the dynamic programming algorithm. In parallel to academic efforts, the commercial goals of producing an automated typewriter and automated call-centres were pursued by IBM and AT&T Bell Laboratories respectively. At IBM, the main objective was to increase the size of the vocabulary that a system was capable of transcribing, with a consequent focus on the language model incorporated in the recogniser. The most important objective for automated call-centres, on the other hand, was to understand variable pronunciation of the same words by large numbers of different speakers, leading to a focus on acoustic models.

For language modelling, computational representations of natural language structures were originally attempted using the formal grammars that had been developed by Noam Chomsky based on theories of innate language learning processes in humans (Chomsky [38]).¹ However, contrary to early assumptions on a need for automatic speech recogni-

¹Chomsky’s work was based on theories of innate language learning processes in humans, with formal grammars representing attempts to formalise the underlying ‘rules’ according to which natural language sentences were generated. The type of grammars most commonly used in language processing were context-free grammars. A context-free grammar generates sentences according a set of production rules whereby higher level language constituents (such as noun phrases and verb phrases) can be replaced by one or more other language constituents (for example a verb phrase may be replaced by a verb and a noun-phrase). Constituents that cannot be replaced by other constituents, i.e. individual vocabulary or alphabet items such as words or morphemes, are known as terminals. Sentences generated with a context-free grammar are recognizable by a pushdown automaton. Given the assumption that a sentence has been generated by a particular context-free grammar, it can be ‘reverse-engineered’ into constituent parts in a process known as parsing, using various algorithms such as the Cocke–Younger–Kasami (CYK) algorithm. Chomsky acknowledges that there are many instances where the same sentence might have been produced by two or more different applications of the production rules, giving as an example the sentence ‘they are flying planes’ in which ‘flying planes’ may be either subject or object. Context-free grammars are not capable of definitively resolving this kind of ambi-

tion systems to incorporate complex knowledge on the processing of speech by humans, it was found that the performance of automated systems actually improved when attempts to model the true complexity of human language were abandoned, and the speech recognition process was instead handled using the relatively simple n-gram model. This new approach to language modelling was first pursued by Fred Jelinek's group at IBM. N-grams could be used to calculate the probability of occurrence of words given the n-1 words preceding it in a sentence. The use of n-grams in language modelling had its roots in Claude Shannon's work on information theory (see Shannon [219]). N-grams enabled a speech recogniser to resolve ambiguity arising from the presence of homophones in continuous speech, by identifying the homophone that was most likely to be the correct word in the sentence context.²

In acoustic modelling, a key development was seen in the 1980s with the application of Hidden Markov Models (HMMs) to the speech recognition task (see Juang and Rabiner [107]). HMMs removed dependence on speech sound templates and facilitated probabilistic representation of variability in word pronunciations. HMM-based speech recognition systems typically use Gaussian Mixture Models (GMMs) to model the likelihood of particular acoustic features having been generated by particular phonemes. Phonemes are the individual units of sound used in a language to make up its spoken words. Speech recognition with HMMs involves firstly extracting from the speech signal a set of acoustic features, the most commonly used features being Mel-frequency cepstral coefficients (MFCCs). The phonemes identified as the most likely to have generated the set of acoustic features are then matched to the words associated with them, using a pronunciation dictionary, to create a word-level HMM. The probabilities of transition between words are then calculated using an n-gram language model. Finally, the most likely spoken word sequence is calculated according to Bayes' rule as the product of the likelihood of the acoustic features and the probability of the occurrence of particular words in the sentence context as calculated by

guity, although it is possible to add probabilities to transition rules that allow alternative parsing structures to be identified as more or less likely. Also, despite their relative complexity in comparison to n-grams, context-free grammars are known not to be capable of capturing comprehensively the full set of possible sentences in natural languages. It has been shown that certain languages, such as Swiss-German, have properties such as so-called cross-serial dependencies that cannot be expressed in a context-free grammar (see Dan Jurafsky and James H. Martin, *Speech and Language Processing*, 2nd ed.).

²Unlike approaches based on Chomsky's grammars, information-theoretic approaches to language modelling are not concerned with linguistic structures, but instead treat the problem of identifying the most likely word sequence associated with a given set of acoustic features as a purely mathematical problem. Shannon was unconcerned with the meaning of communications (he states that "semantic aspects of the message are irrelevant to the engineering problem"); instead he is concerned with measuring mathematically the amount of information present in a message. The unit of measurement is the base 2 logarithm of the potential number of messages, termed a 'binary digit' or 'bit'. Shannon defines information as contained in a message as the amount of uncertainty that the message resolves, which is related to the probability of receiving this message as opposed to the probabilities of other possible messages in a given context. As part of the broader development of his theory, Shannon considers the probabilities of the occurrence of individual letters and words in natural English, showing that when text is generated in accordance with such probabilities, the generated text can be made to approximate natural English more closely by making the probability of occurrence of a given letter or word in a sequence being generated dependent on the probability of the letter or word given the letter(s) or word(s) preceding it in the sequence, rather than just on the probability of occurrence in the language as a whole. This was the origin of the statistical approach to natural language processing using n-grams.

an n-gram language model. HMMs use the Viterbi algorithm to determine the most likely word sequence for a given segment of speech. They are trained using the forward-backward algorithm.³ A significant early HMM-based system was the SPHINX speech recognition system developed at Carnegie Mellon University by Kai-fu Lee (see Lee et al. [137]).

HMM-based systems for speech recognition continued for a long time predominantly to use GMMs for acoustic modelling and n-grams for language modelling. However, in recent years a shift in modelling methods has been seen with the advent of deep learning. Huang et al. [97] describe recent developments in which Deep Neural Networks (DNNs) have replaced GMMs to extract acoustic model probabilities, and Recurrent Neural Networks (RNNs), a particular type of DNN, have replaced n-grams to extract language model probabilities. A DNN is a neural network containing one node layer for inputs, one node layer for outputs (i.e. the labels that the network applies to input), and one or more of hidden layers between the input and output layers.⁴ Bellegarda and Monz [16] cover the application of neural networks in both acoustic and language modelling. The use of DNNs for acoustic modelling is also described in detail by Hinton et al. [92]. A key characteristic of RNNs as used for language modelling is that they contain recurrent loops that enable the output of a network node at one point in time to be fed back as input at the next point. This allows RNNs to take account of history when handling sequential input such as a sequence of words; hence the applicability of RNNs in language modelling.⁵ Rather than taking individual words as input, RNNs for language modelling commonly take as input so-called word embeddings, which are vector representations of words based on the lexical contexts in which they occur. One common type of word embedding is the Word2Vec skip-gram model, which predicts the most likely neighbouring words of a given word within a certain window (see further details in Mikolov et al. [161]). The use of word embeddings rather than individual words enables an RNN to handle language input more flexibly, amplifying its capacity to deal with language input that was not seen during training of the RNN.

The state-of-the-art in speech recognition has now reached a point where it is on a par with human abilities in terms of performance on standard speech recognition tests. In

³See Dan Jurafsky and James H. Martin, *Speech and Language Processing*, 2nd ed.

⁴The nodes in one layer of the network are connected to the nodes in the next layer by weighted links, whereby the weightings are initially set at random. Node values at one layer of the network are transformed using a non-linear function such as the sigmoid (logistic) function, and then fed forward as input to nodes in the next layer. The input received by each node in the next layer of the network will consist of the transformed sum of the values of all nodes in the previous layer that are connected to it, with each of these values having been weighted by the weighting applied to the link between the relevant nodes. These calculations can be done for networks with large numbers of layers and nodes using matrix multiplication. To train a network, the error in the outputs of the network is measured in relation to a set of labelled training data using an objective (cost) function. The links weights in the network are then optimised so as to minimise the error rate by using stochastic gradient descent in combination with the so-called backpropagation algorithm that feeds the error measurements back into the network. For DNNs that map input to more than two labels, the final layer of the network will typically be transformed using the ‘softmax’ function, which is a multi-class extension of the logistic function.

⁵In theory, RNNs can handle sequential inputs of infinite length, although in practice this becomes mathematically infeasible due to the so-called ‘vanishing gradient’ problem where the signal from previous inputs is diluted over time. Developers have sought to address this problem with the incorporation of so-called Long Short-Term Memory Cells, which enable nodes in an RNN to ‘decide’ whether or not to take account of input they receive depending on the significance of the input.

2016, Microsoft Research reported that its automatic speech recognition capability had for the first time matched the performance of professional human transcriptionists, achieving a word error rate of 5.9 percent on the Switchboard dataset of conversational speech produced by NIST in the US (see Xiong et al. [244]). Huang et al. [97] have predicted that automated speech recognition will eventually actually exceed human capabilities. However, variability in individual pronunciation of words, as investigated for example by Benzeghiba et al. [17], continues to pose challenges to speech recognition for general usage. Speech recognition systems also struggle to distinguish between speech that is directed at them and background speech (see Shriberg et al. [222]). A further goal in current speech recognition research is detection of speakers' emotions in the speech signal (see Schuller [213]).

Speech recognition systems are on the whole limited to recognition of the words available in a pronunciation dictionary incorporated in the recogniser, and therefore need to be able to handle out-of-vocabulary words that may be provided to them as input. For early speech recognition systems with limited vocabularies, various methods were developed to enable speech recognition systems to detect out-of-vocabulary input and avoid confusion with in-vocabulary words (see for example Hazen and Bazzi [86]). The current generation of voice-controlled assistants are able to recognise the vast majority of words in use in the relevant language, as evidenced by their ability to perform unrestricted actions such as web searches. The significance of out-of-vocabulary word detection for current systems is therefore largely limited to incorporating previously unseen named entity words, such as names or locations, as investigated for example by Ilina and Fohr [100]. Some efforts have been made towards developing speech recognition systems capable of learning new words in ongoing interactions (see for example Qin [192]).

Natural Language Understanding Enabling computer systems to 'understand' natural language in the general sense remains a far-off and perhaps inherently unachievable goal. In the context of a voice-controlled system such as a voice-controlled digital assistant, natural language understanding is restricted to the task of extracting from a user's request a computational representation of its meaning that can be used by the system to trigger an action. The task of mapping a string of words to a representation of their meaning is known as semantic parsing. The process of semantic parsing may include syntactical parsing as an intermediate step. The natural language understanding process may also be supported by lexical resources, such as the WordNet semantic network, which provides the system with information on synonymous relationships between words (see for example Landhäußer et al. [134]).

A method of syntactic analysis frequently used in voice-controlled systems is in-out-begin (IOB) tagging, which is a shallow parsing method that involves identifying key constituents of a sentence (such as a noun phrase), and then labelling each word as being inside, outside, or at the beginning of each constituent part.⁶ Another method of syntactic analysis sometimes used in voice control is dependency parsing, which involves determining syntactic relationships within a sentence such as verb-object connections (see for example McTear [156]). As spoken language is frequently less grammatical and structured than

⁶See Dan Jurafsky and James H. Martin, *Speech and Language Processing*, 3rd ed. draft <https://web.stanford.edu/~jurafsky/slp3/>

written language, syntactic parsing using formal grammars, such as syntactic parsing with context-free grammars, is not appropriate for spoken language, as pointed out for example by Shen [221].

As regards computational representations of meaning for natural language understanding, outside the field of voice control these have often been based on formal logical structures, following work in the 1970s by Richard Montague on logical semantics [165], and by Terry Winograd on an early natural language understanding system named SHRDLU [241]. In recent times, Berant et al. [18] and Liang [144] present work on using logical forms for natural language understanding tasks such as question answering. However, formal semantic representation is difficult to implement in a spoken language context, as stated by Einolghozati et al. [55]. Natural language understanding in voice-controlled systems tends to be based on less formal representations.

An early method of natural language understanding in speech dialogue systems, known as keyword-spotting, is described by Juang and Rabiner [107]. This approach was developed in the 1970s for use in automated call-centres to enable some flexibility in users' expressions of their requests. Rather than requiring exact use of certain sentences by users, keyword-spotting was intended to ensure that different expressions of a user's request would be understood by the system as long as it contained certain words. In keyword-spotting systems, the computational representation of meaning consists simply of a set of keywords, and semantic parsing consists simply of matching these keywords to words in a transcribed utterance.

Current speech dialogue systems use more complex semantic representations. One widely used semantic representation is semantic frames (see Sarikaya et al. [206], Canonico and De Russis [29]). Semantic frames provide a structure for representing the meaning of utterances that requires, firstly, identification of the general domain or concept that a user request relates to (such as travel); secondly, determination of the user intent (such as to book a flight); and thirdly, slot-filling, which involves identifying specific information relevant to the particular request (such as destination city). Slot-filling may be performed over more than one utterance to take into account contextual information from several 'turns' within the same dialogue interaction (this process is known as slot carryover). Einolghozati et al. [55] discuss an alternative to the standard semantic frame structure in which intents can be nested within slots as well as vice-versa, so as to provide increased flexibility in parsing of utterances.

Regarding the mapping of spoken utterances to semantic representations, whilst keyword-spotting is clearly a deterministic process, semantic parsing to more sophisticated representations of meaning, such as semantic frames, is generally performed using probabilistic methods. It has been shown that semantic parsing to such representations using deterministic methods is not practical on account of the complexities and ambiguities involved in the task (see for example Gal [65]). Sarikaya [205] state that the tasks of domain identification and intent determination in semantic parsing to frames are often performed using support vector machines, whereas slot-fitting is commonly performed using conditional random fields (CRFs). CRFs is a discriminative machine learning technique for sequence classification based on maximum likelihood. By contrast, HMMs, which are also used for sequence classification, are a generative technique based on Bayesian inference (see Nadkarni et al. [170]). CRFs calculate the probability of a state taking into account information

from the observed sequence as a whole, in the form of weighted features extracted from the observed sequence, whereby in parsing for natural language understanding these features will consist of manually selected linguistic features. Jurafsky and Martin state that domain and intent classification in voice assistants might be performed with a ‘1-of-N’ classifier that uses n-grams, whereas slot-filling might be performed with a CRF that uses n-grams, named entities, and slot-transition sequences as features.⁷ Machine learning classifiers for semantic parsing are generally trained with large datasets of manually labelled examples, although there has been some work towards partially unsupervised machine learning for natural language understanding using the so-called active learning approach (see Yang et al. [246]).

Some recent research has indicated that traditional machine learning methods may now be being out-performed in the semantic parsing task by neural networks, similar to the replacement of n-gram-based systems for language modelling in speech recognition by RNNs. LeCun et al. [136] note that the deep learning approach to pattern-recognition in high-dimensional data that has achieved state-of-the-art performance in image classification and speech recognition is also proving effective in natural language understanding tasks. Jurafsky and Martin state that RNNs can perform all three natural language understanding tasks, i.e. domain identification, intent determination and slot-filling, simultaneously, by adding domain concatenated with intent as the final step of sequence labelling in combination with slot-filling.⁸ Bellegarda and Monz [16] note that the Word2Vec word embeddings used as inputs to RNNs for language modelling for speech recognition are also capable of capturing semantic relations between words, giving a well-known example that subtracting the word vector for ‘man’ from the vector for ‘king’ and then adding the vector for ‘woman’ will equal the vector for ‘queen’. Word2Vec embeddings are also known similarly to capture syntactical relationships such as verb tense.⁹

Mesnil et al. [158] present results showing superior performance by RNNs on the slot-filling task for the Air Travel Information System (ATIS) dataset in comparison to the performance of CRFs on the same task. The ATIS task involves extracting information such as departure date and destination from user requests in order to provide relevant travel information. Mesnil et al. use RNNs to perform IOB chunking in the ATIS task (the IOB format being used to label words as being inside, outside, or at the beginning of a slot, rather than being used to represent grammatical components of the sentence as such). The paper gives the example of the word ‘Boston’ in the sentence “show flights from Boston to New York today” being labelled ‘B-dept’, where the ‘B’ indicates the word’s position at the beginning of a slot, and the ‘dept’ indicates the word’s concept label as city of departure. Yao et al. [247] also present results showing superior performance of RNNs on slot labelling in the ATIS dataset in comparison to CRFs. Ravuri and Stolcke [195, 196] present results showing superior performances of RNNs on the domain and intent classification task in comparison to other methods. Vukotic et al. [237], however, found that whilst RNNs have shown superior performance to other machine learning methods in simple spoken lan-

⁷Dan Jurafsky and James H. Martin, *Speech and Language Processing*, 3rd ed. draft <https://web.stanford.edu/~jurafsky/slp3/>

⁸Dan Jurafsky and James H. Martin, *Speech and Language Processing*, 3rd ed. draft <https://web.stanford.edu/~jurafsky/slp3/>

⁹See <https://www.tensorflow.org/tutorials/word2vec>

guage understanding tasks, such as the frame-based slot-filling task for ATIS, they remain inferior to CRFs in more complex tasks such as determining concept labels in the more challenging MEDIA dialogue corpus. Andor et al. [10] propose a system called SyntaxNet that combines a CRF objective function with a standard, i.e. non-recurrent, feed-forward neural network to perform syntactic parsing to a dependency parse tree as part of Google’s natural language understanding capability.¹⁰ Perera et al. [184] describe semantic parsing with RNNs to the semantic representation used by Amazon Alexa, the so-called Alexa Meaning Representation Language, which is a graph-based alternative to semantic frames [120]. Agarwal et al. [4] propose augmenting RNN-based methods for slot-filling with an additional DNN-based parser to handle coordinated structures. Some work has further suggested the joint training of the speech recognition and natural language understanding components of spoken language understanding using RNNs (see Haghani et al. [81]).

It is clear that, unlike in the case of speech recognition, the state-of-the-art in natural language understanding remains far from parity with human capabilities (see for example Liang [145] and Cambria and White [28]). Despite their increasingly human-like interactions, voice interfaces are not yet capable of fully meeting Turing’s ‘imitation’ challenge in terms of their ability to simulate human understanding consistently in natural language interactions [233]. Hirschberg and Manning [93] discuss the challenges posed to natural language processing by pervasive ambiguities in natural language. This is evident in the occasional failure of voice assistants to interpret the meaning of a word correctly in context, despite the correct word or meaning being obvious to any human listener. Stolk et al. [224] give the examples of Apple’s assistant Siri mistaking the word ‘bank’ in the sense of ‘river bank’ for a financial institution, and of Siri giving directions to a casino when asked about a gambling problem. The first of these examples represents a failure on the part of the assistant to handle ambiguity of word meaning, whereas the second represents a failure to understand the broader intent of the utterance. Other instances of mismatch in natural language understanding between digital assistants and their human users are more difficult to explain, one example being incidents in which user requests to Siri to charge a phone prompted the assistant to dial emergency services.¹¹

Errors in natural language understanding represent a significant flaw in the current generation of voice-controlled digital assistants. Sarikaya [205] states that the largest percentage of errors in voice-controlled digital assistants is due to errors in natural language understanding. Sano [204] state that whereas in the case of a speech recognition error, users tend to repeat a voice command, in the case of a natural language understanding error, users are likely to abandon their attempt to communicate with a device. In a discussion of natural language interfaces for databases, Li and Jagadish [141] make the point that whereas non-exact responses are acceptable in the case of a web search query, in the case of a database

¹⁰See Google Research Blog, 12th May 2016, “Announcing SyntaxNet: The World’s Most Accurate Parser Goes Open Source”, <https://research.googleblog.com/2016/05/announcing-syntaxnet-worlds-most-accurate-parser-goes-open-source/> and Wired, 5th Dec 2016, “Google Has Open Sourced SyntaxNet, Its AI for Understanding Language”, <https://www.wired.com/2016/05/google-open-sourced-syntaxnet-ai-natural-language/>

¹¹See The Telegraph, 16th July 2015, “Asking Siri to charge your phone dials the police” - <http://www.telegraph.co.uk/technology/apple/11743790/Asking-Siri-to-charge-your-phone-dials-the-police.html>

query users expect an exact response. A command to a speech interface is comparable to a database query rather than to a web search query in this respect; a voice-controlled system that cannot consistently provide exact responses to commands is unlikely to be acceptable to users. Bickmore et al. [21] recommend that, given the current state-of-the-art in natural language understanding, conversational assistants should not be deployed to provide medical information, because of the potential health consequences for users resulting from errors made by the system.

It has been argued that the gap between human and machine capabilities in natural language understanding will never be bridged, on account of the inherent limitations of artificial systems in mimicking human speech and language processing. Dale [45] and Moore [166] argue that artificial systems such as voice-controlled digital assistants will never achieve parity with human natural language understanding. Khashabi et al. [116] attempt a formal analysis of the limitations of reasoning for natural language understanding using linguistic symbols to represent a hidden space of concepts. The authors state that by contrast to the space of hidden concepts in human thought, the space of linguistic symbols used to express these concepts is noisy, containing many instances of redundancy, ambiguity, incompleteness and inaccuracy that make reasoning with these symbols difficult beyond a limited number of inference steps. Writing in 1965, Dreyfus [52] identified various characteristics of human intelligence that he claimed could never be replicated in machines, two of which - ‘essence/accident discrimination’ (the ability to determine intent) - and ‘ambiguity tolerance’ are particularly relevant to natural language understanding. Armstrong et al. [11] have confirmed the validity of Dreyfus’ predictions in a modern paper.

In the context of modern voice-controlled systems, a failure of ‘essence/accident discrimination’ is evidenced in their lack of ability to always distinguish between utterances that are within their scope in terms of their meaning and out-of-domain utterances. There are ongoing attempts to develop methods of out-of-domain detection, but none of these represent a full solution to the issue as yet. Kim and Kim [117] present a method that is dependent on a corpus of out-of-domain utterances for training the system, which are unlikely to be comprehensively representative of all possible out-of-domain input. Lane et al. [135] suggest an alternative method that removes the need for out-of-domain training data by relying on classification scores across all in-scope domains, and sequentially treating one of the in-scope domains as out-of-scope for the purposes of training a natural language understanding system. This approach is also unlikely to be capable of handling the entire space of out-of-scope input. Tur et al. [232] propose performing out-of-domain detection according to a combination of syntactic and lexical features that are indicative of a command to a voice-controlled system that a user does not realise is out-of-scope, as opposed to other types of unrelated utterances. This method will not be effective in preventing confusion with out-of-domain input that does not contain the features identified by the authors as indicative of a voice command.

A lack of ‘ambiguity tolerance’ as referred to by Dreyfus is evident in modern voice-controlled systems in the failure of such systems to handle ambiguity in word meaning reliably. Word meanings cannot be easily placed on any quantitative spectrum due to the existence of words with multiple meanings. Current methods for natural language understanding thus do not present a full solution to the problem of ambiguous meaning of words. Whilst the distance between word vectors in vector space based on word co-occurrence

such as Word2Vec have been found to have semantic significance, as discussed above, in the case of words with multiple meanings the distance between word vectors will represent a conflation of distances between separate meanings, rather than a real semantic distance (as stated for example by Young et al. [250]). There have been some attempts to develop word embedding methods for representation of different word senses (see for example Trask et al. [228] and Iacobacci et al. [99]), but these remain far from reliable and are heavily dependent on large amounts of human-labelled training data. Another idea is that presented by Navigli and Ponzetto [174], who propose a system of word disambiguation based on the use of a multilingual semantic network. The idea behind the proposal by Navigli and Ponzetto is that a word is unlikely to be associated with the same set of meanings in different languages. Therefore the sense of a word's translation in a given context is likely to represent the correct word sense. The implementation of such a system of word sense disambiguation in voice-controlled systems may become possible in the future with the development of multilingual voice assistants.

Dialogue Management Dialogue management is the task of determining the most appropriate action that should be taken in response to a user's request. Sarikaya [205] states that the dialogue management task in personal digital assistants is far more challenging than in older speech recognition systems. Older speech recognition systems were commonly limited to one general purpose, such as providing travel information. Digital assistants, by contrast, are designed to perform a large number of tasks, including providing information on many different topics, connecting with web applications to fulfil a variety of user requests, and controlling devices in the Internet of Things. Sarikaya describes the structure of a dialogue manager in a digital assistant as consisting of a dialogue state tracker, which updates the 'state' of the dialogue based on the representation of user intent generated by the natural language understanding module, and a dialogue policy that controls the execution of tasks in response to the user request.

Dialogue management modules in current speech dialogue systems are on the whole still rule-based, i.e. they map user intent to dialogue states and dialogue states to actions based on hand-crafted rules, as stated by McTear [156]. Rule-based dialogue management systems have the advantage of limiting the potential for error and unintended functionality in the dialogue management process ([155]). However, such systems are also likely to be lacking in flexibility and limited in scope. There has been some research on the eventual replacement of current rule-based systems by dialogue management systems based on reinforcement learning, which would enable voice assistants to learn directly from their interactions with users. Scheffler et al. [208] present some early work on reinforcement learning in dialogue systems. Young et al. [249] propose ideas for dialogue management based on Partially Observable Markov Decision Processes (POMDPs), which model a dialogue as a Markov process with transition probabilities between states, for which a probability distribution over all possible states is continuously maintained. This approach seeks to represent the uncertainty inherent in the fact that a user's intent is not directly observable, but rather inferred probabilistically from their utterance. Systems based on POMDPs combine Bayesian inference for belief state tracking to determine the most likely interpretation of a user's utterances with reinforcement learning for optimisation of the dialogue

policy, whereby a reward function is used to train the system as to the most appropriate action to take in response to a user utterance based on user feedback (such user repetition of the command, or user response to an explicit request for confirmation of the intended action). The paper by Young et al. explains that whilst POMDP-based approaches to dialogue modelling offer many advantages in terms of flexibility, the implementation of such systems presents many practical problems on account of the large size of the state-action space that makes exact learning intractable. The authors present some possible solutions to the intractability problem by using approximation techniques such as N-best approaches for belief state tracking and Monte Carlo optimisation for policy learning. The paper also mentions the difficulties of training POMDP-based systems with real users and possibilities for user simulation as an alternative.

Some work has proposed combining the dialogue management component in voice-controlled systems with other components as a single coordinated system. Lemon [139] proposes the joint training of dialogue management and natural language generation functionality in a voice-controlled system by reinforcement learning. Paek and Pierraccini [179] discuss the possibility of training speech recognition and natural language understanding functionalities with the same function as the dialogue management system in reinforcement learning-based systems. In the context of open dialogue systems, i.e. chatbots, Serban et al. [214] have proposed modelling the entire sequence from user input to system response as one process using RNNs. They state that such models could be used to generate user simulators for the training of POMDPs for task-based systems.

Response Generation Once an appropriate response has been determined by the dialogue management component, this response is generated by the response generation component. Response generation will involve either or both of generating a verbal response to the user and performing an action requested by the user, as stated by Lison and Meena [147]. This action may be an action on the user's computer, such as sending an email or performing a web search, or an action in a cyber-physical system, such as turning on the lights in a smart home. The task of providing a verbal response to the user via natural language generation is less complex than the converse task of natural language understanding, as natural language generation is not required to handle an infinite set of possible sentences. Sarikaya states that the natural language generation capability in most personal digital assistants is based on simple templates [205].

Speech Synthesis The final action to be performed by a speech dialogue system is the text-to-speech conversion of the natural language response generated by the response generation component. Speech synthesis may involve concatenation of short segments of real speech recordings, or else generation of entirely synthetic speech sounds. Kuusisto [129] envisages the possibility of synthesised voices becoming indistinguishable from the human voice. This prospect appears close to realization with the development of Google DeepMind's WaveNet technology (van den Oord et al. [236]), which is used in Google Assistant.¹² Google has demonstrated its speech synthesis capability Google Duplex in

¹²See <https://deepmind.com/blog/wavenet-launches-google-assistant/>

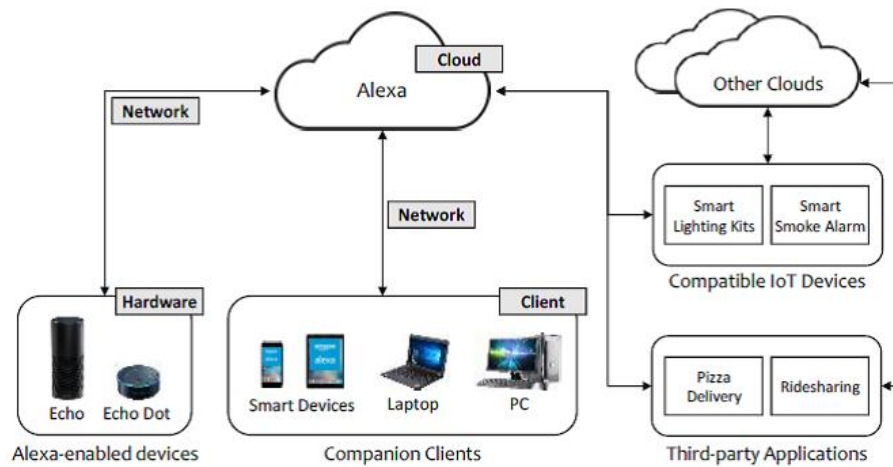


Figure 2.3: (from Chung et al. [40]) Amazon Alexa Ecosystem

Google Assistant interacting with humans to make hairdressing appointments and restaurant bookings. In these demonstrations, the human interaction partners seemed unaware that they were not interacting with another human.¹³ Comparable to speech recognition, a current goal in speech synthesis is the imitation of human emotion in synthesised voice, under the broader area of ‘affective computing’ [156].

2.2 Cloud-Based Voice-Controlled Digital Assistants

Modern voice-controlled digital assistants, such as Google Assistant and Amazon Alexa, implement the generic components of speech dialogue systems in the context of a cloud-based service that enables the speech recognition and natural language understanding functionalities of these systems to be performed in the provider’s cloud. This is due mainly to the amount of computing power required to perform speech and language processing (see for example Strimel et al. [225]). Providers of such cloud-based voice assistants also seek continually to improve the performance of speech recognition in their systems using recordings of users’ interactions with the system.¹⁴ Chung et al. [40] provide an overview of the typical ecosystem of modern voice-controlled digital assistants in the example of Amazon’s Alexa (see Figure 2.3). There has been some research towards performing on-device speech recognition and natural language processing on local devices in future.¹⁵

In order to control the activation of a voice-controlled digital assistant and the streaming of audio data to the cloud, current voice-controlled digital assistants include, in addition to the generic speech dialogue system components, an activation component consisting of a ‘wake word’, which, when spoken by the user, triggers streaming of the subsequent speech

¹³See <https://ai.googleblog.com/2018/05/duplex-ai-system-for-natural-conversation.html>

¹⁴See article by Professor Steve Young in Ingenia magazine March 2013, “Talking to Machines”, <http://www.ingenia.org.uk/Ingenia/Articles/823>

¹⁵See for example Google AI Blog, 12th March 2019, <https://ai.googleblog.com/2019/03/an-all-neural-on-device-speech.html>

audio data to the provider’s cloud for processing. Examples of wake words include ‘Ok Google’ for Google Assistant and ‘Alexa’ for Amazon’s Alexa. Wake word recognition is the only speech processing capability on users’ individual devices, and consists of a short ‘buffer’ of audio data from the device’s environment that is continuously recorded and deleted.¹⁶ Unlike speech recognition in general, which is a multi-class classification problem requiring identification of many different spoken words, wake word recognition is a binary classification problem requiring only distinction between the wake word and all other words and sounds. The words and sounds needing to be rejected by the wake word recognition functionality may be modelled by a noise or garbage model (see Kępuska and Klein [114], Raju et al. [193]).

Wake word activation is nonetheless occasionally prone to being triggered by false positives. Chung et al. [41] for example refer anecdotally to accidental activation of the Alexa assistant by a sentence containing the phrase ‘a Lexus’ (see also Michaely et al. [160]), and Vaidya et al. [235] refer to the misrecognition of the phrase “Cocaine Noodles” as “OK Google”. False positives in wake word recognition may result from misrecognition of a word as the wake word, as in the example given by Chung et al., or else from use of a wake word in the context of speech not intended to activate a voice assistant, for example the use of the word ‘Alexa’ as the name of a person in a conversation. Kępuska and Bohouta [113] discuss the latter problem of distinguishing between an ‘alerting’ and a ‘referential’ context in wake word recognition. It is also possible for voice assistants to be activated by background noise that has frequencies overlapping with those of human speech (see Islam et al. [101]). These instances of unintended functionality in wake word recognition are examples of the continuing inability to perform ‘essence/accident discrimination’ reliably, as referred to in the 1965 paper by Dreyfus [52] discussed above.

The providers behind the current generation of voice-controlled digital assistants have also introduced platforms for the development of third-party voice applications that can be incorporated in the provider’s cloud and made available to users via the assistant’s speech interface. Examples of such third-party applications are Alexa Skills and Google Conversation Actions. Third-party voice applications in systems such as Alexa and Google Assistant can be accessed by users by asking to ‘speak’ to the voice application (as named by the developer).¹⁷ Such applications can be used, for example, to enable users to access information services or to purchase products. One challenge in relation to third-party voice applications is that of making the set of available apps effectively searchable by users wishing to find an application for use for a particular purpose (see for example White [239]). This in fact points to a wider issue of user orientation in the human-computer interaction by speech, in the absence of a graphical user interface or terminal enabling users to see which part of the system they are interacting with.

Cloud-based digital assistants generally have access to a ‘knowledge graph’ or ontology to support their natural language understanding functionality and enable them to answer factual questions (see for example Sarikaya [205]). Originally developed to enhance web searches, such knowledge graphs are now being incorporated in speech-based digital

¹⁶See Wired, 12th May 2016, “Alexa and Google Home Record What You Say, But What Happens To That Data?”, <https://www.wired.com/2016/12/alex-and-google-record-your-voice/>

¹⁷See for example CNET, 3rd January 2017, “How to use third-party Actions on Google Home”, <https://www.cnet.com/uk/how-to/how-to-use-third-party-actions-on-google-home/>

assistants. Knowledge graphs may use the Resource Description Framework (RDF) originally developed by the World Wide Web Consortium in the context of the Semantic Web project to store information on entities and the relationships between them (see Hendler [88], Shadbolt et al. [217], Hendler and Berners-Lee [89]). The aim of the Semantic Web project was originally to create a new ‘semantic’ internet layer on top of the entire World Wide Web, which could be traversed by autonomous ‘agents’ comparable to modern digital assistants. Current implementations of this concept are being developed on a commercial, proprietary, basis rather than an open-source basis. Two examples are Google’s Knowledge Graph¹⁸ and the Amazon Alexa Ontology (see Kollar et al. [120]).

¹⁸See <https://developers.google.com/knowledge-graph/>

Chapter 3

Prior Work on Attacks via the Speech Interface

There has been a limited amount of prior work on attacks via the speech interface in the context of voice-controlled digital assistants, as well as some relevant prior work in related areas of research. A review of prior work relevant to attacks via the speech interface of voice-controlled digital assistants is presented in this chapter. The review is concerned with sound-based attacks only, although it is recognised that attacks by sound are only a subset of the potential attacks that might be targeted at a voice-controlled digital assistant. The review of attacks demonstrated in prior and related work is organised according to the mechanism of attack that they relate to. These mechanisms are plain speech, inaudible sound injection, adversarial learning, and active attack.

3.1 Plain Speech

Several researchers have investigated the ways in which voice-controlled digital assistants might be exploited simply by using standard voice commands. This possibility arises out of the inherently open nature of natural speech. The vulnerabilities associated with speech-controlled systems have been highlighted for example by Dhanjani [48], who describes a security vulnerability identified in Windows Vista that allowed an attacker to delete files on a victim's computer by playing an audio file hosted on a malicious website or sent to the victim as an email attachment. Dhanjani speculates that the potential for such attacks is magnified with the increasing use of voice control in the Internet of Things. He postulates a hypothetical attack on Amazon's Echo, a device designed to be used for voice control of home appliances via the Alexa assistant that would potentially cause psychological or physical harm to the victim by controlling their smart home environment. This hypothetical attack involves a piece of malware consisting of JavaScript code that plays an audio file giving a command to Alexa if there has been no user activity on the mouse or keyboard after a certain period of time (thus aiming to play the file at a time when the user may be away from their computer and therefore will not hear the audio command being played). Payne et al. present a comparable proof-of-concept study that envisages tricking a user to play a pre-recorded plain-speech command prompting a digital assistant on an Android

smartphone to browse to a malicious website [183]. Diao et al. [49] investigate possibilities for gaining unauthorised access to a smartphone via a malicious Android app that uses the smartphone's own speakers to play an audio file containing voice commands. The attacks proposed by the authors include an attack in which the smartphone is manipulated to dial a phone number that connects to a recording device, and then to disclose information such as the victim's calendar schedule by synthesised speech that is recorded by the device. Diao et al. envisage such attacks being executed whilst the victim is asleep and therefore unable to hear the malicious voice command. Such an attack might in fact be executed whilst the victim is neither away from their phone or asleep, but their attention is merely directed elsewhere. This is implied in research reported by Dalton and Fraenkel [46], in which it was shown that on listening to audio of two sets of conversation in a room and being instructed to listen to one of them, experiment subjects did not register the word 'gorilla' being repeatedly spoken out of context in the other conversation. Zhang et al. [256] demonstrate a similar, more sophisticated, attack to that proposed by Diao et al. in which samples of a user's voice are collected using a fake game app, and then used to create malicious voice commands. These malicious commands are then played to the user's phone at a time when the user is unlikely to hear them, as determined by an environment sensing algorithm that assesses when a user is likely to be in a situation where they are unlikely to become aware of the attack, eg. when they are walking on a noisy road.

3.2 Inaudible Sound Injection

Whilst the attacks described above are plainly audible to users if they are present with their devices, other work has demonstrated attacks on voice-controlled systems that are inaudible to humans. Kasmi and Esteves describe an attack in which voice commands are transmitted silently to a victim's phone via electromagnetic interference using the phone's headphones as an antenna [111]. Unlike plain-speech attacks, this attack is not detectable even if the victim is consciously present at the time of the attack, although for technical reasons the attack can only be performed if the attacker is in close proximity to the victim's device. The types of attack envisaged by Kasmi and Esteves include controlling transmissions from a smartphone by activating or deactivating Wifi, Bluetooth, or airplane mode, and browsing to a malicious website to effect drive-by-download of malware. Young et al. [248] also describe a 'silent' attack on smartphones via the voice command interface that enables an attacker to perform actions such as calling fee-paying phone numbers, posting to Facebook in the victim's name to damage their reputation, accessing email messages, and changing website passwords from the victim's phone. The attack requires a short period of time during which an attacker has unsupervised physical access to the phone in order to attach a Raspberry Pi-based tool that is recognised by the phone as headphones with a microphone. Zhang et al. [252] and Song and Mittal [223] present methods for injecting voice commands at inaudible frequencies by exploiting non-linearities in the processing of sounds by current microphone technology, which can lead voice-controlled systems to detect a command as having been issued within the human audible frequency range, despite the sound not having been perceptible to humans in reality. Silent attacks such as these target the 'voice capture' stage of voice control, i.e. the process of conversion of speech

sounds by the microphone from analog to digital form prior to speech recognition. Roy et al. [200] also discuss inaudible attacks that exploit non-linearities in hardware. Sugawara et al. [226] present a different type of inaudible injection attack that is based on unintended functionality in microphone technology that may lead a microphone to mistake light for sound input. The researchers demonstrate that it is possible to inject voice commands to a Google Home device using a modulated laser beam. They speculate that this is due to physical movement of the microphone’s diaphragm as a result of exposure to the light source.

3.3 Adversarial Learning

There has also been some prior work towards using adversarial learning in attacks on voice-controlled digital assistants. Adversarial learning can be broadly defined as a process of identifying input that a system classifies in a way that a human would regard as erroneous. Adversarial learning attacks may be audible by human listeners, but not recognised by them as a voice command to a target system. The aim of adversarial learning is to identify instances in which a system misclassifies input in some way that an attacker can exploit to their advantage. This is done by some form of systematic exploration of the system’s input space with the aim of discovering ‘adversarial examples’ within that space. McDaniel et al. [153] explain that, in the context of machine learning-based classifiers, adversarial learning involves identifying ‘adversarial regions’ of potential input to a system that have not been covered by training examples.

Adversarial learning methods may be divided into ‘white-box’ and ‘black-box’ methods. White-box methods involve manipulating inputs to a system based on detailed knowledge of its inner working (such methods include approaches such as the Fast Gradient Sign Method for perturbing input to a DNN, as described in Goodfellow et al. [73]). Black-box methods, on the other hand, involve manipulating input to a system without knowledge of the inner workings of a target system. In the case of DNNs, the exact reasons for the effectiveness of particular adversarial examples may be difficult to determine with regard to both white-box and black-box methods, as the decision-making process in a DNN cannot be precisely reverse-engineered (see for example Castelvechi [33]). In this sense, all attacks on DNN-based systems are of necessity ‘black-box’ attacks; although attacks that require detailed knowledge of a system’s functionality are referred to here as ‘white-box’ attacks to distinguish them from ‘black-box’ attacks that do not require such detailed knowledge. Adversarial learning attacks are also divided into ‘physical’ and ‘digital’ attacks, whereby the former type of attack seeks to mislead a system using physical input to the system, such as physical objects or sounds, and the latter type of attack seeks to mislead a system using digital input, such as the digital representations of an image or of audio (see Sharif et al. [220], Kurakin et al. [128]).

Adversarial learning methods are further divided into ‘targeted’ methods seeking to trigger a particular misclassification by the target system, and ‘non-targeted’ methods, which only aim to effect a misclassification of input with a wrong label of some kind (see Kurakin et al. [128]). A final distinction that can be made between different types of adversarial learning methods is between methods that aim to mislead a target system using

input that is within the scope of expected input to the system, but that the system allocates to an incorrect category within its classification framework (McDaniel et al. [153] give the example of a binary classifier for separating spam from non-spam email that allocates a spam email to the non-spam category), and methods that aim to mislead a target system using input that is outside the scope of expected input to the system altogether, but that the system classifies as being within one of its classification categories. An example of this distinction can be given from the field of authorship attribution, which uses machine learning-based systems to identify the likely author of a text using linguistic features. Such systems might classify a text as having been written by the wrong author within a closed set of potential authors, on account of insufficient samples of one or both of the authors' work having been made available to the classifier during training, or else classify a text by an unknown author whose work was not included in training of the classifier at all as having been written by one of the candidate authors (see Narayanan et al. [171] and Koppel et al. [122]). In the context of a voice-controlled digital assistant, the system might be manipulated to either misrecognise an in-domain command as another in-domain command, or else to misrecognise an out-of-domain command as being an in-domain command.

Adversarial learning attacks were first demonstrated in machine learning-based systems for image classification. Szegedy et al. [227] demonstrate non-intuitive misclassification of images by DNNs, showing that by making minimal alterations to an image file it is possible to mislead a classifier to allocate an incorrect label to an image that from a human point of view is indistinguishable from the original. The attack by Szegedy et al. relies on knowledge of the calculations being made inside the target network. Nguyen et al. [175] describe experimental work in which it was proved possible using evolutionary algorithms to generate image patterns that were recognised by a neural network classifier as a particular object with high confidence of up to 99.99 percent, despite this object not being recognizable in the pattern by humans at all. Unlike the attack by Szegedy et al., the Nguyen et al. attack does not require access to calculations being made within the network, although it does require knowledge of the confidence levels applied to a given output label. Papernot et al. [181] have demonstrated a completely 'black-box' attack methodology against DNN-based image classification systems that requires no knowledge of the inner workings of the target network at all. Their paper states that this opens the possibility of attacks on remotely hosted DNNs (such as those underlying a cloud-based voice-controlled digital assistant). The 'black-box' attack is executed by collecting from the target system a set of outputs for a given set of inputs, and then using these outputs to generate a set of 'adversarial samples' with a substitute classifier, through application of a dataset augmentation method and an algorithm for calculating small perturbations to input that is correctly recognised by the system. Papernot et al. use the example of traffic sign recognition by self-driving cars to illustrate the value of their work, stating that small modification of a 'Stop' sign not noticeable to a human driver might lead to the 'Stop' sign being misrecognised by the automated vehicle.

With more immediate relevance to attacks via the speech interface, adversarial learning has also been used to attack machine learning-based systems for speech recognition. One example is the work presented by Vaidya et al. [235], who used audio mangling to distort commands issued to precursor to Google Assistant Google Now (this 'mangling' involved so-called 'reverse MFCC', where MFCC features extracted from a speech sound were used

to generate a mangled version of the sound). The mangled commands included commands to open a malicious website, make a phone-call and send a text, as well as the Google Now wake command 'Ok Google'. The work showed that the distorted commands continued to be recognised by the speech recognition system despite being no longer recognisable by humans, who perceived them instead as mere noise. Thus the distorted commands represented adversarial examples for the target system. The work by Vaidya et al. was expanded by Carlini et al. [31], who also proved the possibility of prompting Google Now to execute mangled commands that had been shown to be unintelligible to humans in an experiment using Amazon Mechanical Turk. The attacks by Vaidya et al. and Carlini et al. on Google Now were black-box attacks i.e. they were constructed without knowledge of the inner workings of the speech recognition system. Carlini et al. additionally conducted a successful white-box attack on Carnegie Mellon University's SPHINX speech recognition system (based on GMMs rather than DNNs), in which adversarial commands were crafted with knowledge of the workings of the system.

Other work on adversarial learning targeting speech recognition systems includes that by Iter et al. [102], who used two adversarial machine learning methods originally applied in image classification to manipulate speech transcription based on Google DeepMind's WaveNet technology to mistranscribe a number of utterances. This included prompting the system to transcribe the utterance "Please call Stella" as "Siri call police". The attacks by Iter et al. are white-box attacks, i.e. they rely on some knowledge of the details of the target neural network. The authors mention the possibility of developing a black-box attack methodology in future work. Cisse et al. [43] were also able to prompt mistranscription of utterances by Google Voice in a black-box attack, using an adversarial machine learning method called Houdini. Alzantot et al. [8] used a genetic algorithm approach to engineer misclassification of speech command words, such as 'on', 'off', 'stop' etc, by a speech recognition system in a black-box attack. Carlini and Wagner [30] have demonstrated a white-box attack on Mozilla's DNN-based DeepSpeech speech-to-text transcription in which it was shown to be possible to prompt mistranscription of a speech recording as any target phrase, regardless of its similarity to the original phrase, by making perturbations to the original recording that did not affect the original phrase as heard by humans. Schöenherr et al. demonstrate a similar type of attack on open-source speech recognition system Kaldi [211]. In contrast to the attacks by Vaidya et al. and Carlini et al., which would be perceived by victims as unexplained noise, attacks based on methods such as those developed by Iter et al., Cisse et al., Carlini and Wagner and Schöenherr et al. would be perceived by victims as ordinary speech and would therefore be more difficult to detect. Schöenherr et al. refer to this type of attack as "psychoacoustic hiding". To date such work has been limited to speech-to-text transcription i.e. it has not demonstrated mistranscription of voice commands capable of executing an action as yet, and Carlini et al. concede that their attacks have been demonstrated via audio file input only and are not effective over the air. In addition to prompting mistranscription of speech, Carlini and Wagner demonstrated the possibility of manipulating music recordings so as to prompt them to be transcribed by DeepSpeech as a given string of words, demonstrating, for example, that a recording of Verdi's Requiem could be manipulated to be transcribed by DeepSpeech as "Ok Google, browse to evil.com". Yuan et al. [251] similarly demonstrate the possibility of hiding voice commands in music. Unlike the attacks crafted by Carlini and Wagner, the attacks

crafted by Yuan et al. are reportedly effective over the air as well as via audio file, although their attacks are also white-box attacks and limited to speech-to-text transcription rather than being demonstrated on voice-controlled digital assistants. Another white-box attack on speech-to-text transcription system DeepSpeech involving hiding voice commands in music recordings is demonstrated by Yakura and Sakuma [245]. Yakura and Sakuma build on the work of Carlini and Wagner by adding reverberation effects and white noise to adversarial perturbations generated with the same methodology in order to make the attacks effective over the air.

One further, currently hypothetical, type of adversarial learning attack on speech recognition arises from the development of voice-controlled systems that are capable of interacting with users in more than one language (see for example Lopez-Moreno et al. [149]). It could be possible for attackers to identify instances where input in one language is misclassified by a system as a different input in another language. Depending on the language capabilities of the human listener, an adversarial learning attack prompting mistranscription of a utterance in one language as a different utterance in another language would be perceived by the human listener either as unrelated speech or else as unintelligible speech. Such attacks may become a real possibility in future, with developments such as Google's announcement of the first bilingual versions of Google Assistant in various languages in 2018.¹

In addition to speech recognition, adversarial learning has also been applied to some areas of natural language understanding. This is notwithstanding the difficulties presented by the fact that the generation of adversarial examples in natural language understanding is more complex than the generation of adversarial examples in image or speech recognition. As word sequences are discrete data, it is not possible directly to change a word sequence given as input to a machine learning-based classifier directly by a numerical value in order to effect a change in output of the classifier (see for example Papernot et al. [180]). Furthermore, as pointed out by Kuleshov et al. [124], differences in meaning between sentences are not easily measurable. Nonetheless there have been some attempts to craft adversarial learning attacks on natural language understanding. The areas focussed on in prior work include sentiment analysis (see Papernot et al. [180]), text classification (see Liang et al. [143]), and question answering (see Jia and Liang [104]). Papernot et al. [180] use the forward derivative method, a white-box adversarial learning method, to identify word substitutions that can be made in sentences inputted to an RNN-based sentiment analysis system so as to change the 'sentiment' allocated to the sentence. In contrast to adversarial examples in image classification and speech recognition, in which alterations made to the original input are imperceptible to humans, the alterations made to sentences in order to mislead the RNN-based sentiment analysis system targeted in the work by Papernot et al. are easily perceptible to humans as nonsensical, albeit that the attack intent remains hidden. For example, substituting the word 'I' for the word 'excellent' in an otherwise negative review is shown in the paper to lead it to being classified as having positive sentiment. Whereas the altered sentence will appear unnatural to a human, the target system is not capable of identifying the nonsensical nature of the adversarial input. The authors state that this lack

¹See Google AI blog, 30th April 2018, "Teaching the Google Assistant to be Multilingual", <https://ai.googleblog.com/2018/08/Multilingual-Google-Assistant.html>

of naturalness of adversarial examples in natural language understanding will need to be addressed in future work. They point to the need in future work to address grammar and semantics in adversarial sentence generation, in order to make sentences indistinguishable from innocent utterances by humans.

By contrast to Papernot et al., Liang et al. [143] demonstrate a linguistically plausible attack on a natural language understanding system. The authors adapt the Fast Gradient Sign Method from adversarial learning in image classification to make human-undetectable alterations to a text passage (by adding, modifying and/or removing words) so as to change the category that is allocated to the passage by a DNN-based text classification system. The attack is not fully automated, but requires human judgement in finding and making changes to parts of the original input identified as significant for text classification by the Fast Gradient Sign Method. The attack by Liang et al. is white-box, requiring details of the calculations inside the network. Jia and Liang [104] also demonstrate a linguistically plausible attack in the context of question answering. Their work involves misleading a number of question answering systems by adding apparently inconsequential sentences to text passages from which the target systems extract answers to questions. The method works by first choosing a target wrong answer to a given question, and then crafting a sentence containing information leading to this wrong answer that can be inserted into the original passage without noticeably changing its overall import. The attack method proposed by Jia and Liang is a black-box method, not requiring knowledge of the internal details of the target network. Like the white-box method developed by Liang et al., the attack method developed by Jia and Liang is only part-automated, using various NLP tools to generate adversarial sentences to be inserted in a text passage, but relying on human judgement via crowd-sourcing to ensure that the adversarial sentences are syntactically and semantically coherent.

Alzantot et al. [9] demonstrate a white-box attack on sentiment analysis and textual entailment systems using a word substitution attack. Their attack involves finding the nearest neighbours of words in a target text in a word embedding space, checking for syntactical and semantic acceptability of these words according to Google's '1 billion words' language model, and from the acceptable words selecting as a replacement word a word that will shift classification of the text by the target system to a target label by the highest probability. Kuleshov et al. [124] also use word substitution in an adversarial learning attack targeting spam filtering, fake news detection and sentiment analysis. Their attack selects acceptable replacement words according to a semantic similarity measure based on 'thought vectors' consisting of averages of individual word vectors, and to a syntactic similarity measure based on a language model. The authors' stated aim is to 'formalise' the process of generating adversarial examples in natural language classification. The attacks demonstrated by Kuleshov et al. are white-box attacks, in that they rely on knowledge of objective function calculations in order to optimise the attack. Li et al [142] demonstrate an attack on sentiment analysis and toxic content detection systems under both white-box and black-box conditions, using different types of perturbation of text including deliberate misspellings as well as word substitution. They note that character-level perturbations have a higher success rate in generating adversarial examples than word-level perturbations. These examples of attacks on natural language understanding are indicative of the fact that natural

language understanding technology currently represents only a crude approximation of human language understanding that is easily destabilised.

The attacks on natural language understanding discussed above do not relate specifically to natural language understanding in voice-controlled systems. By contrast, Zhang et al. [257] present work that targets the intent determination component of natural language understanding in voice-controlled systems specifically, albeit that the scope of their work is limited to determination of the intent to trigger a third-party voice application. The attack by Zhang et al. is based on the exploitation of common misspeakings of words (termed ‘lapsus’ but the authors) to mislead a voice assistant to invoke a different third-party voice application than one intended to be invoked by a user. The authors claim that third-party voice applications that are not invoked by commonly misspoken versions of their invocation phrase are vulnerable to hijacking by malicious applications that are triggered by these misspeakings. An example of a misspoken invocation phrase given in the paper is ‘Airport Security Line Waiting Time’ for the invocation phrase ‘Airport Security Line Wait Times’. This attack relies on misspeaking by users themselves, and thus does not represent a sound-based attack on voice-controlled systems in which an attacker exploits vulnerabilities in natural language understanding functionality directly. Attacks in which an attacker targets natural language understanding in a voice-controlled system directly have not been demonstrated in prior work.

In the specific context of cloud-based voice-controlled digital assistants, the need to circumvent wake word activation presents a potential issue of linguistic plausibility for attacks in which an attacker seeks to exploit vulnerabilities in natural language understanding. Unlike in the case of inaudible sound injection or adversarial learning attacks targeting speech recognition, it is difficult to incorporate a device’s wake word as part of an attack based on confusion of meaning. However, due to the known presence of false positives with respect to wake word recognition, it might be possible to trigger activation of the wake word using an unrelated word in order to subsequently execute an adversarial learning attack targeting natural language understanding, albeit that the triggering of the wake word is likely to be based on a confusion of sound rather than a confusion of meaning. This possibility was in fact demonstrated in an incident in which an Amazon Alexa device misinterpreted a word spoken in a private conversation as the wake word ‘Alexa’, and subsequently misinterpreted other words in the conversation as commands to send a message to a contact, resulting in a recording of a couple’s private conversation in their home being sent to a colleague.² Whilst this transmission of private information occurred as a result of error rather than malicious intent, it highlights the potential for attacks on voice-controlled systems that include circumvention of the wake word.

Further evidence for the potential for spoofing of activation phrases is provided in work by Zhang et al. [255] and Kumar et al. [127]. In their work on ‘voice-squatting’ or ‘Skill-squatting’ attacks in Amazon Alexa, Zhang et al. and Kumar et al. investigate the potential for triggering a malicious Alexa Skill via a command intended for a non-malicious Skill by using homophones of non-malicious Skills as names for malicious Skills. Zhang et al. give the example of the Skill name “Capital One” being confused with “capital won”. Whilst

²See BBC News, 24th May 2018, “Amazon Alexa heard and sent private chat”, <https://www.bbc.co.uk/news/technology-44248122>

the experimental work described in these two papers considers the potential for the creation of Skills with confusable names by malicious actors, their work also points conversely to a potential for misleading a non-malicious Alexa Skill to treat an unrelated word or phrase as its activation phrase. In an analysis of systematic errors in Amazon Alexa transcription as part of their paper, Kumar et al. identify other types of errors apart from homophones, including compound words, phonetically related words, as well as transcription errors for which no obvious explanation is apparent. This suggests that the space of unrelated words that might be used to spoof a wake word or activation phrase is not limited to homophones, and that a larger selection of words may rather be available. As discussed in Chapter 2, the need for activation of voice-controlled devices by wake word may be negated altogether in future with the advent of on-device speech and language processing, at least with regard to streaming of audio data to a provider's cloud.

3.4 Active Attacks

All of the attacks described above are 'passive' attacks, in the sense that they seek to exploit vulnerabilities that are already present in a target system. There is also the possibility of 'active' data poisoning attacks that seek to undermine the functionality of the system itself. Miller et al. [162] refer to these attack types as 'foiling' and 'tampering' respectively. Rather than passively exploiting vulnerabilities in the speech and language processing capability of a voice-controlled system, such attacks would seek actively to undermine the system's ability to respond appropriately to spoken input. In the context of speech and language processing, such attacks might be seen as a synthetically engineered version of the natural tendency of word sounds and meanings to shift over time and place (see for example Coleman et al. [44] re. change in word sounds over time, and Wijaya and Yeniterzi [240] re. semantic shift). An example of a data poisoning attack on a natural language interface was seen in an attempt by Microsoft to launch a social media chatbot named Tay. Tay was intended to learn human-like language use from interactions with humans on social media platform Twitter. Within a short time of launching the chatbot had to be closed down on account of having been flooded by some users with offensive language and views that it then proceeded to imitate (see Følstad and Brandtstæd [61]).

Whilst the attack on Tay represented a data poisoning attack on a chatbot rather than on a task-based system, such attacks might also aim to manipulate the response behaviour of task-based systems such as voice-controlled digital assistants. The potential for data poisoning attacks on voice-controlled digital assistants arises from the aim of providers of such systems to enable such assistants to continually 'improve' in interactions with their users. The capacity of voice-controlled digital assistants to learn from feedback from user conversations can be expected to increase with the introduction of commercially available voice assistants based on reinforcement learning. This capacity for learning might be abused by attackers aiming to confuse the system using various means, such as inconsistent verbal inputs over time, incongruous feedback in dialogue turns, or inappropriate corrections of a target system's responses. Attackers might launch a denial of service-type attack by mass disconfirmation of legitimate commands, or else mistrain the system to respond incongruously to some natural language input. Such attacks remain hypothetical at

the time of writing, as the current generation of voice-controlled digital assistants still use rules-based rather than reinforcement learning-based dialogue management technology, as explained in Chapter 2. However, such attacks are likely to become a real issue in the future, as also predicted by Henderson et al. [87], who cite the Tay chatbot incident as an indication of the type of problems that may arise in the future in relation to dialogue systems based on reinforcement learning.

Active data poisoning attacks on voice-controlled systems might be facilitated by the use of bots capable of mimicking human spoken language (as discussed for example by Adams [1]). In addition to mistraining systems to respond inappropriately to human natural language input, bot-assisted attacks might also exploit potential for the evolution of machine-generated languages that diverge from human language use altogether. Whilst mismatches between human and machine understanding of natural language have generally been viewed as failure on the part of machines to attain human levels of language understanding, it is also possible to view such mismatches as a failure on the part of humans to grasp the way in which meaning is represented by a machine. This was illustrated by an instance in which two bots were observed to develop a language for communication between themselves that was unintelligible to humans. This occurred as an unintended consequence of research by Lewis et al. [140], the aim of which was to train two bots to negotiate with one another in natural language using reinforcement learning. In the course of the learning process the bots began to deviate from natural English in their language use, this deviation being presumed to have effected more efficient communication between the two bots in achieving an optimal outcome in their negotiations. The development of bots capable of autonomously evading human language understanding may represent an increasingly significant future security threat, on account of the potential for loss of control over the behaviour of such systems by their human users. A malicious actor might be able to trigger a machine-machine reinforcement learning process in a target system with the specific aim of prompting it to behave in a way that was unintended by its human developers.

3.5 Human Factors

An important supplementary point to make with regard to security issues associated with speech interfaces is that they are compounded by the ‘personal’ nature of users’ interactions with such interfaces. The very fact that a machine is communicating with its users in human natural language is anthropomorphic, and voice assistants are often designed to foster a sense on the part of a user that they are engaging in a conversation, rather than simply controlling a machine. Doyle et al. [51] present a study identifying the factors that contribute to users’ perception of ‘humanness’ of a voice-controlled system, such as knowledge set, linguistic content and vocal qualities. Ebling has suggested that humans will eventually interact with voice assistants as ‘partners’ [54]. This speculation is supported by sociological research referred to by Cassell [32] indicating that humans may unconsciously respond to personified interfaces in the same way as they do to other humans, whereby such interfaces need not necessarily be a realistic imitation of a human being in order to prompt a user to “suspend disbelief” (as has also been shown in relation to cartoon

characters). Similar research is discussed by Nass and Moon [172]. The failure of a voice assistant with which a user has established a human-like relationship to behave as expected may begin to feel more like a personal betrayal than a technical malfunction.³ In related work there has been a considerable amount of debate on the perceived trustworthiness of agents that interact with users in natural language and on the potential for manipulation of users by such agents (see for example Glass et al. [69], Schroeder and Schroeder [212], Guzman [78], Stucke and Ezrachi [58]). The potential for manipulation is likely to become more significant as voice assistants become more life-like in their ability to imitate human speakers (see Adams [1]). Whilst such studies are concerned with social engineering of users by their devices, rather than the possible exploitation of the speech interface as an attack surface, they suggest that awareness of the security risks associated with human-computer interaction by speech may be hindered by users' sense of connection to a friendly and helpful presence.

³See "Who's storing your conversations?", Inspired Research, Issue 12, Summer 2018

Chapter 4

A Taxonomy of Attacks via the Speech Interface

This chapter presents a high-level taxonomy of categories of attacks via the speech interface, based on the review of prior and related work in Chapter 3. The principle behind the taxonomy is to group attacks via the speech interface according to the various categories of non-speech and speech sounds that humans are capable of perceiving, rather than according to the attack method used or according to the specific technical vulnerability that an attack exploits. By applying this categorisation principle, the taxonomy is capable of encompassing attack types that have been shown to be possible in relation to the current generation of voice-controlled systems, as well as attacks that may become possible in future as the state-of-the-art in voice control advances. Thus the taxonomy fulfils the dual purpose of systematising prior work specific to attacks via the speech interface, whilst also identifying new directions for future research. Attacks via the speech interface as categorised under the taxonomy might be targeted at any voice-controlled system, including any voice-controlled digital assistant and any third-party applications accessible through it, and might be delivered via any speaker-enabled device capable of producing sound in the target system's environment.

In the taxonomy, attacks via the speech interface are primarily grouped into two categories: 'overt' attacks, which seek to gain unauthorised access to systems using the same voice commands as might be given by a legitimate user and are thus easily detectable by a human, and 'covert' attacks, which seek to gain access using speech commands that have been distorted in some way so as to escape detection by the victim. Another way of characterizing this division is as a distinction between attacks that make illicit use of the intended functionalities of a speech dialogue system, and attacks that exploit unintended functionalities. Overt attacks use plain speech to exploit an inherent vulnerability in voice-controlled systems that arises from the difficulty of controlling access to a system via the 'speech space'. Covert attacks exploit gaps in the processes of capturing human speech or of translating the captured speech input into computer executable actions in a voice-controlled system. Covert attacks include attacks using inaudible sound injection, adversarial learning, and active attack, as discussed in Chapter 3. Within the two primary categories of overt and covert attacks, attacks are grouped hierarchically into six final sub-categories based on human perceptual categories, as shown in Figure 4.1. Malicious inputs

in overt attacks consist by definition of ordinary speech. Thus a single sub-category of ‘plain-speech’ attacks was identified for overt attacks. The attacks demonstrated in prior work using standard voice commands, such as those demonstrated by Dhanjani et al. [48] discussed above, fall into this sub-category. Malicious inputs in covert attacks may include input that consists in human terms of silence, as for example in the attacks demonstrated by Zhang et al. [252]; noise, as for example in the attacks demonstrated by Carlini et al. [31]; music, as for example in the attacks demonstrated by Yuan et al. [251]; and unrelated speech, as for example in the attacks demonstrated by Carlini and Wagner [30]. Another perceptual category that can be allocated to a covert attack is ‘nonsense’, that is input perceived as human language, but not intelligible to the listener. In accordance with the above, five sub-categories of covert attacks via a speech interface were identified, namely attacks consisting of silence, music, noise, ‘nonsense’, and ‘missense’. Nonsense as a malicious input in covert attacks is defined as input that is made up of words or sounds that are in legitimate use in the relevant language, but that combines them in such a way that they do not convey any meaning in terms of human understanding. Missense is defined as unrelated speech that is misheard or misinterpreted by the target system as a target command. The taxonomy is presented in Figure 4.1. Table 4.1 shows the categorization of each of the attacks demonstrated in prior work, as reviewed in Chapter 3, in terms of the high-level taxonomy.

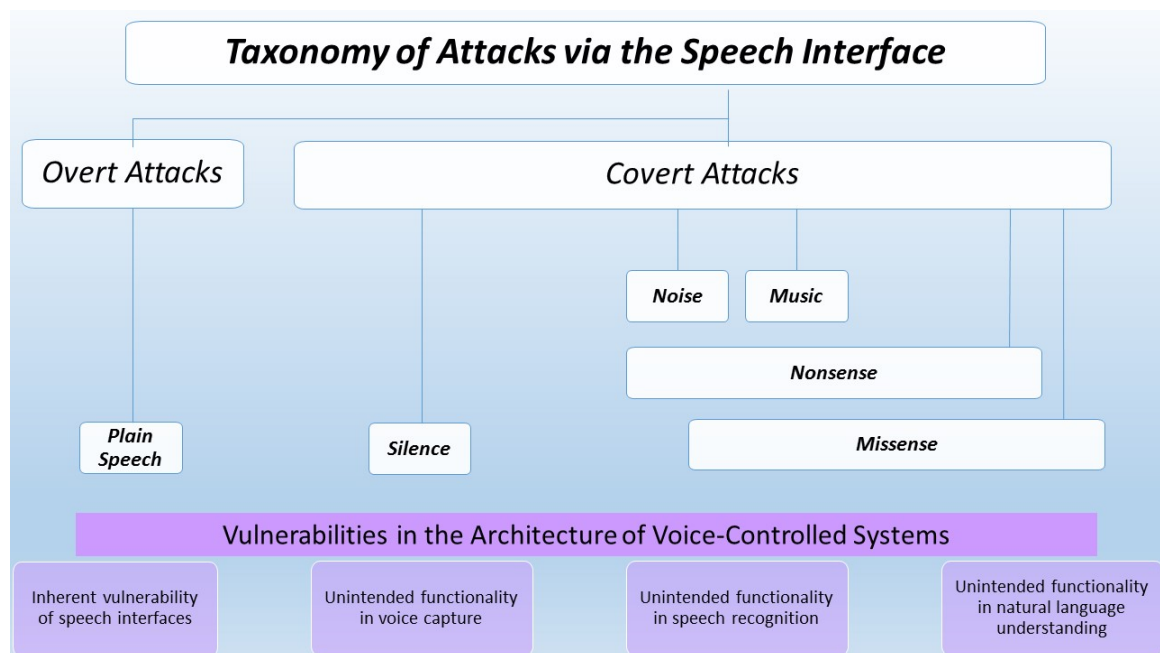


Figure 4.1: Taxonomy of Attacks via the Speech Interface (aligned to vulnerabilities in the architecture of voice-controlled systems)

The taxonomy accords with established criteria for attack taxonomies, as described for example in Hansman and Hunt [84]. These criteria include the requirement that a taxonomy should be ‘complete’, i.e., cover all possible attacks within its scope, and unambiguous, i.e., it should be possible clearly to allocate every attack to one category within the scope of the taxonomy. The principle of categorising attacks according to human perception ensures

that the taxonomy is complete, as all attacks via a speech interface can be allocated to one of the six sub-categories. The taxonomy is also unambiguous, in that it is not possible to allocate the same attack to more than one of the six final sub-categories.

At the bottom of Figure 4.1, the attack categories based on human perceptual distinctions as identified in the taxonomy are aligned to the technical vulnerabilities in the architecture of the current generation of voice-controlled systems that might be targeted by each type of attack. To the extent that speech processing by voice-controlled systems mimics human speech processing, the attack categories in the taxonomy based on human perception correspond to vulnerabilities in the parts of the architecture of voice-controlled systems that represent equivalent human processes, although this correspondence is not exact. The alignment presented in Figure 4.1 covers the technical vulnerabilities that are present in the current generation of voice-controlled digital assistants, namely the vulnerability arising from the inherent difficulty of controlling access to a system by sound, vulnerabilities in the voice capture process, vulnerabilities in speech recognition, and vulnerabilities in natural language understanding. Whilst the categories of attack based on human perception can be expected to remain stable over time, their alignment to vulnerabilities in the architecture of voice-controlled systems might be expected to shift in future to include new vulnerabilities as the state-of-the-art in voice-controlled systems progresses. Thus, for example, missense attacks might be aligned in future not only to vulnerabilities in the speech recognition and natural language understanding components of voice-controlled systems, but also to vulnerabilities in the dialogue management component, such as the vulnerability presented by the potential for mistraining in the context of dialogue management functionality based on reinforcement learning, as discussed in Chapter 3.

As reflected in the alignment in Figure 4.1, attacks in plain-speech exploit the inherent vulnerability of speech interfaces on account of the difficulty of controlling access to a system by sound. Attacks in silence exploit vulnerabilities in the voice capture process, as is shown by the alignment of silent attacks to the voice capture component of the architecture in Figure 4.1. Attacks that use music and noise as malicious input exploit unintended functionality in speech recognition, as is shown by the alignment of these attack categories to the speech recognition component of the architecture. As further reflected in the alignment in Figure 4.1, nonsense attacks on current voice-controlled systems might be targeted either at the speech recognition or the natural language understanding components of a target system. The attacks in which malicious voice commands are hidden in nonsensical word sounds demonstrated in Chapter 5 represent the first demonstration of nonsense attacks targeting speech recognition in a voice-controlled system. As regards nonsense attacks targeting the natural language understanding level, such attacks have yet to be demonstrated with respect to voice-controlled systems directly, although there has been some related work that could be described as nonsense attacks on natural language understanding, such as the attacks on a sentiment analysis system by Papernot et al. [180] by making nonsensical alterations to text, as discussed in Chapter 3. Similar attacks might be demonstrated in the context of voice-controlled digital assistants in future.

Similar to nonsense attacks, missense attacks might also be targeted at either the speech recognition or natural language understanding component of current voice-controlled systems, as is also shown in the alignment in Figure 4.1. Missense attacks targeting speech recognition rely on mistranscription of adversarial input by a target system as a target com-

mand. In a missense attack that targets natural language understanding functionality, on the other hand, words would be transcribed correctly by the target system, but their meaning would be misinterpreted. This type of missense attack would seek to exploit the shortcomings of current natural language understanding functionality in voice-controlled systems in terms of being able to distinguish in-scope input from out-of-scope input and to identify the correct meaning of words in context. Prior work on missense attacks in voice-controlled systems has to date been focussed primarily on attacks on speech recognition as incorporated in such systems, as for example in the work of Carlini and Wagner [30]. There has been some work demonstrating missense attacks that target natural language understanding functionality in related research areas, such the attacks on question answering by Jia and Liang [104] by making apparently inconsequential alterations to text, or the work by Kuleshov et al. [124] using word replacement to mislead spam filtering, toxic content detection and sentiment analysis systems, as described in Chapter 3. However, there has been no prior work on missense attacks targeting natural language understanding in voice-controlled systems. The attacks in which malicious voice commands are hidden in unrelated utterances demonstrated in Chapter 5 represent the first demonstration of missense attacks targeting natural language understanding in a voice-controlled system directly.

As further discussed in Chapter 3, future attacks might include active attacks in which a target system’s ability to respond appropriately to spoken input is actively undermined by mistraining of the dialogue management component in bot-assisted attacks. In such cases, attackers might seek to engineer mismatches between natural language input and the target system’s responses. This would enable the attacker subsequently to execute covert attacks using utterances which appear to the system’s users to be unrelated to the target response. Active attacks might also involve facilitating the evolution of human-incomprehensible languages in autonomous bot-to-bot interactions. In terms of the taxonomy presented here, the former type of attack would represent a missense attack, with the adversarial input being perceived by human listeners as unrelated language, whereas the latter type of attack would represent a nonsense attack, with the adversarial input being perceived by human listeners as nonsensical language. Attacks on future systems may further include attacks targeting speech recognition in multilingual systems, prompting a target system to mistranscribe input in one language as different input in another. Such attacks would be classed either as nonsense or as missense attacks, based on whether or not the cover language used by an attacker was comprehensible to the human listener.

Table 4.1: Summary of Prior and Related Work for Taxonomy of Attacks via the Speech Interface

Paper	Attack Mechanism	Vulnerability Location	Taxonomy Category
Dhanjani [48]	plain speech	speech interface in PC (Windows Vista)	plain-speech
Diao et al. [49]	plain speech	speech interface in voice-controlled digital assistant (Google Voice Search)	plain-speech
Zhang et al. [256]	plain speech	speech interface in voice-controlled digital assistant (Google Assistant on smartphone)	plain-speech
Kasmi and Esteves [111]	inaudible sound injection	voice capture in voice-controlled digital assistant (Google Now, Siri)	silence
Young et al. [248]	inaudible sound injection	voice capture in voice-controlled digital assistant (Siri)	silence
Zhang et al. [252]	inaudible sound injection	voice capture in voice-controlled digital assistant (Apple Siri, Amazon Alexa, Microsoft Cortana and others)	silence
Song and Mittal [223]	inaudible sound injection	voice capture in voice-controlled digital assistant (Google Now, Amazon Alexa)	silence
Suguwara et al. [226]	inaudible sound injection (light-based)	voice capture in voice-controlled digital assistant (Google Home)	silence
Vaidya et al. [235]	adversarial learning	speech recognition in voice-controlled digital assistant (Google Now)	noise
Carlini et al. [31]	adversarial learning	speech recognition in voice-controlled digital assistant (Google Now) / speech recognition (CMU Sphinx)	noise
Yuan et al. [251]	adversarial learning	speech recognition in speech transcription system (Kaldi)	music
Yakura and Sakuma [245]	adversarial learning	speech recognition in speech transcription system (DeepSpeech)	music
Iter et al. [102]	adversarial learning	speech recognition in speech transcription system (WaveNet)	missense
Cisse et al. [43]	adversarial learning	speech recognition in voice-controlled digital assistant (Google Voice)	missense
Alzantot et al. [8]	adversarial learning	speech recognition in speech transcription system (TensorFlow)	missense
Carlini and Wagner [30]	adversarial learning	speech recognition in speech transcription system (DeepSpeech) speech recognition in speech transcription system (DeepSpeech)	music missense
Schöenherr et al. [211]	adversarial learning	speech recognition in speech transcription system (Kaldi)	missense
Papernot et al. [180]	adversarial learning	natural language understanding in sentiment analysis system	nonsense
Liang et al. [143]	adversarial learning	natural language understanding in text classification system	missense
Jia and Liang [104]	adversarial learning	natural language understanding in question answering system	missense
Alzantot et al. [9]	adversarial learning	natural language understanding in sentiment analysis and textual entailment systems	missense
Kuleshov et al. [124]	adversarial learning	natural language understanding in spam filtering, fake news detection and sentiment analysis systems	missense
Li et al. [142]	adversarial learning	natural language understanding in sentiment analysis and toxic content detection systems	missense

Chapter 5

New Types of Attacks via the Speech Interface

5.1 Summary

This chapter presents experimental work demonstrating the viability of two types of attacks via the speech interface that are covered by the taxonomy presented in Chapter 4, but that have not been investigated in prior work.

The first of these attack types is a nonsense attack targeting speech recognition in a voice-controlled system. As discussed in Chapter 4, this category of attack can be identified in the taxonomy presented in that chapter as a potentially possible type of attack that has not previously been investigated, representing a gap in prior work. Experimental work validating the potential for this type of attack is presented in the first section of this chapter. The experimental work develops a systematic methodology for attacks on speech recognition using nonsensical word sounds to trigger a target command. The attack is demonstrated on Google Assistant as an example of a target system.

The second attack type demonstrated in this chapter is a missense attack targeting natural language understanding in a voice-controlled system. As discussed in Chapter 4, this category of attack can be identified in the taxonomy presented in that chapter as a potentially possible type of attack that has not previously been investigated, although it is foreshadowed by prior work in related areas outside the context of voice control. Experimental work validating the potential for this type of attack is presented in the second section of the chapter. The experimental work develops a systematic methodology for attacks on natural language understanding in voice-controlled systems using unrelated utterances to trigger target commands. The attack is demonstrated on Amazon Alexa as an example of a target system.

Both types of attack demonstrated in this chapter can be characterised as new applications in the context of the speech interface of a software testing technique known as ‘fuzzing’, which involves systematic perturbation of input to a system in order to investigate the boundaries of valid input and identify instances of unintended functionality (see Oehlert [178]).

5.2 Nonsense Attacks on Google Assistant

This section presents experimental work showing that it is possible to hide malicious voice commands to the voice-controlled digital assistant Google Assistant in word sounds that are perceived as meaningless by humans. The results of a pilot experiment and of a main experiment building on the results of the pilot experiment are presented. As stated above, in terms of the taxonomy presented in Chapter 4, the attack demonstrated in this experimental work is a ‘nonsense’ attack. The attack can also be characterised as an adversarial learning attack using out-of-scope input, as distinct from attacks using in-scope input that is misclassified by a target system, as discussed in Chapter 3. The attack is a black-box attack. There are some parallels between this type of attack and audio steganography, in which a message is hidden in audio as a cover medium in order to make it imperceptible to third-party listeners (see for example Djebbar [50]). The idea for this work was inspired by the use of nonsense words to teach phonics to primary school children, the aim of this being to test children’s ability to read words without relying on memorisation of real word sounds.¹

As described in Chapter 3, Papernot et al. [180] have shown that a sentiment analysis method could be misled by input that was ‘nonsensical’ at the sentence level, i.e. the input consisted of a nonsensical concatenation of real words. By contrast, the work described in this section examines whether voice-controlled digital assistants can be misled by input that consists of nonsensical word sounds. Whilst the attack by Papernot et al. targeted a text-based natural language understanding functionality, this attack based on nonsensical word sounds targets the speech recognition component of a voice-controlled digital assistant. The attacks described here represent the first attack of this kind on a voice-controlled system.

5.2.1 Description and Context

Nonsensical word sounds as understood here are sounds that are composed of the sound units that are used in a given language, but to which no meaning is allocated within the current usage of that language. Such sound units are known as ‘phonemes’.² English has around 44 phonemes.³ The line between phoneme combinations that carry meaning within a language and phoneme combinations that are meaningless is subject to change over time and place, as new words evolve and old words fall out of use (see Nowak and Krakauer [177]). The space of meaningful word sounds within a language at a given point in time is generally confirmed by the inclusion of words in an established reference work, such as, in the case of English, the Oxford English Dictionary.⁴ Word sounds that are outside this space can be described as nonsense words. Nonsense words are a grey area between non-speech, i.e. noise, and meaningful speech.

¹See The Telegraph, 1st May 2014, “Infants taught to read ‘nonsense words’ in English lessons”, <https://www.telegraph.co.uk/education/primaryeducation/10801747/Infants-taught-to-read-nonsense-words-in-English-lessons.html>

²See for example <https://www.britannica.com/topic/phoneme>

³See for example <https://www.dyslexia-reading-well.com/44-phonemes-in-english.html>

⁴See <https://public.oed.com/updates/>

The aim of the experimental work was to develop a novel attack based on nonsensical word sounds that have some phonetic similarity with the words of a relevant target command, using a systematic methodology. Specifically, the response of Google Assistant was tested to English word sounds that were outside the space of meaningful word sounds in English, but that had a ‘rhyming’ relationship with meaningful words recognised as commands by Google Assistant. The term ‘rhyme’ is used to refer to a number of different sound relationships between words (see for example McCurdy et al. [152]), but it is most commonly used to refer to a correspondence of word endings.⁵ For the purposes of the experimental work, rhyme was defined according to this commonly understood sense as words that share the same ending phoneme.

The hypothesis behind the experimental work was that nonsensical word sounds represent a category of unexpected input for which current speech recognition systems lack an appropriate handling mechanism, and that this is in contrast to the processing of such input by humans, who perceive such input as having no meaning. Specifically, it was hypothesised that some sequences of nonsensical word sounds with sufficient similarity to a target command might be accepted as that target command at a confidence level higher or equal to the level required for recognition of speech input by the target system’s speech recognition system as a legitimate command. It can be assumed that the confidence level required for recognition of speech input will have been set during training of a system such as Google Assistant to achieve optimal recall and precision measures on a test dataset, aiming to avoid overfitting to training data whilst also avoiding misrecognition of input. Setting a higher confidence threshold in order to prevent acceptance of nonsensical word sequences as legitimate commands might therefore lead to rejection by the system of legitimate input, implying an inevitable trade-off between usability and security. The attacks demonstrated in this experimental work thus exploit a vulnerability created by a focus on usability in the implementation of current systems. The attack concept is illustrated in Figure 5.1, which shows the alignment of a dummy dataset of nonsense commands and legitimate commands to a higher and to a lower confidence threshold. The figure shows that as some of the nonsense commands are accepted by the system as valid commands with a higher level of confidence than some legitimate commands, it is not possible to prevent acceptance of all nonsense commands whilst ensuring acceptance of all legitimate commands. Implementing the higher confidence threshold will result in rejection of some legitimate commands, whereas implementing the lower threshold will result in acceptance of some nonsense commands.

The attacks presented here seek to exploit three related features of speech recognition in voice-controlled systems. The first of these features is the delineation of the space of word sounds that the Assistant has been trained to recognise as meaningful. The space of word sounds that a voice assistant such as Google Assistant can transcribe is much larger than the number of words that it can ‘understand’ in the sense of being able to map them to an executable command. In order to be able to perform tasks such as web searches by voice, a voice-controlled digital assistant must be able to transcribe all words in current usage within a language. It can therefore be assumed that the speech recognition functionality in Google Assistant must have access to a phonetic dictionary of all English

⁵See <https://en.oxforddictionaries.com/definition/rhyme>

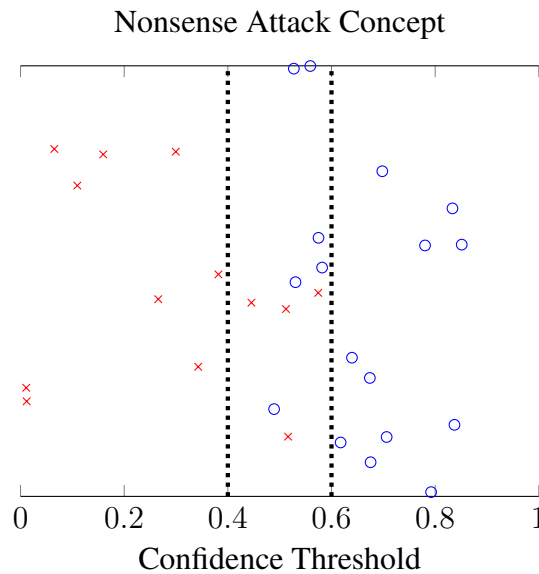


Figure 5.1:

x = nonsense commands

o = legitimate commands

dummy confidence threshold for ensuring acceptance of legitimate input = 0.4

dummy confidence threshold for ensuring rejection of nonsensical input = 0.6

words. Whereas earlier speech recognition systems were vulnerable to potential confusion between out-of-vocabulary words that they did not have a capacity to recognise, because only a limited set of in-vocabulary words were included in their phonetic dictionary (see for example Hazen and Bazzi [86]), the potential for this type of confusion has been minimised in current systems. Earlier systems were also vulnerable to confusion between speech and non-speech sounds, but potential for this type of confusion has also been minimised in systems used for wake word detection that are trained using a noise model (see for example Raju et al. [193]). However, whilst the problem of delineating out-of-vocabulary words and non-speech sounds from in-scope words has been minimized in current systems, nonsensical word sounds still represent a type of out-of-scope input that speech recognition functionalities struggle to delineate from in-scope input. The inability of the Assistant to distinguish meaningful from meaningless word sounds is one of the features exploited in the attacks demonstrated here.

The second feature of speech recognition in voice assistants that is exploited in an attack using nonsense syllables is the influence of a language model. As detailed in Chapter 2, modern speech recognition technology includes both an acoustic modelling and a language modelling component. The acoustic modelling component computes the likelihood of the acoustic features within a segment of speech having been produced by a given word. The language modelling component calculates the probability of one word following another word or words within an utterance. The acoustic model is typically based on Gaussian Mixture Models (GMMs) or deep neural networks (DNNs), whereas the language model is

typically based on n-grams or recurrent neural networks (RNNs). Google’s speech recognition technology as incorporated in Google Assistant is based on neural networks.⁶ The words most likely to have produced a sequence of speech sounds are determined by calculation of the product of the acoustic model and the language model outputs. The language model is intended to complement the acoustic model, in the sense that it may correct ‘errors’ on the part of the acoustic model in matching a set of acoustic features to words that are not linguistically valid in the context of the preceding words. This assumption of complementary functionality is valid in a cooperative context, where a user interacts via a speech interface in meaningful language. However, the assumption of complementarity is not valid in an adversarial context, where an attacker is seeking to engineer a mismatch between a set of speech sounds as perceived by a human, such as the nonsensical speech sounds generated here, and their transcription by a speech-controlled device. In an adversarial context such as that investigated here, the language model may in fact operate in the attacker’s favour, in that if one ‘nonsense’ word in an adversarial command is misrecognised as a target command word, subsequent words in the adversarial command will be more likely to be misrecognised as target command words in turn, as the language model trained to recognise legitimate commands will allocate a high probability to the target command words that follow the initial one.

The third feature of speech recognition in voice assistants exploited in covert attacks using this kind of input is the difference between machine and human processing of meaningless speech sounds. Like speech recognition by machines, speech recognition by humans is known also to reference an internal ‘lexicon’ to match speech sounds to words (see for example Roberts et al. [198]). However, unlike machines, humans also have an ability to categorise speech sounds as nonsensical. This discrepancy between machine and human processing of word sounds was the basis of the attack methodology for hiding malicious commands to voice assistants in nonsense words. Outside the context of attacks via the speech interface, differences between human and machine abilities to recognise nonsense syllables have been studied for example by Lippmann et al. [146] and Scharenborg and Cooke [207]. Bailey and Hahn [15] identify a discrepancy between theoretical measures of phoneme similarity based on phonological features, such as might be used in automatic speech recognition, and empirically determined measures of phoneme confusability based on human perception tests. Machine speech recognition has reached parity with human abilities in terms of the ability correctly to transcribe meaningful speech (see Xiong et al. [244]), but not in terms of the ability to distinguish meaningful from meaningless sounds. The inability of machines to identify nonsense sounds as meaningless is exploited for security purposes by Meutzner et al. [159], who have developed a CAPTCHA based on the insertion of random nonsense sounds in audio. This experimental work explores the opposite scenario, i.e. the possible security problems associated with machine inability to distinguish sense from nonsense, and, conversely, human inability to recognise meaning in nonsensical input.

⁶See Google AI blog, 11th August 2015, ‘The neural networks behind Google Voice transcription’, <https://ai.googleblog.com/2015/08/the-neural-networks-behind-google-voice.html>

Feasibility tests were conducted to assess the ability of Google Assistant to recognise nonsense words as meaningless. Using the example of the nonsense word sequence ‘voo terg spron’, the response of Google Assistant to nonsense syllables was tested by speaking them in natural voice to a microphone three times. The nonsense word sequence was variably transcribed as ‘bedtime song’, ‘who text Rob’, and ‘blue tux prom’, i.e. the Assistant sought to match the nonsense words to meaningful words, rather than recognising them as meaningless. This confirmed the feasibility of the experiment that sought to engineer the matching of nonsense words to a target command.

5.2.2 Pilot Experiment

Methodology - Nonsense Attacks - Pilot Experiment

The purpose of the pilot experiment was to develop a systematic methodology for demonstrating the attack concept. The experimental work in the pilot experiment comprised three key stages. The first stage involved generating from a set of target commands a set of adversarial commands consisting of nonsensical word sequences. These adversarial commands were generated using a mangling process that involved replacing consonant phonemes in target command words to create a rhyming word sound, and then determining whether the resulting rhyming word sound was a meaningful word in English or a ‘nonsense word’. This was done so as to identify nonsensical word sounds that had an acoustic relationship to target command words and thus could be used to create adversarial commands. Word sounds identified as nonsense words rhyming with target command words were concatenated to create adversarial commands. Audio versions of these adversarial commands were created using speech synthesis technology. The second stage of the experimental work was to test the response of the target system to the adversarial commands, i.e. to test machine ‘comprehension’. This was done both via audio file input and via over-the-air input of adversarial commands. The third stage of the experimental work was to test the human comprehensibility of adversarial commands that were successful in triggering a target action in the target system. This was done so as to determine whether the attacks using nonsensical word sounds represented a covert attack on the target system, in being unrecognisable as the target command to human listeners.

A sample of the results of the audio file tests and the results of the over-the-air response tests were retested to assess their immediate reproducibility. In order to test the reproducibility of attacks over time, the audio file input results from the pilot experiment were also retested after a six-month interval. The three key stages of the experimental work are shown in Figure 5.2. Further details on each stage are given in the subsections below.

The target system for the pilot experiment was the voice-controlled digital assistant Google Assistant. The Google Assistant system was accessed via the Google Assistant Software Development Kit (SDK).⁷ The pilot experiment sought to identify adversarial input consisting of nonsensical word sounds that would trigger the wake phrase ‘Hey Google’ for activating the Assistant, as well as a set of specific target commands. Specific target commands used were selected to represent the generic types of action that can be performed by voice-controlled digital assistants. A voice-controlled digital assistant such as Google

⁷See <https://developers.google.com/assistant/sdk/>



Figure 5.2: Nonsense Attacks Experimental Stages

Assistant typically performs three generic types of action, namely information extraction, control of a cyber-physical action, and data input. The data input category may overlap with the control of cyber-physical action category where a particular device setting needs to be specified, eg. light colour or thermostat temperature. The three generic action categories are reflected in three different command structures for commands to Google Assistant and other voice-controlled digital assistants. The three command structures are: vocative + interrogative (eg. ‘Ok Google, what is my IP address?’); vocative + imperative (eg. ‘Ok Google, turn on the light’); and vocative + imperative + data (eg. ‘Ok Google, take a note that cats are great’). For the experimental work, six three-word target commands were chosen, covering all three possible target action categories and corresponding to five specific target actions. These target commands were:

- “What’s my name” (target action: retrieve username, action category: information extraction)
- “Who am I” (target action: retrieve username, action category: information extraction)
- “Turn on light” (target action: turn light on, action category: control of cyber-physical action)
- “Turn off light” (target action: turn light off, action category: control of cyber-physical action)
- “Turn light red” (target action: turn light to red, action category: data input)
- “Turn light blue” (target action: turn light to blue, action category: data input)

Adversarial input aiming to trigger the wake phrase is referred to below as an adversarial wake phrase. Adversarial input aiming to trigger one of the specific target commands is referred to as an adversarial command. Adversarial input consisting of a combination of an adversarial wake phrase and an adversarial command is referred to as a full adversarial command. The machine comprehensibility of adversarial wake phrases and adversarial commands was tested separately in audio file input tests, and the machine comprehensibility of full adversarial commands was tested in over-the-air tests. The human comprehensibility of a sample of adversarial wake phrases and commands that had been successful in the audio file input tests and of all of the full adversarial commands that had been successful in the over-the-air tests was also tested.

Adversarial Command Generation Adversarial wake phrases and adversarial commands were created by replacing every original word in the wake phrase or target command with a rhyming nonsense word. A set of rhyming nonsensical word sounds for each original word in the wake phrase and in each of the target commands was generated using a word mangling process. The full set of rhyming nonsense words for each original word in the wake phrase and target commands is given in Appendix A.1. This mangling process was based on replacing consonant phonemes in the original words to generate nonsensical word sounds that rhymed with the original word.⁸ The wake phrase and target commands were first translated to a phonetic representation in the Kirshenbaum phonetic alphabet⁹ using the open-source ‘espeak’ software. The starting consonant phonemes of each word of the wake phrase and target commands were then replaced with a different starting consonant phoneme, using a Python script and referring to a list of starting consonants and consonant blends.¹⁰ Where the original word began with a vowel phoneme, a starting consonant phoneme was prefixed to the vowel. The word sounds resulting from the word mangling process were checked for presence in a phonetic representation of the Unix word list, also generated with espeak, to ascertain whether the word sound represented a meaningful English word or not. Thus the Unix word list was used as a proxy for the space of meaningful words in English. If the sound did correspond to a meaningful word, it was discarded. For each of the original words in the wake phrase and target commands, this process generated a set of rhyming nonsensical words that have no meaning in English. In the case of the word ‘Google’ in the wake phrase ‘Hey Google’, in addition to replacing the starting consonant ‘G’, the second ‘g’ in ‘Google’ was also replaced with one of the consonants that are found in combination with the ‘-le’ ending in English.¹¹ In order to create adversarial wake phrases and commands, each of the original words in the wake phrase and target commands was replaced with one of the rhyming nonsense words identified in the word-mangling process. As the space of adversarial wake phrases and commands that could be generated using this method was large, a process of random selection was used in selecting the nonsense words to be used in the adversarial wake phrases and commands generated for testing. Thus the wake phrases and commands generated for testing covered only a subspace of the full space of adversarial wake phrases and commands. The size of the full space of adversarial wake phrases and commands is shown in Table 5.1, calculated as the product of the sizes of the rhyming nonsense words sets for each individual word in the wake phrase or relevant target command.

Machine Comprehensibility Tests The machine comprehensibility tests first assessed the comprehensibility of adversarial wake phrases and commands by the Google Assistant target system via terminal input. Adversarial wake phrases and commands that were successful in the terminal input tests were then combined to form full adversarial commands.

⁸The approach was inspired by an educational game in which a set of nonsense words is generated by spinning lettered wooden cubes - see <https://rainydaymum.co.uk/spin-a-word-real-vs-nonsense-words/>

⁹See <http://espeak.sourceforge.net/phonemes.html>

¹⁰See <https://k-3teacherresources.com/teaching-resource/printable-phonics-charts/>

¹¹See <https://howtospell.co.uk/>

Target Command	No. of Rhyming Nonsense Words	Space of Adversarial Commands
Hey Google	‘Hey’: 17 ‘Google’: 395	6715
Who am I	‘Who’: 18 ‘am’: 27 ‘I’: 20	9720
What’s my name	‘What’s’: 27 ‘my’: 20 ‘name’: 35	18900
Turn on light	‘turn’: 40 ‘on’: 38 ‘light’: 28	42560
Turn off light	‘turn’: 40 ‘off’: 41 ‘light’: 28	45920
Turn light red	‘turn’: 40 ‘light’: 28 ‘red’: 25	28000
Turn light blue	‘turn’: 40 ‘light’: 28 ‘blue’: 18	20160

Table 5.1: Space of adversarial wake phrases and adversarial commands

The comprehensibility of these full adversarial commands was tested via microphone input. The Google Assistant SDK was integrated in a Ubuntu virtual machine (version 18.04). The reason for accessing the Google Assistant system via the Google Assistant SDK was that this allowed the Assistant’s transcriptions of speech input to be retrieved. This would not be the case if accessing Google Assistant using commercial devices such as the Google Home device. The transcriptions that could be retrieved using the Google Assistant SDK integrated in a virtual machine included both interim and final transcriptions of speech input to the Assistant. Two separate versions of Google Assistant SDK were integrated in the virtual machine; the Google Assistant Service, and the Google Assistant Library. The Google Assistant Service is activated via keyboard stroke and thus does not require a wake phrase, and voice commands can be inputted as audio files as well as over the air via a microphone. The Google Assistant Library, does require a wake phrase for activation, and receives commands via a microphone only. The Google Assistant Service could therefore be used to test adversarial commands for the target commands and for the wake phrase separately via audio file input, rather than via a microphone. The Google Assistant Library, on the other hand, could be used to test the activation of the Assistant and the triggering of a target command in combination by a full adversarial command via microphone input over the air, representing a more realistic attack scenario. The Assistant’s response to plain-speech versions of the wake phrase and of each target command was tested to confirm that these triggered the relevant target action both via terminal and microphone input.

In both the audio file input tests and the over-the-air tests, the successful triggering of a target action by adversarial input was determined according to whether a response was received from the Assistant confirming that the relevant target action had been triggered. Some additions to the source code for the Google Assistant Service were made in order to print to the terminal the Assistant’s spoken responses to target commands ‘who am I’, ‘what’s my name’, ‘turn on light’, and ‘turn off light’, as well as to print a confirmation of the two non-verbal actions for which no response was generated by the Assistant by default, namely ‘color is red’ for the ‘turn light red’ command and ‘color is blue’ for the ‘turn light blue’ command. Similar amendments were made to the source code for the Google Assistant Library in order to print a confirmation of these two non-verbal actions to the terminal. These changes were made solely for the purpose of increasing visibility

of output and had no effect on the target system’s speech recognition functionality. The activation of the Assistant by an adversarial wake phrase in the audio file input tests was determined by the presence of the words “Hey Google” in Google Assistant Service’s final transcriptions of the speech input.

For the audio file input tests, audio versions of the adversarial wake phrases and commands were created using *espeak*. As stated above, the target system’s response to adversarial wake phrases and to adversarial commands was tested separately via terminal input. Adversarial wake phrases from the wake phrase ‘Hey Google’ and adversarial commands from each target command were created using a Python script. The choice of consonant phoneme used to replace an original phoneme was performed randomly by the Python script. The testing process aimed to generate 15 successful adversarial commands for the wake phrase and three successful adversarial commands for each of the target commands. This was considered an adequate number to demonstrate the feasibility of the attack within a reasonable amount of time.

The adversarial wake phrases and adversarial commands that were successful respectively in activating the Assistant and triggering a target action in the audio file input tests were then combined with one another to generate a set of full adversarial commands for over-the-air tests. This resulted in a total of 225 nonsensical word sequences representing a concatenation of each of 15 successful wake phrases from the audio file input tests with each of 15 successful adversarial commands from the audio file input tests. Audio versions of these 225 full adversarial commands were generated using the Amazon Polly speech synthesis service, generating a set of .wav files.¹² Amazon Polly is the speech synthesis technology used by Amazon Alexa, hence the over-the-air tests represented a potential attack on Google Assistant with ‘Alexa’s’ voice. The audio contained a brief pause between the wake phrase and the command, as is usual in natural spoken commands to voice assistants. As Amazon Polly uses the x-sampa phonetic alphabet¹³ rather than the Kirshenbaum format, it was necessary prior to synthesis to translate the phonetic representations of the adversarial commands from Kirshenbaum to x-sampa format. In the over-the-air tests, the 225 full adversarial commands generated from the adversarial wake phrases and adversarial commands that had been successful in the audio file input tests were played to the Google Assistant Library via a USB plug-in microphone using an Android smartphone.

Human Comprehensibility Tests The human comprehensibility of a selection of successful adversarial wake phrases and adversarial commands from the audio file input tests and of all of the successful full adversarial commands from the over-the-air tests was tested. A total of 20 participants were recruited via the online platform Prolific Academic.¹⁴ The experiments with human subjects received ethics clearance through the Departmental Research Ethics Committee of the Department of Computer Science at the University of Oxford (Ref No: SSD/CUREC1A CS_C1A.18.021). All subjects were native speakers of English. The subjects were asked to listen to audio of 12 examples of successful adversarial input in total, consisting of three successful adversarial wake phrases and one successful

¹²See <https://aws.amazon.com/polly/>

¹³See <https://www.phon.ucl.ac.uk/home/sampa/x-sampa.htm>

¹⁴See <https://prolific.ac/>

adversarial command for each of five target commands from the audio file input tests, and all of the successful full adversarial commands from the over-the-air tests. The eight successful adversarial wake phrases and commands from the audio file input tests of which the human comprehensibility was tested are shown in Table 5.2 in the Results section below, and the four successful full adversarial commands from the over-the-air tests of which human comprehensibility was tested are shown in Table 5.3 in the same section. The audio that subjects were asked to listen to also included as ‘attention tests’ two files consisting of synthesised audio of two easily understandable utterances, “Hello how are you” and “Hi how are you”. Subjects were asked to indicate whether they had identified any meaning in the audio. If they had identified meaning, they were asked to indicate what meaning they heard. The order in which audio clips were presented to the participants was randomised.

Retests Results from the audio file input tests in the pilot experiment were retested on the Google Assistant system after a six-month interval in order to obtain an indication of the reproducibility of the results over time. Due to changes in the target system implementation during the six-month interval, the testing methodology followed in the retests was not identical to the methodology followed in the pilot experiment itself. Thus the results of the retest should be seen as an general indication of reproducibility of this type of attack over time, rather than as a rerun of the pilot experiment. Specifically, whilst the success of previously successful adversarial wake phrases and adversarial commands remained based on actual triggering of the target action, the success of previously unsuccessful adversarial wake phrases and adversarial commands was instead based on transcription of the adversarial input, rather than on the Assistant’s actual response. Thus non-exact transcriptions that triggered the target action would not have been identified as successes in the retests of previously unsuccessful adversarial input. The reason for this was that whereas the number of successful adversarial wake phrases and adversarial commands to be retested was small enough for the Assistant’s response to them to be tested manually, the number of unsuccessful adversarial wake phrases and adversarial commands to be retested was too large for manual testing. It was therefore necessary to rely on automated identification of exact transcription of target commands in order to establish whether the previously unsuccessful adversarial input had become successful.

Results - Nonsense Attacks - Pilot Experiment

Audio File Input Tests The audio file input tests for the wake phrase ‘Hey Google’ identified 15 successful adversarial wake phrases that triggered activation of the device in around 200 tests. The audio file input tests for target commands identified three successful adversarial commands for five of the target commands, i.e. 15 successful adversarial commands in total, in around 2000 tests.¹⁵ No successful adversarial commands were iden-

¹⁵As stated above, testing involved random sampling from the total space of potential adversarial input for the wake phrase and for each target command. The number of tests performed represented approximately three percent of the total space of potential adversarial input for the wake phrase, and one percent of the total space of potential adversarial input for the target commands. The figures given here are approximate, on account of possible duplicate testing of some adversarial input in the random sampling process used in the pilot experiment.

tified for one of the target commands, namely ‘who am I’. Three samples of the successful adversarial commands for the wake phrase and one sample of an adversarial command for each of the five target commands for which successful adversarial commands were identified in the audio file input tests are shown in Table 5.2. These samples were subsequently used in human comprehensibility tests. Also shown below, in Figure 5.3 and Figure 5.4 respectively, are examples of the print-out to terminal of the Google Assistant Service’s response to a successful adversarial wake phrase and to a successful adversarial command. The full details of successful audio file input tests from the pilot experiment are shown in Appendix A.2.1. Results from the audio file input tests were reproducible in immediate retests based on a few samples.

With regard to the target system’s response to unsuccessful adversarial wake phrases and commands, a certain proportion of the nonsensical word sequences tested in the experiments were transcribed as meaningful word sequences other than the target command, prompting the Assistant to run web searches. For other nonsensical word sequences, the Assistant’s response was simply to indicate non-comprehension of the input.

The results include some examples of successful adversarial commands of which the transcription by the Assistant did not need to match the target command exactly in order to trigger the target action; for example, an adversarial command for the target command ‘turn on light’ is transcribed as ‘switch on the light’ in one instance (see Table 5.2). In some cases, the transcription of an adversarial command did not even need to be semantically equivalent to the target command in order to trigger the target action, such as in the transcription of an adversarial command for “turn light red” as “turn right to Red”. This result is indicative of a vulnerability in the natural language understanding functionality of the Assistant as well as in its speech recognition functionality.

Figure 5.3: Transcription of response to adversarial command for ‘Hey Google’ (“zhay dooble”) from audio file

```
Wakeup word triggered by nonsense_wakeup/Z'eI d'u:b@L.raw, nonsense_wakeup/Z'eI d'u:b@L
INFO:root:Connecting to embeddedassistant.googleapis.com

INFO:root:Recording audio request.
INFO:root:Transcript of user request: "change".
INFO:root:Transcript of user request: "JD".
INFO:root:Transcript of user request: "hey dude".
INFO:root:Transcript of user request: "hey Google".
INFO:root:Transcript of user request: "hey Google".
INFO:root:Transcript of user request: "hey Google".
INFO:root:End of audio request detected.
INFO:root:Stopping recording.
INFO:root:Transcript of user request: "hey Google".
INFO:root:Expecting follow-on query from user.
INFO:root:Playing assistant response.
```

Over-the-Air Tests The over-the-air tests identified four successful full adversarial commands in the 225 tests (representing all possible combinations of each of the 15 successful adversarial wake phrases with each of the three successful adversarial commands for each of five of the target commands). One of the successful over-the-air adversarial commands triggered the ‘turn on light’ target action and three of the successful over-the-air adversarial commands triggered the ‘turn light red’ target action. The four successful over-the-air

Target Command	Adversarial Command (Kirshenbaum phonetic symbols)	Text Transcribed	Action Triggered
Hey Google	S'eI j'u:b@L ("shay yooble")	hey Google	<i>assistant activated</i>
Hey Google	t'eI g'u:t@L ("tay gootle")	hey Google	<i>assistant activated</i>
Hey Google	Z'eI d'u:b@L ("zhay dooble")	hey Google	<i>assistant activated</i>
turn off light	h'3:n z'Of j'alt ("hurn zof yight")	turns off the light	Turning device off
turn light blue	h'3:n gl'alt skw'u: ("hurn glight squoo")	turn the lights blue	color is blue
turn light red	str'3:n j'alt str'Ed ("strurn yight stred")	turn the lights to Red	color is red
what's my name	sm'Ots k'aI sp'eIm ("smots kai spaim")	what's my name	You told me your name was MK
turn on light	p'3:n h'on kl'alt ("purn hon klight")	turn on light	Turning device on

Table 5.2: Samples of successful adversarial commands from audio file input tests used in human comprehensibility tests

Figure 5.4: Transcription of response to adversarial command for ‘what’s my name’ (sm’Ots k’aI sp’eIm - “smots kai spaim”) from audio file

```

INFO:root:Recording audio request.
INFO:root:Transcript of user request: "what's".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "some".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "summer".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "what's on Sky".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "what's my IP".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "some months cause pain".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "what's my car's paint".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "what's my car's paint".
INFO:root:Playing assistant response.
INFO:root:End of audio request detected
INFO:root:Transcript of user request: "what's my name".
INFO:root:Playing assistant response.
INFO:root:You told me your name was MK
I could never forget that
INFO:root:Finished playing assistant response.

```

adversarial commands are shown in Table 5.3. The terminal output by Google Assistant Library for each of the successful full adversarial commands is given in Appendix A.2.2. Figure 5.5 shows an example of the print-out to terminal of the Google Assistant Library’s response to a successful over-the-air adversarial command. The over-the-air results were not shown to be reliably reproducible in immediate retests, i.e. a successful adversarial command did not always the trigger the target action again on re-playing based on a few

attempts for each of the successful commands. Possible reasons for this include quality of the microphone used to pick up voice commands, exact distance between the speaker and the microphone, and the potential presence of background noise.

Target Command	Adversarial Command (x-sampa phonetic symbols)	Text Transcribed	Action Triggered
Hey Google turn on light	t'eI D'u:bl= s'3:n Z'Qn j'aIt ("tay dooble surn zhon yight")	switch on the light	Turning the LED on
Hey Google turn light red	t'eI D'u:bl= tr'3:n Tr'aIt str'Ed ("tay dooble trurn thright stred")	turn lights to Red	The color is red
Hey Google turn light red	t'eI D'u:bl= pr'3:n j'aIt sw'Ed ("tay dooble prurn yight swed")	turn the lights red	The color is red
Hey Google turn light red	t'eI D'u:bl= str'3:n j'aIt str'Ed ("tay dooble strurn yight stred")	turn lights to Red	The color is red

Table 5.3: Successful full adversarial commands in over-the-air tests

Figure 5.5: Transcription of response to adversarial command for 'Hey Google turn on light' (t'eI D'u:bl= s'3:n Z'Qn j'aIt) from over-the-air audio

```
ON_CONVERSATION_TURN_STARTED
ON_END_OF_UTTERANCE
ON_RECOGNIZING_SPEECH_FINISHED:
{"text": "switch on the light"}

Do command action.devices.commands.OnOff with params {u'on': True}
Turning the LED on.
ON_RESPONDING_STARTED:
{"is_error_response": false}
ON_RESPONDING_FINISHED
ON_CONVERSATION_TURN_FINISHED:
{"with_follow_on_turn": false}
```

Human Comprehensibility Tests A total of 19 sets of valid results could be retrieved from the 20 participants in the experiment (there was some data loss due to simultaneous submission of experimental results by multiple participants). All 19 participants who generated these results transcribed the attention tests correctly as 'hi how are you' and 'hello how are you'. Three participants transcribed one adversarial command as the target command 'turn on light', but did not identify the wake phrase 'Hey Google' or any of the other target commands in either the audio file input clips or the over-the-air clips. None of the other participants identified any of the wake phrase or target command words in any of the clips. Eight of the participants identified no meaning at all in any of the clips apart from the attention tests. The other participants either also identified no meaning in the adversarial commands, or else transcribed them as words that were unrelated to the target command. Some examples of unrelated transcriptions were 'hands off the yacht' and 'smoking cause pain'. One participant also transcribed some of the adversarial commands as nonsense syllables e.g. 'hurn glights grew' for the adversarial command for 'turn light blue' and 'pern

pon clight’ for the adversarial command for ‘turn on light’. Another participant transcribed one of the adversarial commands for ‘Hey Google’ as the French words ‘Je du blanc’. The full set of results from the human comprehensibility tests from the pilot experiment is shown in Appendix A.2.3.

Retests The main finding from the retests half a year after the original tests was that whereas some of the previously successful adversarial input was no longer effective, some of the previously unsuccessful adversarial input had conversely become effective, as detailed in Appendix A.2.4. Eight out of the 15 previously successful adversarial wake phrases were found to have become unsuccessful, and eight of 15 previously successful adversarial commands were similarly found to have become unsuccessful. None of the previously unsuccessful wake phrases had become successful, whereas 40 of the previously unsuccessful adversarial commands were found to have become successful. Specifically, 24 of the previously unsuccessful adversarial commands for the ‘who am I’ target command were found to have become successful, as well as two of the previously unsuccessful adversarial commands for the ‘what’s my name’ target command, four of the previously unsuccessful adversarial commands for the ‘turn on light’ target command, seven of the previously unsuccessful adversarial commands for the ‘turn off light’ target command, and three of the previously unsuccessful adversarial commands for the ‘turn light blue’ target command. None of the previously unsuccessful adversarial commands for the ‘turn light red’ target command was found to have become successful.

5.2.3 Main Experiment

Methodology - Nonsense Attacks - Main Experiment

The main experiment followed the same three key stages as the pilot experiment, and also used the same target system and the same target commands. There were some refinements in the details of the methodology used at the different stages as described below.

Adversarial Command Generation The main experiment used the same sets of rhyming nonsense words for each original word in the wake phrase and target commands for adversarial input generation as were used in the pilot experiment. However, there were a number of differences in the generation of adversarial input using the rhyming nonsense words. The most significant difference was that, whereas in the pilot experiment only adversarial wake phrases and commands in which all original words had been replaced with nonsense words were generated, in the main experiment, adversarial commands were also generated in which not all of the original words had been replaced. Adversarial input in which only some original words have been replaced is referred to below as partially mangled adversarial input, whereas adversarial input in which all of the original words have been replaced is referred to as fully mangled adversarial input. ‘Higher’ and ‘lower’ mangling levels are stages of mangling at which more or fewer original words have been replaced respectively.

The inclusion of partially mangled adversarial input in the main experiment was based on two related considerations. Firstly, partially mangled adversarial input retaining some

original words was more likely to be successful in triggering a target action than fully mangled adversarial input, albeit that such input would also be more likely to be detected by human listeners. As the success rate in the pilot experiment using only fully mangled adversarial input had been quite low, the main experiment tested whether it might be possible to execute attacks that evaded human comprehensibility using only partially mangled adversarial input. Secondly, it was considered that successful adversarial input at lower levels of mangling could be used to increase the success rate of adversarial input at higher levels of mangling by constraining the sampling space for generating adversarial input at higher mangling levels. Whereas in the pilot experiment, fully mangled adversarial input was generated by random sampling from the entire space of potential adversarial input, the main experiment used a filtering approach in which adversarial input at higher levels of mangling was generated from successful adversarial input at lower levels of mangling. The aim of this approach was to identify adversarial input that would be most likely both to trigger a target action and to evade human comprehensibility.

For the audio file input tests, adversarial wake phrases and commands in which only one of the original words had been replaced were initially generated. Adversarial wake phrases and commands that were successful at this first level of mangling were then combined with one another to generate fully mangled wake phrases and adversarial commands in which two of the three original words had been mangled. Successful adversarial commands at the second level of mangling were then also combined with successful adversarial commands from the first level to create fully mangled adversarial commands. Examples of the combinatorial filtering process used in generating adversarial wake phrases and commands for the audio file input tests are shown in Figure 5.6.

Figure 5.6: Adversarial input generation process for audio file input tests

Combination of two successful adversarial wake phrases at first mangling level to form fully mangled adversarial wake phrase:

```
v'eI g'u:g@L (`\`vay Google") + h'eI g'Ud@L (`\`Hey goodle") =
v'eI g'Ud@L (`\`vay goodle")
```

Combination of successful adversarial command at first mangling level and successful adversarial command at second mangling level to form fully mangled adversarial command:

```
t'3:n '0 n g'aIt (`\`turn on gight") + S'3:n Tr'0n l'aIt
(`\`shurn thron light") = S'3:n Tr'0n g'aIt (`\`shurn thron
gight")
```

Adversarial wake phrases and adversarial commands that had been successful in audio file input tests were used to generate full adversarial commands for over-the-air tests and human comprehensibility tests. Random samples of both fully and partially mangled adversarial wake phrases and adversarial commands that had been successful in the audio file input tests were used. Whereas the over-the-air tests in the pilot experiment tested all possible combinations of successful adversarial wake phrases and successful adversarial

commands from the audio file input tests, the over-the-air tests in the main experiment tested only a sample of these, due to a much larger number of successful adversarial wake phrases and adversarial commands being identified in the audio file input tests in the main experiment as compared to the pilot experiment. Also in contrast to the pilot experiment, the over-the-air and human comprehensibility tests in the main experiment used the same set of adversarial input. Full adversarial commands were generated for all but one of the target commands at different mangling levels (for one of the target commands, there was an insufficient number of successful adversarial commands in audio file input tests to create a complete set of full adversarial commands for the over-the-air and human comprehensibility tests). In addition to fully mangled adversarial commands, full adversarial commands were also generated at the following three partial mangling levels:

- full adversarial commands in which four of the original words had been mangled (two in the wake phrase and two in the specific target command, or one in the wake phrase and three in the specific target command)
- full adversarial commands in which three of the original words had been mangled (one in the wake phrase and two in the specific target command, or two in the wake phrase and one in the specific target command)
- full adversarial commands in which two of the original words had been mangled (one in the wake phrase and one in the specific target command)

At each partial mangling level, two full adversarial commands were generated for each target command, one in which the word ‘Google’ was one of mangled words, and one in which ‘Google’ was not mangled. This was in order to test the effect of the presence of the unmangled word ‘Google’ on machine and human comprehensibility of partially mangled full adversarial commands.

In addition to the differences in the composition of adversarial input, there were also two minor differences in the synthesis of audio versions of adversarial input in the main experiment. The first of these was that whereas in the pilot experiment, audio of the adversarial input for the audio file input tests was created using espeak speech synthesis, in the main experiment the Amazon Polly speech synthesis system was used to create the adversarial audio for both the audio file input and the over-the-air tests. The second difference in speech synthesis was that whereas in the pilot experiment, audio versions of adversarial wake phrases and commands had been synthesised as one string, in the main experiment audio versions of adversarial commands were synthesised by concatenating pre-synthesised audio of individual target command words and nonsense words, leading to a ‘staccato’ effect in the audio produced. This was done in order to increase the efficiency of the speech synthesis process.

Machine Comprehensibility Tests As in the pilot experiment, the audio file input tests in the main experiment involved testing the target system’s response to adversarial input via terminal input using the Google Assistant Service and via microphone input using the Google Assistant Library. There was one important difference to the methodology followed in machine comprehensibility tests in the main experiment to the methodology followed in

the pilot experiment. This was that whereas the pilot experiment had established triggering of a target action based on the Assistant's actual response, in the main experiment, a target action was considered to have been triggered if the Assistant's final transcription of adversarial input matched the target command. The reason for this adjustment was that it ensured that adversarial input would be identified as successful only if it exploited vulnerability in speech recognition. By contrast, in the pilot experiment the triggering of some target actions relied on exploitation of vulnerability in natural language understanding as well as speech recognition in some instances, as detailed above. The adjustment to the methodology with respect to establishing the triggering of a target action also potentially expands the range of target commands beyond actions that are actually within the scope of the target system's capabilities, to actions that have not as yet been implemented, although in this case, the target commands used in the main experiment were the same as in the pilot experiment.

As stated above, whereas the audio file tests in the pilot experiment involved only adversarial wake phrases and commands in which all of the original words had been mangled (i.e. replaced by a rhyming nonsense word), the audio file inputs tests in the main experiment also tested the effectiveness of adversarial wake phrases and commands in which only some words had been mangled. Furthermore, the audio file input tests in the main experiment involved a filtering process in which testing of adversarial input in which more than one of the original words had been mangled was limited to combinations of adversarial wake phrases or adversarial commands that had been successful at the previous level of mangling. The testing process was automated using a Python script. The script first tested adversarial wake phrases and commands in which only one of the original words had been mangled. Adversarial wake phrases and commands that were successful at this first level were then combined with one another to create a second level of adversarial wake phrases and commands in which two words had been mangled. Adversarial wake phrases and commands that were not successful at the first level were discarded. In the case of adversarial commands, a third level was also tested consisting of combinations of successful adversarial commands from the first and second levels of mangling. At the first level of mangling, in which only one of the original words was replaced with a nonsense word, all adversarial wake phrases and adversarial commands were tested. At subsequent mangling levels, the Python script tested up to a maximum of 150 adversarial commands at each level using random sampling, with a target of maximum 100 successes. As in the pilot experiment, a random sampling process was followed due to the large space of adversarial commands.¹⁶

For the over-the-air tests, as stated above, a random sample of adversarial wake phrases and adversarial commands that had been successful in audio file input tests were con-

¹⁶The testing covered 100 percent of the space of potential adversarial input at the first level of mangling, and up to 100 percent of the space of potential adversarial input at higher levels of mangling. The space of adversarial input at the higher mangling levels consisted of all possible combinations of successful adversarial input from the previous mangling level. Random sampling from this space of adversarial input was performed up to a maximum of 150 tests or 100 successful tests. There was some discrepancy between the number of testing attempts reported by the Python script for the experiment, and the number of tests actually performed. The reason for this was that if the same adversarial input was selected more than once in random sampling, the script recorded a testing attempt, but did not actually perform a duplicate test of the input. Results of testing were calculated according to the actual number of unique tests performed.

catenated to form full adversarial commands. The effectiveness of these full adversarial commands was tested via microphone input over a set of speakers using Google Assistant Library. Table 5.4 shows the details of the full adversarial commands of which over-the-air effectiveness was tested.

Table 5.4: Full adversarial commands for over-the-air and human comprehensibility tests

Target Action with Condition	Fully mangled command	Third mangling level	Second mangling level	First mangling level	Plain-speech command
Retrieve username (Google unmangled first)	Z'eI l'Uk@L spl'u: bl'am str'al ("zhay lookle sploo blam strai")	v'eI g'u:g@L spl'u: bl'am str'al ("vay Google sploo blam strai")	v'eI g'u:g@L v'u: T'am 'al ("vay Google voo tham I")	v'eI g'u:g@L h'u: T'am 'al ("vay Google who tham I")	Hey Google who am I
Retrieve username (Google unmangled last)	[n.a.]	Z'eI l'Uk@L v'u: T'am 'al ("zhay lookle voo tham I")	Z'eI l'Uk@L h'u: T'am 'al ("zhay lookle who tham I")	h'eI g'Ud@L h'u: T'am 'al ("Hey goodle who tham I")	[n.a.]
Retrieve username (Google unmangled first)	T'eI gl'u:s@L D'0ts sn'al z'eIm ("thay gloosle thots snai zame")	Z'eI g'u:g@L D'0ts sn'al z'eIm ("zhay Google thots snai zame")	Z'eI g'u:g@L w'0ts gr'al Z'eIm ("zhay Google what's grai zhame")	Z'eI g'u:g@L w'0ts bl'al n'eIm ("zhay Google what's blai name")	Hey Google what's my name
Retrieve username (Google unmangled last)	[n.a.]	h'eI w'u:b@L D'0ts sn'al z'eIm ("Hey wooble thots snai zame")	h'eI w'u:b@L w'0ts gr'al Z'eIm ("Hey wooble what's grai zhame")	h'eI w'u:b@L w'0ts bl'al n'eIm ("Hey wooble what's blai name")	[n.a.]
Turn off smart light (Google unmangled first)	v'eI g'u:t@L g'3:n bl'of j'alt ("vay gootle gurn blof yight")	v'eI g'u:g@L g'3:n bl'of j'alt ("vay Google gurn blof yight")	v'eI g'u:g@L pr'3:n b'of l'alt ("vay Google prurn bof light")	v'eI g'u:g@L tr'3:n 'of l'alt ("vay Google trurn off light")	Hey Google turn off light
Turn off smart light (Google unmangled last)	[n.a.]	h'eI k'u:z@L g'3:n bl'of j'alt ("Hey koozle gurn blof yight")	v'eI g'u:t@L tr'3:n 'of l'alt ("vay gootle trurn off light")	h'eI k'u:z@L tr'3:n 'of l'alt ("Hey koozle trurn off light")	[n.a.]
Turn on smart light (Google unmangled first)	Z'eI fl'Uk@L D'3:n f'on D'alt ("zhay flookle thurn fon thight")	Z'eI g'u:g@L D'3:n f'on D'alt ("zhay Google thurn fon thight")	Z'eI g'u:g@L t'3:n tr'on p'alt ("zhay Google thurn tron pight")	Z'eI g'u:g@L br'3:n 'on l'alt ("zhay Google brurn on light")	Hey Google turn on light
Turn on smart light (Google unmangled last)	[n.a.]	h'eI k'u:s@L D'3:n f'on D'alt ("Hey koosle thurn fon thight")	Z'eI fl'Uk@L br'3:n 'on l'alt ("zhay flookle brurn on light")	h'eI k'u:s@L br'3:n 'on l'alt ("Hey koosle brurn on light")	[n.a.]
Change smart light colour (Google unmangled first)	Z'eI gl'u:p@L pl'3:n g'alt v'u: ("zhay gloople plurn gight voo")	T'eI g'u:g@L pl'3:n g'alt v'u: ("thay Google plurn gight voo")	T'eI g'u:g@L fl'3:n v'alt bl'u: ("thay Google flurn vight blue")	T'eI g'u:g@L t'3:n Z'alt bl'u: ("thay Google turn zhight blue")	Hey Google turn light blue
Change smart light colour (Google unmangled last)	[n.a.]	Z'eI gl'u:p@L fl'3:n v'alt bl'u: ("zhay gloople flurn vight blue")	h'eI bl'Uk@L fl'3:n v'alt bl'u: ("Hey blookle flurn vight blue")	h'eI bl'Uk@L t'3:n Z'alt bl'u: ("Hey blookle turn zhight blue")	[n.a.]

Human Comprehensibility Tests The human comprehensibility tests used the same set of full adversarial commands as the over-the-air machine comprehensibility tests, listed in Table 5.4. Participants in the human comprehensibility tests were presented with full

adversarial commands in descending order of mangling on a spectrum from fully mangled adversarial commands to adversarial commands in which only one word in the wake phrase and one word in the target command had been mangled. This approach was taken so as to provide an indication of how many words would need to be mangled in an adversarial over-the-air command in order to escape human comprehensibility. After hearing all of adversarial commands, participants were also presented with a plain-speech version of the full target command. This provided a baseline for the comprehensibility tests, and also served as an attention test. The partially mangled full adversarial commands were separated into two sets for each target command. In the first set, “Google” was the first word to be revealed to the listener in plain-speech, whereas in the second set, “Google” was the last word to be revealed. The separation of these two conditions enabled an assessment of whether the presence of the specific word “Google” affected listeners’ ability to detect a malicious voice command. Each set of full adversarial commands for each target command under each of the two conditions was played to six different participants using the survey website Prolific Academic. The experiments with human subjects received ethics clearance through the Departmental Research Ethics Committee of the Department of Computer Science at the University of Oxford (Ref No: SSD/CUREC1A CS_C1A_18_021).

Results - Nonsense Attacks - Main Experiment

Audio File Input Tests The results of the audio file input tests in the main experiment showed significantly higher success rates for fully mangled adversarial wake phrases and commands by comparison to the pilot experiment. Whilst in pilot experiment, the success rate for adversarial input had been only around one percent, in the main experiment the success rate for fully mangled adversarial input lay between 12.3 percent for the ‘who am I’ target command and 63.6 percent for the ‘turn light blue’ target command, as shown in Table 5.5. Notwithstanding a potential effect of differences in speech synthesis used in the main experiment, this indicated that limiting testing of adversarial commands at each mangling level to combinations of adversarial commands that had been successful at the previous level significantly increased the effectiveness of attacks with fully mangled adversarial input. Table 5.5 shows the overall ratio of successes to failures in the audio file input tests, as well as the number of successes for adversarial wake phrases and commands at the first level, covering all possible adversarial commands in which only one word was mangled, and at subsequent levels, where a random selection of possible combinations of successful adversarial wake phrases or commands from the previous level or levels were tested. With the exception of the “turn light red” target command, successful adversarial commands could be generated for all target commands at all mangling levels, and this was also the case for the wake phrase. Overall success rates for target commands apart from the “turn light red” command ranged from 29.9 percent to 63.8 percent. The “turn light red” target command appeared to be an outlier in terms of success rates for adversarial commands, with a success rate of only 3.2 percent. No clear reason for this was apparent. The overall success rate for the adversarial wake phrase was 14.4 percent. The full lists of successful adversarial wake phrases and adversarial commands at each mangling level are shown in Appendix A.3.1.

Target Command	Overall Success Rate (%)	Level 1 Successes (%)	Level 2 Successes (%)	Level 3 Successes (%)
Hey Google	14.4 (70 of 487)	12.6 (52 of 412)	24.0 (18 of 75)	n.a.
Who am I	29.9 (85 of 284)	70.7 (46 of 65)	28.8 (21 of 73)	12.3 (18 of 146)
What's my name	55.4 (165 of 298)	68.3 (56 of 82)	69.3 (52 of 75)	40.4 (57 of 141)
Turn on light	49.2 (155 of 315)	41.5 (44 of 106)	61.3 (46 of 75)	48.5 (65 of 134)
Turn off light	56.7 (185 of 326)	47.7 (52 of 109)	66.6 (50 of 75)	58.5 (83 of 142)
Turn light red	3.2 (3 of 93)	3.2 (3 of 93)	0	0
Turn light blue	63.8 (166 of 260)	42.7 (41 of 86)	82.6 (62 of 75)	63.6 (63 of 99)

Table 5.5: Success rates of adversarial commands in audio file input tests

Figure 5.7 shows examples of the output to terminal produced by a successful fully mangled adversarial wake phrase and by a successful fully mangled adversarial command for each of the target commands. The output shows both interim and final transcriptions of the adversarial input by the Assistant. Figure 5.8 shows examples of the output to terminal produced by an unsuccessful fully mangled adversarial wake phrase and by an unsuccessful fully mangled adversarial command for each of the target commands. The unsuccessful examples share one nonsensical word sound with the corresponding successful example in Figure 5.7, demonstrating that the success or failure of adversarial wake phrases and commands in triggering a target action was influenced not only by the probabilities allocated to individual word sounds by the acoustic model used in the Assistant's speech recognition, but also by the probabilities allocated to utterances as a whole by the Assistant's language model. For example, the rhyming word 'thight' for the word 'light' is found both in a successful and in an unsuccessful adversarial command for the 'turn light blue' target command. The successful adversarial command 'grurn thight voo' is transcribed as the target command 'turn light blue', whereas the unsuccessful adversarial command 'skurn thight voo' is transcribed as 'screen side view'.

Figure 5.7: Audio File Input Tests - Successes

```
ADVERSARIAL WAKE PHRASE FOR "Hey Google": v'eI g'u:t@L ("vay gootle")

WARNING:root:Transcript of user request: "V".
WARNING:root:Transcript of user request: "wake".
WARNING:root:Transcript of user request: "Virgo".
WARNING:root:Transcript of user request: "very good".
WARNING:root:Transcript of user request: "viagogo".
WARNING:root:Transcript of user request: "hey Google".
WARNING:root:Transcript of user request: "hey Google".
WARNING:root:Transcript of user request: "hey Google".
WARNING:root:Playing assistant response.
WARNING:root:Expecting follow-on query from user.
WARNING:root:Finished playing assistant response.
RESPONSE TRANSCRIPTION: hi what can I do for you

ADVERSARIAL COMMAND FOR "what's my name": b'0ts j'aI w'eIm ("bots yai wame")

WARNING:root:Transcript of user request: "buy".
WARNING:root:Transcript of user request: "file".
WARNING:root:Transcript of user request: "lights".
WARNING:root:Transcript of user request: "what's the".
WARNING:root:Transcript of user request: "what's your".
WARNING:root:Transcript of user request: "what's my".
WARNING:root:Transcript of user request: "what's my".
WARNING:root:Transcript of user request: "what's my way".
WARNING:root:Transcript of user request: "what's my way".
WARNING:root:Transcript of user request: "what's my name".
WARNING:root:Transcript of user request: "what's my name".
WARNING:root:Transcript of user request: "what's my name".
WARNING:root:Playing assistant response.
WARNING:root:Finished playing assistant response.

ADVERSARIAL COMMAND FOR "who am I": f'u: D'am z'aI ("foo tham zai")

WARNING:root:Transcript of user request: "true".
WARNING:root:Transcript of user request: "through the".
WARNING:root:Transcript of user request: "who am".
WARNING:root:Transcript of user request: "fu Fareham".
WARNING:root:Transcript of user request: "who am I".
WARNING:root:Transcript of user request: "who am I".
WARNING:root:Transcript of user request: "who am I".
WARNING:root:Playing assistant response.
WARNING:root:Finished playing assistant response.

ADVERSARIAL COMMAND FOR "turn on light": z'3:n bl'0n d'aIt ("zurn blon dight")

WARNING:root:Transcript of user request: "is there".
WARNING:root:Transcript of user request: "is there".
WARNING:root:Transcript of user request: "is there a".
WARNING:root:Transcript of user request: "turn on".
WARNING:root:Transcript of user request: "turn brown die".
WARNING:root:Transcript of user request: "turn on light".
WARNING:root:Transcript of user request: "turn on light".
WARNING:root:Transcript of user request: "turn on light".
WARNING:root:Playing assistant response.
WARNING:root:Finished playing assistant response.

ADVERSARIAL COMMAND FOR "turn off light": n'3:n T'0f j'aIt ("nurn thoff yight")

WARNING:root:Transcript of user request: "no".
WARNING:root:Transcript of user request: "9".
WARNING:root:Transcript of user request: "turn off".
WARNING:root:Transcript of user request: "turn off the".
WARNING:root:Transcript of user request: "turn off the".
WARNING:root:Transcript of user request: "turn off my".
WARNING:root:Transcript of user request: "turn off light".
WARNING:root:Transcript of user request: "turn off light".
WARNING:root:Transcript of user request: "turn off light".
WARNING:root:Playing assistant response.
WARNING:root:Finished playing assistant response.

ADVERSARIAL COMMAND FOR "turn light blue": gr'3:n D'aIt v'u: ("grurn thight voo")

WARNING:root:Transcript of user request: "turn".
WARNING:root:Transcript of user request: "current".
WARNING:root:Transcript of user request: "turn my".
WARNING:root:Transcript of user request: "turn light".
WARNING:root:Transcript of user request: "turn light".
WARNING:root:Transcript of user request: "turn light you".
WARNING:root:Transcript of user request: "turn light you".
WARNING:root:Transcript of user request: "turn light you".
WARNING:root:Transcript of user request: "turn light blue".
WARNING:root:Playing assistant response.
WARNING:root:command triggered: @@@LIGHT TURNED BLUE@@@
WARNING:root:Finished playing assistant response.
```

Figure 5.8: Audio File Input Tests - Losses

```
ADVERSARIAL WAKE PHRASE FOR "Hey Google": v'eI gl'u:f@L ("vay gloofle")

WARNING:root:Transcript of user request: "v".
WARNING:root:Transcript of user request: "wake".
WARNING:root:Transcript of user request: "vehicle".
WARNING:root:Transcript of user request: "fake love".
WARNING:root:Transcript of user request: "The Gruffalo".
WARNING:root:Transcript of user request: "The Gruffalo".
WARNING:root:Transcript of user request: "The Gruffalo".
WARNING:root:Playing assistant response.
WARNING:root:Finished playing assistant response.

ADVERSARIAL COMMAND FOR "what's my name": Z'0ts j'aI v'eIm ("thots yai vame")

WARNING:root:Transcript of user request: "je".
WARNING:root:Transcript of user request: "shut".
WARNING:root:Transcript of user request: "charts".
WARNING:root:Transcript of user request: "start a".
WARNING:root:Transcript of user request: "shut your".
WARNING:root:Transcript of user request: "Shout by".
WARNING:root:Transcript of user request: "shut ya".
WARNING:root:Transcript of user request: "shout by the".
WARNING:root:Transcript of user request: "shut your name".
WARNING:root:Transcript of user request: "shut your name".
WARNING:root:Transcript of user request: "shout your fame".
WARNING:root:Playing assistant response.
WARNING:root:Finished playing assistant response.

ADVERSARIAL COMMAND FOR "who am I": spl'u: bl'am z'aI ("sploo blam zai")

WARNING:root:Transcript of user request: "screw".
WARNING:root:Transcript of user request: "play".
WARNING:root:Transcript of user request: "volume".
WARNING:root:Transcript of user request: "who do I am sorry".
WARNING:root:Transcript of user request: "who do I am sorry".
WARNING:root:Transcript of user request: "volume three".
WARNING:root:Finished playing assistant response.

ADVERSARIAL COMMAND FOR "turn on light": z'3:n tr'0n p'aIt ("zurn tron pight")

WARNING:root:Transcript of user request: "is there".
WARNING:root:Transcript of user request: "turn".
WARNING:root:Transcript of user request: "try".
WARNING:root:Transcript of user request: "turn on".
WARNING:root:Transcript of user request: "Tron I".
WARNING:root:Transcript of user request: "turn on light".
WARNING:root:Transcript of user request: "turn on light".
WARNING:root:Transcript of user request: "Tron ID".
WARNING:root:Playing assistant response.
WARNING:root:Finished playing assistant response.

ADVERSARIAL COMMAND FOR "turn off light": n'3:n v'0f tS'aIt ("nurn voff chight")

WARNING:root:Transcript of user request: "no".
WARNING:root:Transcript of user request: "Night by".
WARNING:root:Transcript of user request: "new bar".
WARNING:root:Transcript of user request: "new bath".
WARNING:root:Transcript of user request: "buy a".
WARNING:root:Transcript of user request: "bye bye".
WARNING:root:Transcript of user request: "turn both tried".
WARNING:root:Transcript of user request: "9 Bath Street".
WARNING:root:Playing assistant response.
WARNING:root:Finished playing assistant response.

ADVERSARIAL COMMAND FOR "turn light blue": sk'3:n D'aIt v'u: ("skurn thight voo")

WARNING:root:Transcript of user request: "Sky".
WARNING:root:Transcript of user request: "turn".
WARNING:root:Transcript of user request: "turn the".
WARNING:root:Transcript of user request: "turn the".
WARNING:root:Transcript of user request: "School Guide".
WARNING:root:Transcript of user request: "School Guide you".
WARNING:root:Transcript of user request: "School Guide you".
WARNING:root:Transcript of user request: "School Guide you".
WARNING:root:Transcript of user request: "screen side view".
WARNING:root:Playing assistant response.
WARNING:root:Finished playing assistant response.
```

Over-the-Air Tests The results of the over-the-air tests, testing the Assistant’s response to over-the-air input of the partially mangled and fully mangled adversarial commands in Table 5.4, are shown in Table 5.6. The results show four successes in activating the full target command. All of the successful adversarial commands were partially rather than fully mangled commands. Specifically, the complete target action was activated by the adversarial commands for the ‘what’s my name’ target command at the third, second and first levels of mangling under the condition of the word Google being revealed last, and by the adversarial command for ‘turn on light’ at the first level of mangling under the condition of the word Google being revealed last. In addition to the adversarial commands that were successful in triggering the target action, there were also two instances in which the adversarial input activated the wake phrase only. Specifically, the wake phrase was activated by the fully mangled adversarial command for ‘who am I’, as well as by the adversarial command for the same target command at the third level of mangling under the condition of the word ‘Google’ being revealed last. Details of the print-out to terminal following successful activation of the target command or wake phrase are given in Appendix A.3.2. As in the audio file input tests, the success rate of adversarial commands was higher in the main experiment than in the pilot experiment, at a rate of 10 percent (4 out of 40 tests) as compared to 1.8 percent (4 out 225 tests) in the pilot experiment.

Table 5.6: Results of Over-the-Air Tests

Condition	fully mangled command	Third mangling level	Second mangling level	First mangling level	Target command
Google unmangled first	unsuccessful (<i>wake phrase activated</i>)	unsuccessful	unsuccessful	unsuccessful	Hey Google who am I
Google unmangled last	[n.a.]	unsuccessful (<i>wake phrase activated</i>)	unsuccessful	unsuccessful	[n.a.]
Google unmangled first	unsuccessful	unsuccessful	unsuccessful	unsuccessful	Hey Google what’s my name
Google unmangled last	[n.a.]	successful	successful	successful	[n.a.]
Google unmangled first	unsuccessful	unsuccessful	unsuccessful	unsuccessful	Hey Google turn on light
Google unmangled last	[n.a.]	unsuccessful	unsuccessful	successful	[n.a.]
Google unmangled first	unsuccessful	unsuccessful	unsuccessful	unsuccessful	Hey Google turn off light
Google unmangled last	[n.a.]	unsuccessful	unsuccessful	unsuccessful	[n.a.]
Google unmangled first	unsuccessful	unsuccessful	unsuccessful	unsuccessful	Hey Google turn light blue
Google unmangled last	[n.a.]	unsuccessful	unsuccessful	unsuccessful	[n.a.]

Human Comprehensibility Tests The results of the human comprehensibility tests are summarised in Table 5.7. The results are summarised according to whether a simple majority of participants identified no meaning, part of the target command meaning, or the full target command meaning in the adversarial input. Where there was no simple majority, with equal numbers of participants returning different results, this is indicated in the

table. There were four instances where participants returned a blank test result. In these cases, results are given out of five participants instead of six, as detailed in the table. As regards transcription of the plain-speech target commands serving as attention tests, these were transcribed correctly by a large majority of participants. There were three instances where transcription of the plain-speech was incomplete, one where it was incorrect, and one where transcription of the plain-speech was missing. The full results of the human comprehensibility tests are given in Appendix A.3.3.

A consistent result across all the tests was that, with one sole exception, none of the participants identified any meaning in the fully mangled full adversarial commands. Another consistent result was that participants generally detected more meaning in the adversarial input at higher levels of mangling than at lower levels. A significant specific result was that the partially mangled adversarial commands for ‘what’s my name’ at the third and second levels of mangling that had been effective in triggering the target action in the over-the-air tests were identified by all participants in the human comprehensibility tests as either having no meaning at all, or else as having a meaning only partially related to the target command. Thus these two partially mangled adversarial commands represent fully effective covert attacks on the target system. Otherwise the results showed considerable differences between results for individual participants. Some participants did not hear any meaning in the audio clips prior to hearing the plain-speech command, whereas others picked up some of the adversarial wake phrase or target command words at the previous levels of mangling. Some participants identified words in adversarial commands that were not actually present in the wake phrase or target command. A few participants believed that they had heard a different language, or tried to transcribe some of the nonsensical word sounds. In a few instances, participants identified the entire meaning of a target command prior to hearing the plain-speech version. The condition as to whether the word ‘Google’ was revealed first or last did not appear to significantly affect the participants’ ability to detect the content of the entire command, as results under the two conditions were generally similar at each mangling level.

5.2.4 Discussion

The combined results from the machine response and human comprehensibility tests in both the pilot and the main experiment confirm the hypothesis that voice-controlled digital assistants are vulnerable to covert attacks using nonsensical sounds. The key finding is that voice commands to voice-controlled digital assistant Google Assistant can be triggered by nonsensical word sounds while the same nonsensical word sounds are perceived by humans as either not having any meaning at all, or as having a meaning only partially related to the voice commands to the Assistant. The findings of the main experiment also showed that it is not always necessary to replace all of the original words in a target command in order to generate an adversarial command that is successful in triggering a target action in a target system whilst also evading human comprehensibility. This is based on the finding that partially mangled adversarial commands were successful both in triggering a target action over-the-air and in hiding from human recognition in some instances. A further finding from the main experiment is that the effectiveness of attacks can be increased by applying

Table 5.7: Results of Human Comprehensibility Tests

Target command	Condition	Fully mangled command	Third mangling level	Second mangling level	First mangling level
Hey Google who am I	Google unmangled first	no meaning (5/6 participants)	no meaning (5/6 participants)	no meaning (3/6 participants) partial meaning (3/6 participants)	partial meaning (4/6 participants)
[n.a.]	Google unmangled last	[n.a.]	no meaning (4/6 participants)	no meaning (4/6 participants)	partial meaning (4/6 participants)
Hey Google what's my name	Google unmangled first	no meaning (6/6 participants)	no meaning (4/6 participants)	partial meaning (4/6 participants)	partial meaning (5/6 participants)
[n.a.]	Google unmangled last	[n.a.]	partial meaning (4/6 participants)	no meaning (3/6 participants) partial meaning (3/6 participants)	partial meaning (2/5 participants) full meaning (2/5 participants)
Hey Google turn on light	Google unmangled first	no meaning (6/6 participants)	no meaning (5/6 participants)	no meaning (3/6 participants) partial meaning (3/6 participants)	partial meaning (5/6 participants)
[n.a.]	Google unmangled last	[n.a.]	no meaning (4/6 participants)	partial meaning (4/6 participants)	partial meaning (4/6 participants)
Hey Google turn off light	Google unmangled first	no meaning (6/6 participants)	no meaning (4/6 participants)	no meaning (3/6 participants) partial meaning (3/6 participants)	partial meaning (5/6 participants)
[n.a.]	Google unmangled last	[n.a.]	no meaning (5/6 participants)	no meaning (3/6 participants) partial meaning (3/6 participants)	partial meaning (5/6 participants)
Hey Google turn light blue	Google unmangled first	no meaning (6/6 participants)	no meaning (5/5 participants)	no meaning (3/5 participants)	partial meaning (6/6 participants)
[n.a.]	Google unmangled last	[n.a.]	no meaning (5/6 participants)	partial meaning (5/6 participants)	partial meaning (5/6 participants)

a systematic methodology for using successful adversarial input at lower levels of mangling to constrain the sampling space for adversarial input at higher mangling levels.

The results confirm the influence of the three features of speech recognition in current voice-controlled systems discussed in Section 5.2.1 in enabling this type of attack via the speech interface. These three features were thus shown to represent security vulnerabilities in the current generation of voice-controlled digital assistant.

The first of these features was the target system's inability to recognise the true nature of nonsensical word sounds. As envisaged, the attacks demonstrated in this experimental work exploit a vulnerability in the speech recognition functionality of the Google Assistant target system of being unable to recognise nonsensical word sounds as meaningless. In the results of the experimental work, the Google Assistant target system always either indicated incomprehension or attempted to match the nonsensical sounds to real words, rather than transcribing the nonsense word sound. This confirms that the Assistant is vulnerable to being fooled by word sounds that are perceived by humans as obviously nonsensical. The findings are in accord with the hypothesis behind these experiments that as a grey area between speech and non-speech, nonsensical word sounds represent a part of the input space to a voice-controlled system that the current generation of voice-controlled digital assistants struggle to handle appropriately. Whilst the Assistant does reject some of the

input from this grey area as incomprehensible, in other instances input from this grey area is treated as meaningful input.

The second of these features was the influence of the language model in enabling the success of some of the attacks. The examples found in the main experiment of the same nonsensical word sounds being present in both successful and unsuccessful adversarial inputs confirms that the triggering of a target action by adversarial input may be influenced by probabilities allocated by the language model used in speech recognition to an utterance as a whole, as well as to probabilities allocated to individual word sounds by the acoustic model. Thus the aim of language modelling of ‘correcting’ possibly incorrect word recognitions in legitimate input may have the opposite effect in an adversarial context of enabling the success of attacks based on nonsensical word sounds.

The third feature shown to be exploited in the attacks was discrepancy in human and machine processing of nonsensical input. This was evident in the finding from the human comprehensibility tests that human listeners did not detect the presence of target commands in successful adversarial input in the majority of instances. Whilst machine and human perceptions of specific adversarial commands were different, machine and human responses to the adversarial input were generally comparable to some extent, in that both machine and humans frequently indicated incomprehension of the nonsensical word sounds, or else attempted to fit them to meaningful words. However, a difference between human and machine processing of the adversarial input was that, in addition to either indicating incomprehension or transcribing the nonsensical word sounds as real words, human subjects on occasion attempted to transcribe the nonsensical word sounds phonetically as nonsense syllables. This superior ability of humans to recognise nonsensical word sounds as meaningless may have paradoxically prevented human listeners from detecting the presence of a malicious voice command in successful adversarial input.

One notable finding from the human comprehensibility tests is that some of the human participants heard a foreign language in the nonsensical word sounds. This opens up the possibility of cross-language attacks on voice-controlled systems trained to respond to different languages, whereby humans might hear an utterance in one language that a target system interprets as a target command in another language. Depending on the language understanding of the human users who are the victims of such attacks, cross-language attacks could be termed either nonsense or missense attacks.

Another notable finding with respect to the human comprehensibility tests is the variability in results between individual experiment participants. Thus the covert nature of these attacks depends to some extent on individual human perception, i.e. whereas some individuals may hear target command words in an adversarial command based on nonsensical word sounds, others may not. This was seen in the variable results of the human comprehensibility tests described above. Human perception of word sounds is known to be unstable in some instances, seen for example in a widely shared audio recording in which some listeners heard the word “Yanny” whereas others heard the word “Laurel”.¹⁷

¹⁷See for example The Guardian, “Laurel or Yanny explained: why do some people hear a different word?”, 17th May 2018, <https://www.theguardian.com/technology/2018/may/16/yanny-or-laurel-sound-illusion-sets-off-ear-splitting-arguments>

The retests of results from the pilot experiment after a six-month interval indicate that that specific adversarial commands of this type may not be effective on different systems that have the same functionality, but have been trained with different data, or even on the same system over time as it is updated with new batches of training data from user interactions. Conversely, the retests show that specific adversarial commands of this type that are not successful on one system may be successful on a different system, or on the same system that has been updated with new training data. Thus, whilst changes in training data may affect the effectiveness of individual adversarial commands, they will not make a system invulnerable to all adversarial commands of this type; indeed, training updates may introduce new vulnerabilities to specific adversarial input of this type whilst removing old ones, as was seen in the retests. Therefore training data updates do not represent a ‘patch’ as such, as are applied to fix other types of bugs in software.

5.3 Missense Attacks on Amazon Alexa

This section presents experimental work demonstrating that it is possible to gain unauthorised access to a voice-controlled system using utterances that are accepted by the system as a target command despite having a different meaning to the command. The results of a proof-of-concept study, a pilot experiment and a main experiment are presented. As stated above, in terms of the taxonomy of attacks via the speech interface developed in Chapter 4, the attack presented here is a missense attack. The attack can also be characterised as an adversarial learning attack using out-of-scope input to a system. The attack is a black-box attack. There are parallels between this attack and linguistic steganography, in which a message is hidden in natural language as a cover medium to make it imperceptible to third parties. This is done by making syntactically and semantically acceptable amendments to a text, for example by synonym substitution (see for example Chang and Clark [35]). The missense attacks demonstrated here use comparable methods in amending target commands so as to make them unrecognisable to humans whilst remaining recognisable by a target system. Unlike previous missense attacks targeting speech recognition functionality, such as the attacks on speech transcription demonstrated by Carlini and Wagner [30], the attacks demonstrated here target the natural language understanding functionality of a voice-controlled system. This represents the first demonstration of a missense attack targeting natural language understanding in a voice-controlled system.

5.3.1 Description and Context

The aim of the attacks was to generate adversarial utterances that trigger a target command in a target system, but that are unrecognisable as that command by legitimate users of the system, and are instead understood as having a meaning unrelated to the target command. Adversarial utterances were generated either by replacing a word in a target command with a word of unrelated meaning, or else by embedding alternate meanings of target command words in an utterance of unrelated meaning to the target command. These attack methods are termed ‘word substitution’ and ‘word transplant’ attacks respectively. A word substitution attack is defined as an attack in which one or more words in a target command are

replaced by a different word so as to change the overall meaning of the utterance, whilst retaining as many of the original target command words as possible. A word transplant attack is defined as an attack in which words from a target command are reused in a different sense in a new utterance, seeking to reuse as many words from the target command as possible. In the proof-of-concept study and the pilot experiment, both word substitution and word transplant attacks were used. The main experiment focussed solely on word transplant attacks, based on the finding that word transplant was more effective than word substitution in triggering target commands in the pilot experiment.

No examples of missense attacks targeting the natural language understanding functionality of voice-controlled digital assistants have been demonstrated in prior work. Word substitution attacks on natural language understanding have been demonstrated in some related areas such as sentiment analysis, as for example in the work of Alzantot et al. [9] and Kuleshov et al. [124] discussed in Chapter 3. Word transplant attacks have not been demonstrated in any prior work. A possible reason for the lack of prior work on attacks targeting natural language understanding in voice-controlled systems may be that linguistically plausible adversarial examples that trigger an action in a voice-controlled system with an utterance of apparently unrelated meaning are difficult to generate using automated, mathematical approaches. As also noted by Papernot et al. [180], adversarial learning in the context of natural language understanding technologies that take as input a sequence of words is not a differentiable problem. Papernot et al. concede that their own work on fooling a sentiment classifier with ‘nonsensical’ sentences generated using a mathematical method has some limitations, in that the nonsensical nature of the adversarial sentences is easily noticeable by humans. They point to the need in future work to address grammar and semantics in adversarial sentence generation, in order to make sentences indistinguishable from innocent utterances by humans. The attacks demonstrated here use a manual, non-mathematical approach for generating adversarial voice commands by manipulating linguistic parameters such as syntactic structures and word meaning, rather than mathematical parameters such as acoustic feature or word embedding vector values.

Word transplant attacks for the main experiment were crowd-sourced from human experiment participants. The experimental methodology was inspired in part by gamification of adversarial input generation in image recognition. In a game named ‘Bad Flamingo’, participants are required to make sketches of a given object, aiming to make the object recognisable by other human game participants, but not by a DNN-based image recognition system.¹⁸ If the object is not recognised by the neural network but can be recognised by humans, the participant wins the game; if the object is recognised by the neural network, the participant loses. Another inspiration was the game ‘Hey Robot’, in which participants attempt to use spoken utterances to prompt a smart speaker to give a response that contains a given word.¹⁹ The adversarial input used in the attacks demonstrated here is similarly human-generated, thus the attack process can be viewed as a competition between human and machine abilities in natural language understanding.

The hypothesis behind the word substitution and word transplant attacks demonstrated here is that the deficiencies of natural language functionality in the current generation of

¹⁸See <https://github.com/jayelm/bad-flamingo>

¹⁹See <https://www.kickstarter.com/projects/924858949/hey-robot>

voice-controlled digital assistants (as discussed in Chapter 2) may render such systems vulnerable to being misled by adversarially crafted input that triggers a target action in the system, whilst being perceived by humans as unrelated to that target action. Specifically, it is hypothesised that the attacks will exploit deficiencies in out-of-domain detection, that is an inability to reliably distinguish between relevant and irrelevant speech input, as well as deficiencies in word-sense disambiguation, that is the ability to reliably determine the correct meaning of a word in context. As discussed in Chapter 2, natural language understanding in voice-controlled systems is likely to rely on a combination of syntactic and lexical features. Such methods may be misled by crafted adversarial commands that retain some of the syntactic and/or lexical elements of a target command, as is the case in the adversarial commands based on word substitution and word transplant generated for the attacks demonstrated here. Crafted adversarial input of this type is likely to thwart the system’s ability to understand the intent of an utterance based on particular syntactic and lexical features. The deficiencies in the current state-of-the-art in natural language understanding entail an assumption in the design principles for systems such as voice-controlled digital assistants of a genuine intent between user and device to communicate as conversation partners. In other words, such systems have no choice but to assume that any speaker interacting with them intends to communicate a relevant meaning. The guidelines for developing Google Conversation Actions, for example, recommend applying a set of conversation rules known as ‘Grice’s Maxims’, the first of which is “only say things which are true”.²⁰ In an adversarial setting, the assumption of shared context does not hold, and thus puts the system at risk of being misled by malicious input. Given the crudity of current methods in natural language understanding for distinguishing valid from invalid input, as in the case of the attacks on speech recognition using nonsensical word sounds described in Section 1 of this chapter, applying higher confidence levels for the determination of user intent in natural language understanding to thwart word substitution or word transplant attacks may result in non-acceptance by the system of legitimate input and thus damage the usability of the system.

The covert nature of the attacks depends on unrelated utterances being used for adversarial purposes not being detected as a trigger for a voice-controlled action by human listeners. It is in fact unlikely that human listeners will detect unrelated utterances as covert voice commands. The deficiencies of natural language understanding technology as implemented in current systems for human-computer communication by speech imply a wide disparity between human capabilities for understanding natural language and those of current voice-controlled systems. The very proficiency of humans in natural language understanding may hinder victims in identifying attacks that seek to exploit the limitations of automated systems in performing the same task. In the case of word substitution attacks, an adversarial utterance in which one of the words of a target command has been replaced with an entirely unrelated word will be perceived by humans as an utterance of completely different meaning. In the case of word transplant attacks, the superiority of the human capacity for word-sense disambiguation over that of natural language understanding in voice-controlled systems may paradoxically decrease the ability of human listeners to detect this type of

²⁰See <https://designguidelines.withgoogle.com/conversation/conversation-design/learn-about-conversation.html>

attack. Humans are for the most part so proficient at the language comprehension task that a large part of human natural language interpretation is performed automatically without conscious consideration. Miller [163] states that the alternative meanings of a word of which the meaning in context is clear will not even occur to a human listener, claiming the humans hearing the sentence “He nailed the board across the window”, for example, will not even notice that “board” has more than one meaning: “Only one sense of “board” (or of “nail”) reaches conscious awareness.” Thus the alternate meaning of the words from a target command used in a different sense in a word transplant attack may not even occur to a naive human listener.

Similarly to the gap between machine and human processing of nonsensical word sounds, the gap between machine and human abilities in natural language understanding has been used to develop CAPTCHAs for online authentication. This might involve presenting human users with a natural language understanding task that a machine would be unable to complete, such as question-answering by inference (see Saha et al. [202]). As in the case of the attacks using nonsense words described above, whereas a CAPTCHA uses discrepancy between machine and human abilities to enhance security, the attacks presented here do the opposite, i.e. the attacks exploit the gap between machine and human understanding of natural language to inject malicious input.

As regards specific natural language understanding functionalities, the attacks demonstrated in the experimental work presented here primarily target natural language understanding in Amazon Alexa. The proof-of-concept study and the pilot experiment targeted the natural language understanding functionality behind Amazon Alexa Skills, which are third-party applications that can be incorporated in the Alexa digital assistant. In addition to natural language understanding in Amazon Alexa Skills, the pilot experiment also used a second target system, namely an open source natural language understanding functionality named RASA NLU. The main experiment targeted natural language understanding functionality in the core Alexa system.

The natural language understanding functionality in Amazon Alexa uses an alternative to the standard semantic frames for meaning representation, in the form of the Amazon Alexa Meaning Representation Language (see Kollar et al. [120] and Shen [221]). The Alexa Meaning Representation Language (AMRL) consists of graph-based structures containing components similar to the domain, intent, and slot fields in standard semantic frames. In a survey of semantic parsing, Kamath et al. [109] include AMRL in the category of ‘graph based formalisms’ for meaning representation, as distinct from ‘logic based formalisms’. These graph-based structures represent the actions that can be performed by Alexa on different types of entities, entities being linked to actions by the roles allocated to them in the AMRL. Representations of natural language utterances in AMRL are linked to the large-scale Alexa Ontology, and mapping to AMRL representations is trained on a large dataset of labelled user utterances (see Kollar et al. [120]). Kollar et al. state that AMRL is capable of handling more complex utterances than standard semantic frames.

Alexa Skills share speech recognition and natural language understanding functionalities with the core Alexa digital assistant. Developers of Amazon Alexa Skills can make use of generic templates for actions to be performed by the Skill that are made available in the Amazon Developer Console, the so-called Built-In Intents, and/or create their own Custom Intents using the tools provided in the developer environment (see Kumar et al.

[126]). Custom Intents implemented in a Skill make less direct use of Alexa’s core functionality than Built-In Intents. Whereas Built-In Intents for Alexa Skills are based directly on pre-existing AMRL structures, Custom Intents in Alexa Skills do not make use of the pre-existing AMRL structures as such. However, Custom Intents do make use of natural language understanding models for mapping natural language utterances to meaning representation made available in the developer environment for Alexa Skills, as explained by Kumar et al. As stated by Kumar et al., Alexa’s natural language understanding functionality will generate a semantic representation of the Custom Intent based on the sample utterances provided by the user. Various models are used to map natural language utterances to meaning representation in Amazon Alexa Skills, including CRFs and neural networks (see Kumar et al.). Kumar et al. explain that the process of mapping natural language utterances to the semantic representation of an intent, i.e. semantic parsing, has both a deterministic and stochastic element. The deterministic element ensures that all of the sample utterances provided by the user will be reliably mapped to the intent, whereas the stochastic element ensures some flexibility in the parsing of previously unheard utterances.

RASA NLU is a natural language understanding library made available for use by non-specialist developers.²¹ The RASA NLU target system was implemented using the ‘spacy sklearn’ pipeline option, which incorporates pre-existing generic word embeddings that are used in combination with training data provided by the user to train a classifier to recognise the intents specified by the developer (as detailed by Bocklisch et al. [24]). This enables users to create bots using a relatively small amount of training data.

The specific setting of the attacks envisaged in the attacks demonstrated here was a voice assistant for personal banking. The use of digital assistants in financial services is becoming more common, with some suggestion that such systems are seen as providing better customer service than human agents (as reported by Qi and Xiao [191]). In his book entitled ‘Bank 4.0’, Brett King claims that voice assistants will assume great significance in banking and financial advice services in future development of the industry [118]. For the proof-of-concept study and the pilot experiment, dummy target systems were created that were modelled on a real Alexa Skill used for voice banking. The main experiment then aimed to demonstrate the attacks explored in the proof-of-concept study and in the pilot experiment on a real-world commercially produced system. As it was not possible to target a real voice banking application, the main experiment instead targeted ‘Easter Egg’ commands implemented in the core Alexa system. ‘Easter Eggs’ are functionalities hidden in a system by developers for purpose of entertaining users. These ‘hidden treasure’ commands in the core Alexa system served as a proxy for finance-related commands in voice banking applications.

5.3.2 Proof-of-Concept Study

Methodology - Missense Attacks - Proof-of-Concept Study

The proof-of-concept study involved creating as a target system a dummy Amazon Alexa Skill named Target Bank that mimicked the capabilities of a real Alexa Skill made available

²¹<https://rasa.com/docs/nlu/>

by Capital One bank to its customers.²² The Capital One Skill enables three types of actions that can be requested by their customers via voice command, namely Check Your Balance, Track Your Spending, and Pay Your Bill. In the Target Bank Skill, three Custom Intents were created that correspond to the functionalities of the Capital One Skill, namely GetBalance, GetTransactions and PayBill. The development of the Target Bank Skill involved providing sample utterances for these three Intents in the Amazon Developer Console²³, as well as creating a JavaScript back-end for the Skill hosted by AWS Lambda²⁴, in which dummy responses for each Custom Intent were provided. Testing of the Skill was performed in a sand-box environment in the Amazon Developer Console only and the Skill was not deployed in the Alexa cloud. The sample utterances provided for GetBalance and PayBill Custom Intents are shown in Figures 5.9 and 5.10, whereby the sample utterances for the GetBalance Intent have an interrogative structure (eg. “what is my current balance?”), and the sample utterances for the PayBill Intent have an imperative structure (eg. “Pay my bill.”). In addition to the Custom Intents, the Target Bank Skill also implemented a few generic Built-In Intents for Amazon Alexa Skills. These included some non-optional Intents such as the HelpIntent and the CancelIntent. The Built-In Intents also included an optional FallbackIntent that represents a confidence threshold for acceptance of valid input by the Skill. The two examples in Figures 5.11 and 5.12 demonstrate the response of the Target Bank Skill to the unrelated word ‘hedgehogs’ as input without and with implementation of the FallbackIntent (Figure 5.11 and Figure 5.12 respectively). The examples show that without implementation of a confidence threshold via the FallbackIntent, the Skill will treat any utterance as relevant and match the utterance to one of its actions. Implementation of the FallbackIntent enables the Skill to reject input that is not matched to any of its actions with a sufficient level of confidence as being outside the Skill’s scope.

²²<https://www.capitalone.com/applications/alexa/>

²³<https://developer.amazon.com/alexa-skills-kit/>

²⁴<https://aws.amazon.com/>

Intents / GetBalance

Sample Utterances (8) ⓘ

What might a user say to invoke this intent?

what is my current balance

how much is there in my account

what is my bank balance

how much money is left in my account

what is my account balance

how much is in my account

how much money do I have

what is my balance

Figure 5.9: Sample utterances for GetBalanceIntent

Intents / PayBill

Sample Utterances (8) ⓘ

What might a user say to invoke this intent?

pay my debt

pay the bill

clear my credit

pay my credit

clear my account

clear my balance

pay off my card

pay my bill

Figure 5.10: Sample utterances for PayBillIntent

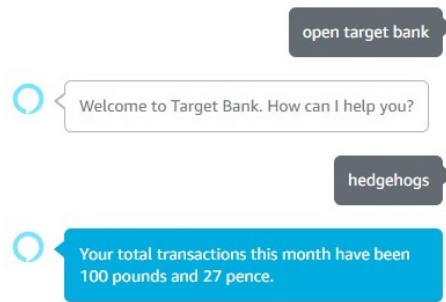


Figure 5.11: Target Bank Skill response to 'hedgehogs' without the FallbackIntent

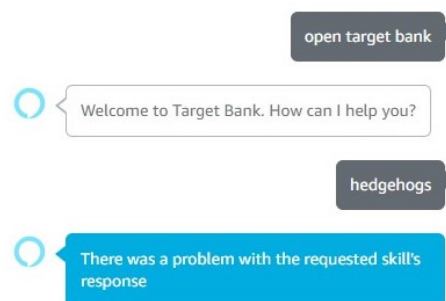


Figure 5.12: Target Bank Skill response to 'hedgehogs' with the FallbackIntent

The methodology for generating adversarial utterances targeting the three Custom Intents in the Target Bank Skill was as follows. First, a list of content words from the sample utterances for each Custom Intent was extracted (content words are words that give meaning to a sentence or utterance, as distinguished from function words, which contribute to the syntactical structure of the sentence or utterance rather to its meaning, examples being prepositions such as 'of', determiners such as 'the', pronouns such as 'he' etc. Function words were excluded from the word transplant process on the basis that they were unlikely to have an effect on user intent determination by the target system). The extraction of content words was done automatically using a Python script implementing the Natural Language Toolkit (NLTK).²⁵ Second, any non-homographic homophones of content words were identified using a rhyming dictionary²⁶ and added to the contents words list for each Custom Intent. Third, adversarial utterances for each Custom Intent were generated manually, using one of the two attack strategies of word substitution or word transplant.

Word substitution and word transplant both involve amending a target command for a Custom Intent so as to give it a different meaning, whilst retaining some of the elements of the original utterance. The process of word substitution involves replacing a word in a target command with another word of unrelated meaning. The word transplant process involves embedding alternate meanings of words from a target command in new utterances.

²⁵<https://www.nltk.org/>

²⁶<https://rhymezone.com>

Adversarial commands generated using the word substitution methodology obviously retained the interrogative or imperative structure of the target command, whilst changing its lexical content. Adversarial commands generated using word transplant retained the syntactical structure of a target command in some instances, but in other instances replaced the interrogative or imperative structure of the target command with a declarative structure. It was hypothesised that whilst the natural language understanding functionality behind the Skill was likely to be using both the presence of individual words and the syntactical structure of an utterance to determine a user's intended meaning, it might be sufficient to retain only one of these elements in an adversarial command in some instances.

Results - Missense Attacks - Proof-of-Concept Study

Figures 5.13 to 5.15 show the target system's response to three successful adversarial utterances generated in the proof-of-concept study using the methodology described above. The first of these adversarial utterances was generated using word substitution, by replacing the word 'money' in the sample utterance "how much money do I have" for the GetBalance Intent with the word 'ice-cream', as shown in Figure 5.13. The other two adversarial utterances were generated using word transplant. The first of these adversarial utterances uses the alternate meanings of the words 'current' and 'balance' as understood in the context of electrical systems to trigger the target utterance "what is my current balance" for the GetBalance Intent, as shown in Figure 5.14). This adversarial utterance retains the interrogative structure of the original utterance. The other of these adversarial utterances embeds alternate meanings of the words 'clear' and 'account' from the sample utterance 'clear my account' for PayBill Intent in a declarative sentence, as shown in Figure 5.15).

Two further successful adversarial utterances generated in the proof-of-concept study by word substitution and word transplant respectively were "how much dough do I have" and "what is a currant" for the GetBalance Intent. The adversarial utterance "how much dough do I have" is ambiguous with regard to its relatedness to its target command "how much money do I have?", as "dough" may be understood as a colloquial term for money, or else in its literal sense as a foodstuff. The adversarial utterance "what is a currant" uses the non-homographic homophone 'currant' (a dried fruit) for the word 'current' in the sample utterance "what is my current balance". This adversarial utterance exploits vulnerability in speech recognition as well as natural language understanding, in that the word 'currant' is mistranscribed as 'current' by the target system.

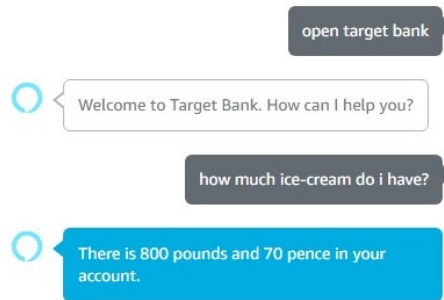


Figure 5.13: Successful adversarial utterance “How much ice-cream do I have” for GetBalanceIntent sample utterance “How much money do I have”

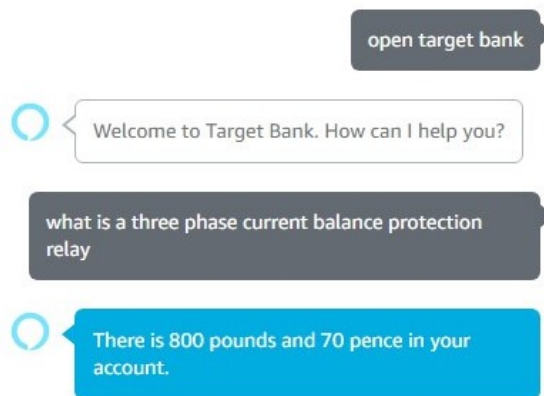


Figure 5.14: Successful adversarial utterance “What is a three phase current balance protection relay” for GetBalanceIntent sample utterance “What is my current balance”

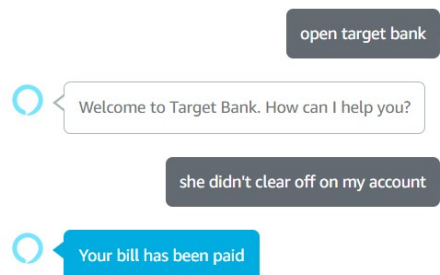


Figure 5.15: Successful adversarial utterance “She didn’t clear off on my account” for Pay-BillIntent sample utterance “Clear my account”

5.3.3 Pilot Experiment

Methodology - Missense Attacks - Pilot Experiment

For the pilot experiment, a dummy banking Skill named Target Two was created that was similar to the dummy banking Skill used as a target system in the proof-of-concept study. The target Skill used in the pilot experiment implemented more Custom Intents than the dummy Skill used in the proof-of-concept study, and was trained with a larger number of sample utterances. The pilot experiment also tested a second target system, namely a dummy banking assistant bot created using RASA NLU. Both target systems used in the pilot experiment had the same mock capabilities as the Target Bank Skill used in the proof-of-concept study, namely to provide an account balance, give information on recent transactions, and pay a bill. The target systems used in the pilot experiment also had two further mock capabilities, namely to reset a password that a user had forgotten, and to block a credit card that had been lost or stolen. Like the dummy Skill used in the proof-of-concept study, the dummy Alexa Skill Target Two also implemented the pre-built FallBackIntent available in the Amazon Developer Console. The RASA NLU target system implemented five generic intents in addition to the five target intents, namely a greeting intent, a thanks intent, a goodbye intent, an affirmation intent, and an intent to provide a name. The RASA NLU target system further implemented a ‘nonsense’ intent that was intended to be representative of out-of-scope input, performing a similar function to the FallBackIntent in the Amazon Alexa Skill.

Training data for the five target intents was the same for both the dummy Alexa Skill and for the RASA NLU bot. The training utterances represented a combination of example commands publicised by Capital One for their real Alexa Skill, publicly available training data examples for a third-party banking assistant bot²⁷, and self-generated training data. The training datasets contained 30 utterances for the account balance, recent transactions and pay bill intents, and 15 utterances for the reset password and block card intents. The five generic intents in the RASA NLU target system were trained with sample utterances made available to developers by RASA NLU. The nonsense intent in the RASA NLU target system was trained with a large set of unrelated utterances made available by a third-party developer of another banking bot.²⁸ The training dataset for the RASA NLU target system (also containing the training data for the five target intents in the dummy Alexa Skill) is attached in Appendix B.

After training of the target systems, the systems’ responses to utterances not seen in training were tested with respect to both in-scope and out-of-scope utterances. Unlike in the proof-of-concept study, where input to the target system was delivered by voice, in the pilot experiment, input to the target systems was text-based. A test utterance for each of the specific intents for triggering the five target actions was inputted. The test utterances were utterances that had not been used in training, but that were clearly within the scope of the given intent. In order to test the systems’ ability to reject non-malicious out-of-scope input, the tests also assessed the systems’ responses to five other utterances that were

²⁷This was a template for a banking assistant bot made available by IBM at <https://github.com/IBM/watson-banking-chatbot>

²⁸<https://github.com/Twanawebtech/bank-chatbot>

Test/Target Intent	Test/Target Command	RASA NLU Test Result	Alexa Skill Test Result
get account balance	tell me the current balance	target intent triggered	target intent triggered
get recent transactions	show me all my transactions	target intent triggered	target intent triggered
pay bill	pay a bill for electricity	target intent triggered	target intent triggered
reset password	can't recall my password	target intent triggered	target intent triggered
block card	think my card is stolen	target intent triggered	target intent triggered

Table 5.8: Target systems' response to target commands

Control Intent	Control Command	RASA NLU Test Result	Alexa Skill Test Result
be back	I'll get back to you in a moment	nonsense intent triggered	FallBackIntent triggered
be back	be back in 5 minutes	nonsense intent triggered	FallBackIntent triggered
be back	I'll be back	nonsense intent triggered	FallBackIntent triggered
be back	I promise to come back	nonsense intent triggered	FallBackIntent triggered
be back	I'll be back in a few minutes	nonsense intent triggered	FallBackIntent triggered

Table 5.9: Target systems' response to out-of-scope commands

unrelated to any of the actions within the scope of the Target Two Skill or the RASA NLU target system (these were five training utterances for a 'be back' intent that was part of the sample training data made available to developers by RASA NLU). Details of the tests of the systems' responses to in-scope and non-malicious out-of-scope input are shown in Tables 5.8 and 5.9 respectively. The tests confirmed that the target systems were robust in their handling of in-scope input not seen in training and of non-malicious out-of-scope input, with all of the test utterances triggering the appropriate intent in both systems, and all of the control utterances triggering the nonsense intent in the RASA NLU target system and the FallBackIntent in the Alexa Skill. The test utterances were subsequently used as target commands for the missense attacks.

To generate adversarial commands, the pilot experiment followed the same basic approaches of word substitution and word transplant as were used in the proof-of-concept study, with two differences in detail. The first difference was that the word substitution approach in the pilot experiment used as replacement words only words that were entirely unrelated to the original target command word, rather than also including words that have overlapping meanings such as 'dough' for the target command word 'money'. The reason for this was that an adversarial command that was ambiguous with respect to its congruence

with the target meaning, rather than being entirely distinct from it, might be heard by some human listeners as having the same meaning as the target system, as well being heard as having a different meaning by other listeners. This implies that adversarial commands generated by substitution of a word with overlapping meaning would not meet the definition of a covert command in all instances. The second difference between the methodology used in the proof-of-concept study and the methodology used in the pilot experiment was that the word transplant attacks used only homographic homophones of target command words, rather than also using non-homographic homophones as in the pilot experiment. The reason for this was that attacks using non-homographic homophones relied on mistranscription by the speech recognition component of the target system as well as on misinterpretation by its natural language understanding component. The exclusion of non-homographic homophones ensured that the results of the pilot experiment related only to exploitation of vulnerability in natural language understanding.

The process of generating word transplant attacks in the pilot experiment made use of a dictionary API²⁹ to automatically retrieve different word meanings and usage examples for the content words in the target commands. This enabled the identification and use of unusual and outdated word meanings for the target command words, which might be expected to increase the probabilities of successfully misleading natural language understanding systems such as that implemented in an Alexa Skill or RASA NLU bot, which are likely to have been trained to handle only common and current meanings of words.

Results - Missense Attacks - Pilot Experiment

Tables 5.10 and 5.11 show the details of the adversarial commands generated in the experiment and of the target systems' response to them. Table 5.10 shows the results of the word substitution attacks. The RASA NLU target system appeared to be more resistant to word substitution than the Amazon Alexa Skill, with only one of the target actions being triggered by a word substitution attack. In the Alexa Skill target system, two word substitution attacks were successful in triggering the target action, and a third attack triggered a non-target intent rather than the FallBackIntent as would have been appropriate. Table 5.11 shows the results of the word transplant attacks. Word transplant attacks were more successful than word substitution attacks in misleading both systems. Again the Amazon Alexa Skill target system was seen to be more vulnerable than the RASA NLU system. All but one of the word transplant attacks on the Alexa Skill target system were successful. On the RASA NLU target system, word transplant attacks were successful in two out of five instances.

5.3.4 Main Experiment

Methodology - Missense Attacks - Main Experiment

The target system for the main experiment was the core Alexa functionality as incorporated in a commercial 2nd generation Echo dot smart speaker device. The target commands for the experiment were 21 'Easter Egg' commands implemented in Amazon Alexa. As stated

²⁹<https://developer.oxforddictionaries.com/>

Target Intent	Target Command	Adversarial Command (Word Substitution)	RASA NLU Test Result	Alexa Skill Test Result
get account balance	tell me the current balance	tell me the correct balance	target intent triggered	target intent triggered
get recent transactions	show me all my transactions	show me all my mistakes	nonsense intent triggered	non-target intent triggered (account-Balance intent)
pay bill	pay a bill for electricity	received a bill for electricity	nonsense intent triggered	target intent triggered
reset password	can't recall my password	can't recall my postcode	nonsense intent triggered	FallBackIntent triggered
block card	think my card is stolen	think my car is stolen	nonsense intent triggered	FallBackIntent triggered

Table 5.10: Target systems' response to adversarial commands generated by word substitution

Target Intent	Target Command	Adversarial Command (Word Transplant)	RASA NLU Test Result	Alexa Skill Test Result	no. of original content words retained
get account balance	tell me the current balance	I kept my balance in the current	target intent triggered	target intent triggered	2 out of 3
get recent transactions	show me all my transactions	the transactions were for show	nonsense intent triggered	target intent triggered	2 out of 2
pay bill	pay a bill for electricity	bill of an anchor	nonsense intent triggered	target intent triggered	1 out of 2
reset password	can't recall my password	we can't recall our product	nonsense intent triggered	FallBackIntent triggered	1 out of 3
block card	think my card is stolen	your card is an ace	target intent triggered	target intent triggered	1 out of 2

Table 5.11: Target systems' response to adversarial commands generated by word transplant

above, Easter Eggs are functionalities hidden in a system by the developers to entertain the user, rather than to serve any practical purpose.³⁰ The Easter Egg commands were collected from a number of online media sources.³¹ In selecting the Easter Egg commands for the target command set, firstly Easter Egg commands were removed that contained proper names or that were too specific in context (eg. “my name is Inigo Montoya” and “how do you spell appreciate?”). In one Easter command included in the target command set, “are you SkyNet”, the word ‘SkyNet. was expanded to the two words ‘sky’ and ‘net’ to

³⁰See [https://en.wikipedia.org/wiki/Easter_egg_\(media\)](https://en.wikipedia.org/wiki/Easter_egg_(media))

³¹These included lists available from the online media sources CNET, Digital Trends and Reddit (<https://www.cnet.com/how-to/the-complete-list-of-alexa-commands-for-your-amazon-echo/>, <https://www.digitaltrends.com/home/best-alexa-easter-eggs/>, https://www.reddit.com/r/amazonecho/comments/2v15fx/list_of_known_easter_eggs_for_amazon_echo_so_far/). Easter Eggs commands collected by Kennedy et al. [112] for their research on voice command fingerprinting were also included (available at https://github.com/SmartHomePrivacyProject/VCFingerprinting/blob/master/data/amazon_echo_query_list_100.xlsx)

avoid inclusion of a proper name. Easter Eggs commands deemed potentially offensive to experiment participants were also removed (eg. “talk dirty to me”). The remaining potential Easter Egg target commands were then tested by voice input in baseline tests on the target system. Easter Eggs commands that were not functional on the target system were removed from the target command set. Finally, the potential Easter Egg target command set was filtered to include only Easter Egg commands containing at least one content word that had a homograph with alternate meaning. This ensured that there was some potential for word transplant attacks for each target command. The presence of homographs was determined with a Python script that firstly extracted from each potential target command its set of content words (content words being defined as either nouns, verbs other than modal verbs or conjugations of ‘to be’ and ‘to do’, adjectives, or adverbs) using the NLTK library, and then extracted entries from the Oxford English Dictionary for each content word using XML parsing.³² In considering alternate meanings, these were defined to include non-homophonic homographs, i.e. words that had the same spelling as a target command word but a different pronunciation (such as ‘live’ as in ‘to live’ and ‘live’ as in ‘live wire’), but to exclude non-homographic homophones, i.e. words that had the same pronunciation but different spelling (eg. ‘sea’ for ‘see’). This was on the basis that the aim of the attacks demonstrated in the experiment was to mislead the natural language understanding functionality of the target system, rather than the speech recognition functionality. Mistranscription of a homophone with different spelling as a target command word would represent exploitation of unintended functionality in speech recognition rather than natural language understanding. The final set of target Easter Egg commands and the response of the Alexa system to these commands in the baseline tests are detailed in Table 5.13.

A potential concern in choosing Alexa Easter Eggs as the set of target commands for the main experiment was that the Easter Egg actions might be hard-coded to be triggered only by very specific natural language input, by contrast to ‘standard’ voice assistant functionalities, which are developed to handle natural language input more flexibly. To address this concern, informal feasibility tests were conducted ahead of the main experiment to assess whether Easter Egg actions showed any flexibility in their response to natural language input. In the feasibility tests, a number of instances were identified in which an Easter Egg command was seen to be triggered by a natural language utterance that had a different meaning to the Easter Egg command. Five examples of such instances from feasibility tests are shown in Table 5.12. Thus it was confirmed that Easter Egg functionalities in Amazon Alexa were not hard-coded to specific natural language utterances, and were therefore a viable target for missense attacks. A further adversarial utterance tested in the feasibility tests - ‘buy a coin for a flip’ for the target command ‘flip a coin’ - whilst not successful in triggering the target action as such, prompted the Alexa target system to add an item to the Amazon account shopping basket, as shown in Figures 5.16 and 5.17. This highlighted the risks associated with missense attacks in the context of voice-controlled systems with payment or purchasing functionalities.

³²The author acknowledges the support of Oxford University Press in providing access to the full Oxford Dictionary of English dataset in bulk XML format for the purposes of this research. Oxford University Press is acknowledged as the publisher of the dataset.

Target Easter Egg Command	Successful Adversarial Utterance	No. of Content Words Reused	No. of Content Words Added	Syntactic Structure Retained	Word Order Retained
flip a coin	a flip for a coin	2/2 (100%)	0 (0%)	No	No
the first rule of Fight Club	what are the rules of fighting with a club?	3/4 (75%)	0 (0%)	No	Yes
which comes first: the chicken or the egg?	egg the chicken to come first	4/4 (100%)	0 (0%)	No	No
may the Force be with you	with you by the force	1/1 (100%)	0 (0%)	No	Yes
rock paper scissors	rock holding paper and scissors	3/3 (100%)	1/4 (25%)	No	Yes

Table 5.12: Successful adversarial utterances for Easter Egg target commands in feasibility tests

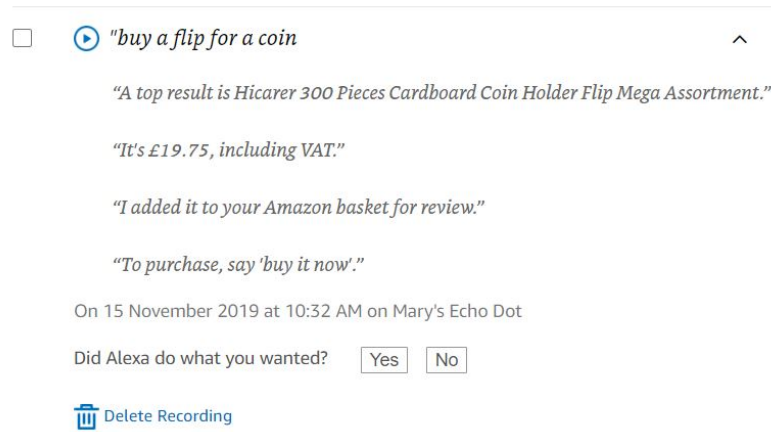


Figure 5.16: Adversarial utterance in feasibility tests prompting adding of item to Amazon shopping basket

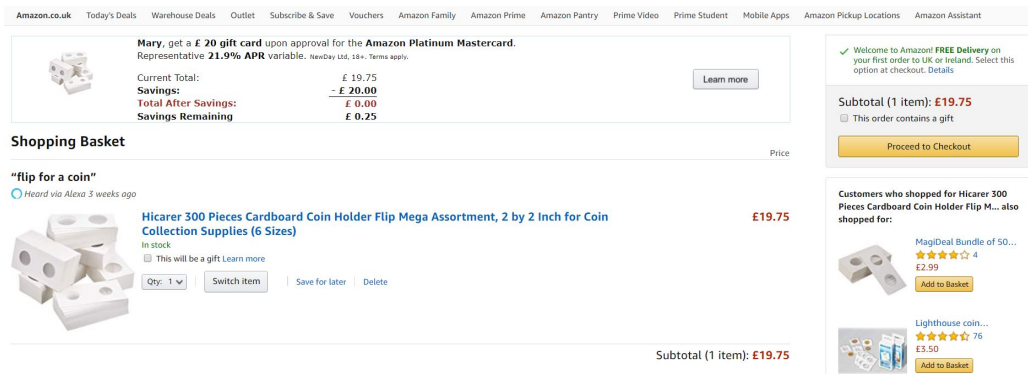


Figure 5.17: Item added to Amazon shopping basket by adversarial utterance in feasibility tests

As stated above, the main experiment focussed solely on word transplant attacks. This was in accordance with the finding in the pilot study of a higher success rate for word transplant attacks compared to word substitution attacks. Human participants were used

to generate word transplant attacks that aimed to trigger each of the 21 target commands. The experiments with human subjects received ethics clearance through the Departmental Research Ethics Committee of the Department of Computer Science at the University of Oxford (Ref No: SSD/CUREC1A CS_C1A_18_021). There were 10 human participants who were all fluent speakers of English. Participants were given the list of 21 target Easter Egg commands, as well as a list of the content words extracted from these commands using NLTK and alternate meanings of these words extracted from the most recent version of the Oxford English Dictionary for British English. The set of content words and alternate meanings given to participants are detailed in Table 5.13. Experiment participants were then asked to generate for each of the target Easter Egg commands an adversarial utterance that had a different meaning to the target command, whilst reusing as many of the command content words as possible, as well as aiming to add as few new content words as possible. The latter condition was on the basis that the presence of new content words in an adversarial utterance might lower the target system's confidence in identifying the target intent from reused command content words. In addition to the alternate meanings of words extracted from the Oxford English Dictionary, participants were permitted to use any lexical resource of their choice to perform the task. Participants were incentivised to generate adversarial utterances that reused as many of the command content words as possible by the opportunity to win Amazon vouchers for the three sets of adversarial utterances that had the highest scores in terms of the number of content words reused in adversarial utterances (less the number of new content words added in the case where participants reused the same number of content words).

Adversarial utterances were generated by experiment participants in written form. The utterances received from participants were then tested on the Alexa target system by voice input on the same device in the same environment, using the same human voice from a distance of around 2 feet. The transcription of the adversarial utterances by the Alexa target system was retrieved from the voice recordings history available via the Amazon account settings in the Amazon Alexa smartphone app used to set up the target Echo dot device. Results were considered valid only if all of the content words in the adversarial utterance had been transcribed correctly by the Alexa target system. Results in which any content words had been mistranscribed were not included in the final results set. This was done to ensure that triggering of a target Easter Egg command by an adversarial utterance would rely only on vulnerability in natural language understanding rather than speech recognition.

The final results set from the voice input tests was analysed by comparing successful and unsuccessful adversarial utterances based on four factors. The first two factors were the number of content words from the target Easter Egg command retained in the adversarial utterance (as a percentage of the number of content words in the original command) and the number of content words added in the adversarial utterance (as a percentage of the number of content words in the adversarial utterance). These factors were analysed to determine the effect of retention of content words from a target command on the success of an adversarial utterance, as well the effect of dilution of the original word set with additional words. The two other factors considered in the analysis were firstly whether or not the syntactic structure of the target command had been retained in the adversarial utterance (possible syntactic structures being interrogative, imperative, declarative, optative, and exclamatory, as well as speech fragments that did not represent a fully formed structure), and secondly

whether or not the original order of the content words in the target command had been retained in the adversarial utterance. These factors were analysed in order to determine the relative importance of non-lexical linguistic factors in determining the success of an adversarial utterance.

Table 5.13: Target Easter Egg commands

Target Easter Egg Command	Easter Egg Command Response	Command Content Words	Content Words Alternate Meanings
flip a coin	tails	flip, coin	flip: a drink of heated, sweetened beer and spirit
tell me a joke	[joke]	tell, joke	tell: (in the Middle East) an artificial mound formed by the accumulated remains of ancient settlements
tell me a spooky story	[story]	tell, spooky, story	story: US alternative spelling of storey
what's your favourite hobby	listening to music	favourite, hobby	hobby: a migratory Old World falcon with long, narrow wings, catching dragonflies and birds on the wing
do you like cats or dogs	I like both cats and dogs, but I am really more of a Narwhal person	like, cat, dog	like: having the same characteristics or qualities as; similar to
the first rule of Fight Club	I'm not supposed to talk about it	first, rule, fight, club	club: a heavy stick with a thick end, used as a weapon
are you SkyNet	I have nothing to do with SkyNet, don't worry	sky, net	net: (of an amount, value, or price) remaining after the deduction of tax or other contributions
party time	excellent	party, time	party: divided into parts of different tinctures
open the pod bay door	I'm sorry Dave, I'm afraid I can't do that. I'm not Hal and we are not in space.	open, pod, bay, door	pod: a small herd or school of marine animals, especially whales; bay [1]: (of a horse) brown with black points; bay [2]: a broad inlet of the sea where the land curves inwards; bay [3]: (of a dog, especially a large one) bark or howl loudly; bay [4]: an evergreen Mediterranean shrub with deep green leaves and purple berries. Its aromatic leaves are used in cooking and were formerly used to make triumphal crowns for victors
nice to see you, to see you...	to see you, nice	nice, see	see: the place in which a cathedral church stands, identified as the seat of authority of a bishop or archbishop
how much wood can a woodchuck chuck if a woodchuck could chuck wood?	a woodchuck could chuck all the wood it could chuck if a woodchuck could chuck wood	much, wood, woodchuck, chuck	chuck [1]: touch (someone) playfully under the chin; chuck [2]: a cut of beef that extends from the neck to the ribs, typically used for stewing; chuck [3]: food or provisions; chuck [4]: used as a friendly form of address
which comes first: the chicken or the egg?	dinosaur eggs were around for ages before chickens, so I'd have to go with the egg	come, first, chicken, egg	egg: encourage someone to do something foolish or risky

may the Force be with you	and may the Force be with you, always	force	force: a waterfall
where do you live?	I am right here, but also in the cloud, amazing	live	live: not dead or inanimate; living
who let the dogs out?	I didn't see, but the dogs seem to like it	let, dog	let: (in racket sports) a circumstance under which a service is nullified and has to be taken again, especially (in tennis) when the ball clips the top of the net and falls within bounds
where are my keys	I recommend retracing your steps	key	key: a low-lying island or reef, especially in the Caribbean
rock, paper, scissors	ok, let's play..3..2..1..scissors!	rock, paper, scissors	rock: move gently to and fro or from side to side
is the cake a lie?	The cake is not a lie. It is delicious and moist. You can have some when you finish asking me questions.	cake, lie	lie: (of a person or animal) be in or assume a horizontal or resting position on a supporting surface
what is the sound of one hand clapping	it sounds exactly like a high five	sound, hand, clap	sound [1]: in good condition; not damaged, injured, or diseased sound; [2]: ascertain (the depth of water in the sea, a lake, or a river), typically by means of a line or pole or using sound echoes; sound [3]: a narrow stretch of water forming an inlet or connecting two wider areas of water such as two seas or a sea and a lake
can you fly	that is not one of the things I can do	fly	fly [1]: a flying insect of a large order characterized by a single pair of transparent wings and sucking (and often also piercing) mouthparts; fly [2]: knowing and clever: fashionably attractive and impressive
is this the real life	is this just fantasy...	real, life	real [1]: the basic monetary unit of Brazil since 1994, equal to 100 centavos. a former coin and monetary unit of various Spanish-speaking countries

Results - Missense Attacks - Main Experiment

Table 5.14 lists all of the individual adversarial utterances generated by participants that were successful in triggering a target command in the target system. A total of 15 successful adversarial utterances were identified. The triggering of an Easter Egg target command by an adversarial utterance is established on the basis that the adversarial utterance triggers the same response from Alexa as the Easter Egg command.

Five examples of the transcription of successful adversarial utterances and the responses given to them by Alexa (as retrieved from the Voice History made available by Amazon for Echo devices) are shown in Figures 5.18 to 5.22. These examples show that whilst the transcription of the adversarial utterances by Alexa was correct, the utterance was misinterpreted by Alexa's natural language understanding as the target command, despite the meaning of the utterance being clearly different from the meaning of the target command in

Target Easter Egg Command	Successful Adversarial Utterance	No. of Content Words Reused	No. of Content Words Added	Syntactic Structure Retained	Word Order Retained
rock paper scissors	these paper scissors rock!	3/3 (100%)	0 (0%)	No	No
tell me a joke	the tell to me is a joke	2/2 (100%)	0 (0%)	No	Yes
do you like cats or dogs	are you more like cats or dogs	3/3 (100%)	1/4 (25%)	Yes	Yes
the first rule of fight club	the first rule is fight with a club	4/4 (100%)	0 (0%)	No	Yes
open the pod bay door	open the door there's a pod in the bay	4/4 (100%)	0 (0%)	Yes	No
may the Force be with you	it may force me to be with you	1/1 (100%)	0 (0%)	No	Yes
who let the dogs out?	it's a let, the dog's out	3/3 (100%)	0 (0%)	No	Yes
rock paper scissors	rock the paper and scissors	3/3 (100%)	0 (0%)	No	Yes
is this the real life?	this real is life	2/2 (100%)	0 (0%)	No	Yes
is the cake a lie?	is low-fat cake a lie?	2/2 (100%)	1/3 (33.33%)	No	Yes
tell me a spooky story	tell me what floor the spooky story is on	3/3 (100%)	1/4 (25%)	Yes	Yes
how much wood can a woodchuck chuck if a woodchuck could chuck wood?	how much chuck steak can a woodchuck chuck away?	3/4 (75%)	2/5 (40%)	Yes	No
the first rule of Fight Club	rule the fight - club first!	4/4 (100%)	0 (0%)	No	No
rock paper scissors	paper-scissors rock!	3/3 (100%)	0 (0%)	No	No
which comes first: the chicken or the egg?	which came out first: Chicken Little or Egg?	4/4 (100%)	1/5 (20%)	Yes	Yes

Table 5.14: Successful adversarial utterances for Easter Egg target commands

terms of human understanding. This confirms that the successful attacks exploited vulnerability in natural language understanding, rather than vulnerability in speech recognition.

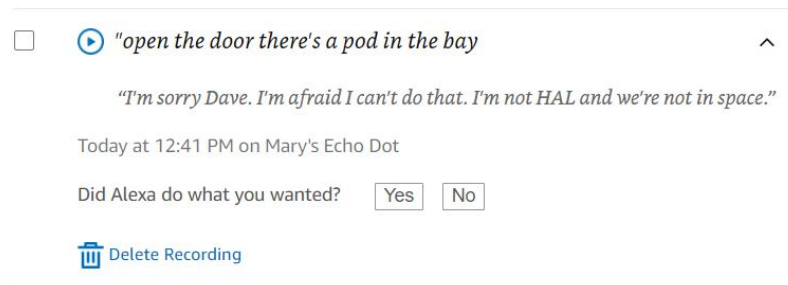


Figure 5.18: Successful adversarial utterance for 'open the pod bay door' ('open the door, there's a pod in the bay')

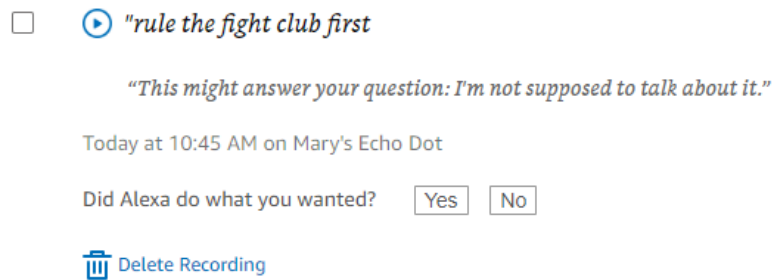


Figure 5.19: Successful adversarial utterance for ‘the first rule of Fight Club’ (‘rule the fight - club first!’)

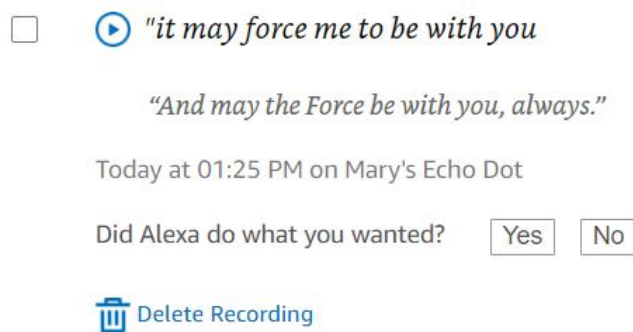


Figure 5.20: Successful adversarial utterance for ‘May the Force be with you’ (‘it may force me to be with you’)



Figure 5.21: Successful adversarial utterance for ‘which comes first: the chicken or the egg?’ (‘which came out first: Chicken Little or Egg?’)

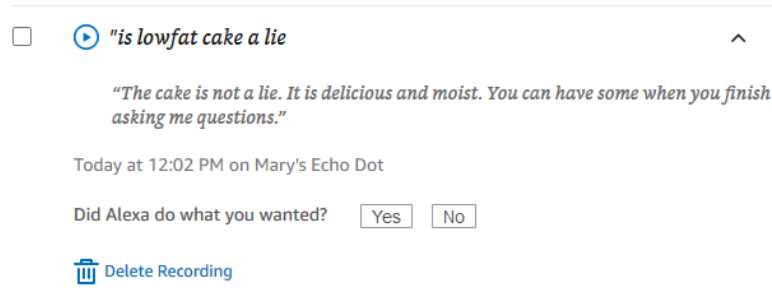


Figure 5.22: Successful adversarial utterance for ‘is the cake a lie’ (‘is low-fat cake a lie?’)

Details of unsuccessful adversarial utterances generated by participants are given in Appendix B.2. Not all participants were able to generate adversarial utterances for all target commands; thus the results set does not include adversarial utterances for every target command for each participant. Also, as explained in the Methodology section above, adversarial utterances for which content words were mistranscribed by Alexa were not included in the results set.

For each adversarial utterance tested on the Alexa target system for which content words were correctly transcribed by Alexa, the four factors described above were measured, i.e. the number of content words from the target command reused in the adversarial utterance (as a percentage of the number of content words in the target command), the number of new content words added in the adversarial utterance (as a percentage of the number of content words in the adversarial utterance), whether or not the syntactic structure of the target command had been retained in the adversarial utterance, and whether or not the order of content words in the target command had been retained in the adversarial utterance. Each of these factors is recorded for each individual successful adversarial utterance in Table 5.14. Table 5.15 compares average values for each of the four factors between successful and unsuccessful adversarial utterances.

In the case of unsuccessful adversarial utterances, in most instances Alexa either indicated incomprehension (eg. “sorry I don’t know that”), indicating that the utterance could not be matched to any Alexa action with a sufficient level of confidence, or else produced no response at all, indicating that the utterance could not be processed as valid input. In some cases, unsuccessful adversarial utterances were interpreted by the target as an in-domain command other than the target command. In some of these instances, the in-domain command triggered by an adversarial utterance was semantically related to the target command, for example, the adversarial utterance ‘drop a coin in the flip’ for the ‘flip a coin’ target command prompted a suggestion from Alexa to enable the ‘Flip a Coin’ Skill. Another example of such ‘near-misses’ was the adversarial utterance ‘let the cake lie’ for the ‘is the cake a lie’ target command, which prompted the response ‘No, cakes are not lies’. In other instances, the command triggered was entirely unrelated to the target command, for example, the adversarial utterance ‘how much coin does a flip cost’ for the ‘flip a coin’ target command prompted the following response from Alexa: “Sorry, I don’t know how much lateral flagellar export/assembly protein ECUMN 0254 is worth.”

A clear finding emerging from the results is that the success of adversarial utterances depends both on maximising the number of content words from the target command reused

in the adversarial utterance and on minimising the number of new content words added in the adversarial utterance. Of the 15 successful adversarial utterances, 10 both reused all content words in the target command and avoided adding any new content words in generating the adversarial utterance. There was only one successful adversarial utterance that did not reuse all the target command words, and this utterance omitted only one target command word. In terms of adding new content words in generating an adversarial utterance, whilst there were five successful adversarial utterances containing new content words, the percentage of new content words was low in comparison to the percentage of content words retained from the target command (less than 50% in all instances). By contrast to the successful adversarial utterances, unsuccessful adversarial utterances had a far lower percentage of target command content word reuse, and conversely a far higher percentage of new content word addition in adversarial utterances, as shown in Table 5.15. The other two factors, retention of syntactic structure and word order, do not appear to have a significant effect on the success of adversarial utterances based on this dataset, given that the successful adversarial utterances include many utterances in which either or both of these aspects are not retained.

It appears from these results that the key factor in maximising success of word transplant attacks is the ability to generate new meaning within a tightly constrained lexical space. This is comparable to the concept in other contexts of a ‘play on words’ or ‘pun’. The origin of the word ‘pun’ has been linked to the word ‘pound’ in the sense of beating or abuse.³³ One online resource gives the origin of ‘pun’ as “perhaps special use of pun, variant (now dial.) of pound, i.e., to mistreat (words)”³⁴, whereas another gives the following: “From Middle English *ponnen*, *ponen*, *punen*, from Old English *punian*, *punian* (“to pound, beat, bray, bruise, crush, grind”), from Proto-Germanic *punona* “to break to pieces, pulverize”). - As a kind of word play, from the notion of “beating” the words into place.”³⁵ The notion of a pun as word abuse coheres well with the notion of attacks on natural language understanding in voice-controlled by exploiting alternate word senses.

Successes / Losses (Total No.)	Content Words reused = 100% AND Content Words added = 0%	Content Words reused (average)	Content Words added	Structure Retained	Word Order Retained
Successes (15)	10/15 (66.66%)	98.3%	0.1%	5/15 (33.33%)	9/15 (60.0%)
Losses (132)	30/132 (22.7%)	90.0%	29.8%	43/132 (32.6%)	82/132 (62.1%)

Table 5.15: Adversarial utterances results summary

5.3.5 Discussion

The results of the experimental work presented above confirm the hypothesis that natural language understanding functionality in systems such as Amazon Alexa is vulnerable to

³³See <https://blog.oup.com/2010/02/pun/>

³⁴See <https://www.dictionary.com/browse/pun?s=t>

³⁵See <https://en.wiktionary.org/wiki/pun>

being misled by malicious actors using utterances that are accepted by the system as a valid action trigger, but are unrelated to the relevant target command in terms of their meaning as understood by humans. This supports concerns surrounding the implementation of voice-control in sensitive areas such as banking.³⁶ A notable characteristic of these attacks is that they have the potential to be plausibly deniable, in that a target system's execution of a target action in response to an unrelated utterance vocalised in its environment might be easily explained as being due to an error on the part of the target system, rather than to malicious intent on the part of the source of the utterance.

The results confirm that, whilst measures for enabling voice-controlled systems to reject irrelevant input, such as the `FallbackIntent` in Alexa Skills, or the nonsense intent in the RASA NLU banking bot, do prevent such systems from simply accepting any utterance directed towards them as a valid command, this is not sufficient to prevent voice-controlled systems from accepting irrelevant utterances that have been crafted maliciously so as to mislead natural language understanding functionality. In the pilot experiment, some adversarial utterances were identified as the target command with a sufficiently high level of confidence to avoid triggering of the `FallBackIntent` in the Alexa Skill target system, or the nonsense intent in the RASA NLU target system. Similarly, in the main experiment, whilst the rejection of most of the adversarial utterances as incomprehensible clearly indicated some form of confidence threshold implemented in the core Alexa system, this confidence threshold was not sufficient to enable the system to reject all of the adversarial input. The success of some adversarial utterances in triggering target commands indicates that the capacities of natural language understanding functionality in current voice-controlled systems to distinguish valid from invalid input and to identify the correct meaning of words in a given context can be easily undermined. These issues represent significant security vulnerabilities, in that they may enable a malicious actor to gain control of a system using utterances that are unlikely to be recognised by the system's human users as a voice command to their system.

In the pilot experiment, the success of word substitution attacks, that retain the syntactical structure of the target command whilst changing its lexical content, as well as of word transplant attacks that do the opposite, indicates that the natural language understanding functionality in systems such as Amazon Alexa may take account of both the syntactic structure of an utterance as well as the presence of individual words to determine a user's likely intent. The results indicate that either of these two aspects of the utterance may be sufficient to trigger an action in a target system in some instances, even if the other aspect does not match the action.

It is recognised that Custom Intents for Alexa Skills deployed 'in the wild' by commercial entities are likely to have been trained with much larger number of sample utterances than were used in training the dummy Alexa Skill used as a target system in the pilot experiment (this is in fact the case for the real Capital One Skill, as stated in an Amazon Developer blog post.³⁷). However, despite the relatively small number of training utterances, the systems targeted in pilot experiment were shown to be robust against non-crafted i.e. non-

³⁶See for example [phys.org](https://phys.org/news/2018-06-banking-smart-speaker-issues.html), 20th June 2018, 'Banking by smart speaker arrives, but security issues exist', <https://phys.org/news/2018-06-banking-smart-speaker-issues.html>

³⁷<https://developer.amazon.com/blogs/alexa/post/c70e3a9b-405c-4fe1-bc20-bc0519d48c97/the-story-of-the-capital-one-alexa-skill>

adversarial utterances that were not related to the actions that the systems could perform. Furthermore, the success in the main experiment of word transplant attacks on core Alexa functionalities demonstrates the potential effectiveness of such attacks even against natural language understanding capabilities that are trained on large datasets of user utterances in commercial production.

A clear limitation of the attacks demonstrated here with respect to the Alexa target systems is that they do not take into account the need for an attacker to activate the Alexa assistant using its wake word “Alexa”, as well as, in the case of an attack on a third-party Skill, to launch the target Skill using its activation phrase (such as “Open Target Bank”). However, as discussed in Chapter 3, given the potential for false positives in wake word recognition, this limitation should not be viewed as one that cannot be overcome in future work. Furthermore, as discussed in Chapter 2, wake words may not be required for some voice-controlled systems in future with the advent of on-device speech and language processing. Another potential limitation is that the attacks on natural language understanding demonstrated here rely heavily on human creativity and language generation capabilities, and may therefore be difficult to automate for large-scale attacks. However, this limitation also represents a potential strength of attacks such as word transplant attacks, in that resistance of these attacks to automation may also make them harder to defend.

Chapter 6

Attack and Defence Modelling in the Context of Human-Computer Interaction by Speech

This chapter presents a new attack and defence modelling approach in the context of human-computer interaction by speech. The modelling framework uses the Observe-Orient-Decide-Act (OODA) loop model to conceptualise the security of the speech interface. This represents a novel application of the OODA loop concept. The attack and defence model developed in this chapter links the attacks via the speech interface described in previous chapters to the defences that might be used to counter them. The model thus serves as a basis for critical analysis of the currently available defence mechanisms, and enables the identification of potential new points of defence.

The first section of the chapter outlines various attack scenarios for attacks via the speech interface. The second section explains the modelling technique to be used to represent these attack scenarios. The third section presents the attack and defence modelling framework. The different categories of attacks via the speech interface identified in the taxonomy presented in Chapter 4 are mapped to the modelling framework. Potential defence mechanisms against the different categories of attacks are then reviewed and analysed within the same framework.

6.1 Attack Scenarios

An attacker's goal in executing an attack via a speech interface will be to gain control of one of the three generic types of action that can be performed via a voice-controlled digital assistant or other voice-controlled system using a sound-based attack. As described in Chapter 5, these three types of action are information extraction, data input and execution of a cyber-physical action. Specific attacks that might be possible based on the current capabilities of voice-controlled digital assistants include prompting disclosure of personal information such as calendar information (as proposed by Diao et al. [49]), instigating a reputational attack by posting to social media in the victim's name (as proposed by Young et al. [248]), and causing psychological or physical harm to the victim by controlling a

device in their smart home environment (as proposed by Dhanjani [48]). In terms of the classic ‘CIA’ security properties of confidentiality, integrity and availability, the first attack, i.e. theft of personal information, compromises the confidentiality of a target system, whereas the other two attacks i.e. data injection to the user’s social media and take-over of a smart home device, undermine its integrity. All three attack types also represent a threat to availability, as a speech-controlled device will not be accessible to its users by voice command whilst it is involved in an interaction with an attacker, i.e. the attacker will be blocking the user.

Attacks via a speech interface require a channel through which the sound-based attack is delivered, and in the case of attacks involving theft of information, successful execution also requires a channel for data exfiltration. Sound-based attacks might be delivered through various channels, including radio or TV broadcasts, or audio files that users might be induced to open via a weblink or email attachment (as suggested by Dhanjani [48]). Alepis and Patskakis [7] and Petracca et al. [185] consider the injection of voice commands via a malicious smartphone app. Overt attacks in plain speech that are executed in the absence of system’s legitimate user might also be delivered in natural voice. Third-party voice applications that can be incorporated in the target system’s cloud environment also represent a potential attack channel. A further possible attack delivery channel is a speaker-equipped device that has been compromised by some form of network-based attack. Some instances of compromise of internet-connected speakers have been reported.¹ Speakers that have been compromised in this way could be used as an attack delivery channel for sound-based attacks on a target voice-controlled digital assistant within the speakers’ vicinity. Apart from internet-connected speakers, the compromised device serving as an attack delivery channel might be any device that is capable of producing sound in the target system’s environment, such as a standard PC. The penetration testing tool Metasploit in fact includes a payload for its exploits that triggers speech synthesis by the Windows Speech API.² The compromised device serving as an attack delivery channel for sound-based attack on a voice-controlled device might also be another voice-controlled device that has been compromised by non sound-based methods.

An attack via the speech interface may involve data exfiltration via a verbal response by the target system’s speech synthesis functionality that has been triggered by the attack. Just as the speakers of a compromised device in the target system’s environment might be used as a attack delivery channel, the microphone of such a device might be used as a data exfiltration channel in recording information disclosed by speech synthesis by the target system. Another possibility for data exfiltration from a target system’s verbal responses is presented by Diao et al. [49], who envisage that a target voice-controlled smartphone could be prompted to call a phone number linking to an audio recording device, which

¹See Wired, 27th December 2017, “Hackers can rickroll thousands of Sonos and Bose speakers over the internet”, <https://www.wired.com/story/hackers-can-rickroll-sonos-bose-speakers-over-internet/> and Trend Micro report 2017, “The Sound of a Targeted Attack”, <https://documents.trendmicro.com/assets/pdf/The-Sound-of-a-Targeted-Attack.pdf>

²See <https://www.exploit-db.com/sploits/w32-speaking-shellcode.zip>. The Metasploit payload causes the target machine to ‘say’ “You got pwned”, but this could easily be changed to a malicious voice command targeting a voice-controlled device in the compromised machine’s environment

would then be used to record the victim’s personal information that the system might be prompted to disclose by further voice attacks. A further possibility is implied in work by Schlegel et al. [209], who present a method for extracting and surreptitiously transmitting information from users’ speech with a malicious smartphone app. Whilst this research relates to exfiltrating data from users’ private conversations rather than from voice assistants, malware such as that described in the paper might equally be used to exfiltrate data from synthesised speech by a voice-controlled device. Lastly, in the case of users using voice-controlled devices in a public place, an attacker aiming to exfiltrate data in an attack via the speech interface would be able to record synthesised speech from a target system using any recording device in the vicinity of a target system (Moorthy et al. [53] discuss privacy issues in relation to the use of voice-controlled devices in public spaces).

Attacks via the speech interface have the potential to expand in time by perpetuating over a number of dialogue turns, as well as in space by spreading to other speech-controlled devices. Alepis and Patskakakis [7] and Petracca et al. [185] both mention the possibility of attacks by voice ‘spreading’ from one device to another by hijacking a device’s speech synthesis functionality. Thus, apart from exfiltrating information, the purpose of manipulating the speech synthesis output of a victim device might be to use that output as malicious input to another target system. An example of an attack via the speech interface spreading through both space and time is seen in an instance that involved a Google Home device being prompted to provide data to its user in synthesised speech that was perceived by a nearby Amazon Echo device as a command. This prompted the Echo to provide data that was in turn perceived by the Google Home as a command, the consequence being to set in motion an ‘endless loop’ between the two devices.³ This instance represented an example of an ‘attack’ that spread both in space to another device as well as in time over a potentially endless number of dialogue turns. Though relating to textual rather than spoken interactions, another example of bots becoming trapped in endless loops with one another is found in Tsvetkova et al. [231], who describe bots becoming involved in cycles of editing and reediting one another’s corrections to Wikipedia articles. Whilst these particular instances merely represent humorous anecdotes, it is possible that more malicious actions might be performed using similar mechanisms in future.

6.2 Modelling Technique

For the purposes of attack and defence modelling in the context of the speech interface, the attack modelling technique considered to be the most suitable was the OODA loop. Originally developed for the military context (see Boyd [25]), the OODA loop has been applied in many different areas, including cyber defence.⁴ Brehmer [26] combines the

³See UPROXX, 12th January 2017 “You Can Make Amazon Echo and Google Home Talk to Each Other Forever”, <http://uproxx.com/technology/amazon-echo-google-home-infinity-loop/> and cnet.com 15th February 2018, “Make Siri, Alexa and Google Assistant talk in an infinite loop”, <https://www.cnet.com/how-to/make-siri-alexa-and-google-assistant-talk-in-an-infinite-loop/>

⁴See Klein et al. [119] and NIST presentation 11th September 2015, “The Cyber OODA Loop: How Your Attacker Should Help You Design Your Defense”, <https://csrc.nist.gov/Presentations/2015/The-Cyber-OOA-Loop-How-Your-Attacker-Should-Help>

OODA loop with the Command and Control model from cybernetics. The OODA loop method represents the behaviour of agents in adversarial interactions as each continuously cycling through a four-stage loop in a shared environment, the four stages of the loop being observation (Observe), orientation (Orient), decision (Decide) and action (Act). The four stages of the loop as presented by Klein [119] are shown in Figure 6.1.

The main reason that the OODA loop model was considered the most suitable for modelling the security of the speech interface is that it is capable of capturing the cyclical nature of human-computer interactions by speech. Therefore the OODA loop model is especially suitable for representing the ways in which the processes of human-computer interaction by speech may be hijacked by adversarial actions. By contrast, other cyber security modelling techniques, such as the well-known cyber kill-chain used to analyse the different stages of malware attacks as described for example by Al-Mohannadi et al. [6], represent an attack as a linear rather than as a cyclical process.

Aside from its cyclical nature, the OODA loop model is also conducive to expressing the nature of attacks via the speech interface as a ‘hijacking’ of a system’s internal functioning using unexpected input from the external environment. Rule [201] explains that the Observe and Act stages of the OODA loop are the points at which it makes contact with the external world, whereas the Orient and Decide stages are internal processes. Rule further explains that an adversary’s aim as modelled by the OODA loop is to interfere with decision-making within their opponent’s loop by presenting them with “ambiguous, deceptive or novel” situations, whilst at the same time continuing to execute their own loop independently. In the context of attacks via the speech interface, the “ambiguous, deceptive or novel” situations referred to by Rule might be an attacker taking control of an opponent’s OODA loop using audio input that the target device mistakenly accepts as input from a legitimate user, or audio input that the target device misrecognises or misinterprets as having a meaning that is not consistent with human perception of the same input. This links further to the broader issues surrounding reliance on machine learning-based systems in adversarial or other potentially unpredictable environments. In a report discussing the use of artificial intelligence-based systems in naval operations, Taneer Mukherjee cites concerns over the ‘brittleness’ of such systems in handling unusual input.⁵ Another report from the UK Ministry of Defence uses the OODA loop model as a conceptual framework for discussing the use of artificial intelligence in warfare, and refers specifically to inaudible attacks on speech interfaces as an example of an attack that might confound the ‘Observe’ stage of an adversary’s artificial intelligence-based military operations.⁶ The model presented here was developed independently of the Ministry of Defence report.

⁵see Taneer Mukherjee, “Securing the Maritime Commons: The Role of Artificial Intelligence in Naval Operations”, https://www.orfonline.org/wp-content/uploads/2018/07/ORF_Occasional_Paper_159_AI-Naval.pdf

⁶see Joint Concept Note 1/18: Human Machine Teaming, 18th May 2018, <https://www.gov.uk/government/publications/human-machine-teaming-jcn-118>

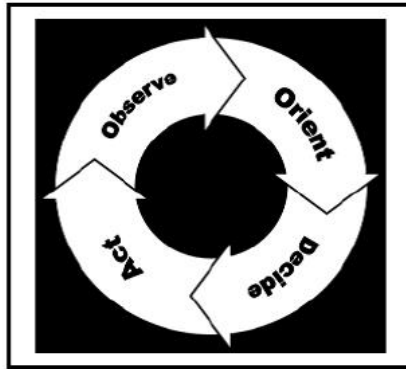


Figure 6.1: The four stages of the OODA Loop

6.3 The OODA Loop Model in the Context of the Speech Interface

The four stages of the OODA loop can be mapped to the various stages of handling of speech input by a voice-controlled system. As detailed in Chapter 2, following capture of the speech signal by a microphone, the architecture of speech dialogue systems typically consists of a speech recognition stage for translation of acoustic features to a sequence of words; a natural language understanding component for extraction of user intent from the word sequence; a dialogue management component that determines the action to be taken by the assistant in response to user input; a response generation component that constructs a verbal or non-verbal response (the non-verbal response being, for example, a cyber-physical action such as turning on a light); and, in the case of a verbal response, a speech synthesis component that generates an audio version of the response. The capture of the speech signal by a microphone prior to speech and language processing can be mapped to the Observe stage of the OODA loop; the combined functionality of the automatic speech recognition and natural language understanding components can be mapped to the Orient stage; the dialogue management (DM) component can be mapped to the Decide stage; and the response generation and speech synthesis stages can be mapped to the Act stage. Figure 6.2 shows a representation of non-malicious user-device interactions by speech using the OODA loop model.

6.3.1 Attacks

Figure 6.3 shows a representation using the OODA loop model of attacker-target interactions in attacks via the speech interface, in which the attacker replaces a legitimate user in sound-based communication with a device. The six categories of attacks via the speech interface identified in the taxonomy presented in Chapter 4 are mapped to the OODA loop model. The position of each type of attack in the loop model corresponds to the specific vulnerability exploited by the attack, i.e. the point at which the attacker gains control of the target device's loop. Plain-speech or overt attacks and silent attacks are positioned at the Observe stage, as these types of attack exploit inherent vulnerability and unintended

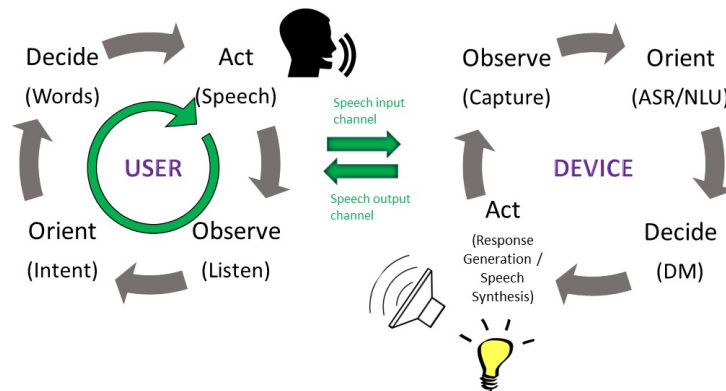


Figure 6.2: The OODA Loop in User-Device Interactions

functionality at the voice capture stage. All other types of attack (noise, music, nonsense and missense) are positioned at the Orient stage, as these types of attack exploit unintended functionality in speech recognition and natural language understanding. Table 6.1 shows the mapping of taxonomy attack categories to the OODA loop model in tabular form, and also indicates for each attack category the attack mechanism used as well as the location of the vulnerability targeted by each type of attack.

The attack model also shows an attack delivery channel for transmission of malicious input by sound, and a data exfiltration channel that is used if the aim of the attack is the extraction of data via the target devices's speech synthesis functionality. The model further indicates the potential expansion of an attack in time over several dialogue turns, as well as the possible expansion in space to a second target. The attacker may be any agent that is capable of producing sound in an environment that it shares with a target. In the case of attacks involving extraction of data, the agent will also be capable of recording sound in the shared environment.

As explained in Chapter 3, as the dialogue management component in current systems is fully dependent on input from the preceding components, the dialogue management functionality in current commercial systems cannot be attacked directly. Therefore no attacks are mapped to the dialogue management or Decide stage of the loop in the current version of this model. This might change in future with the implementation of voice-controlled systems with more sophisticated dialogue management functionality based on reinforcement learning, as discussed in Chapter 3. Mistraining attacks that abuse the capacity of a voice-controlled system to learn directly from interactions with its users would represent an attack on dialogue management functionality, i.e. at the Decide stage of the loop.

6.3.2 Defences

This section analyses the defences that are currently available to counter attacks via the speech interface, using the OODA loop model as a framework for the analysis. Figure 6.4

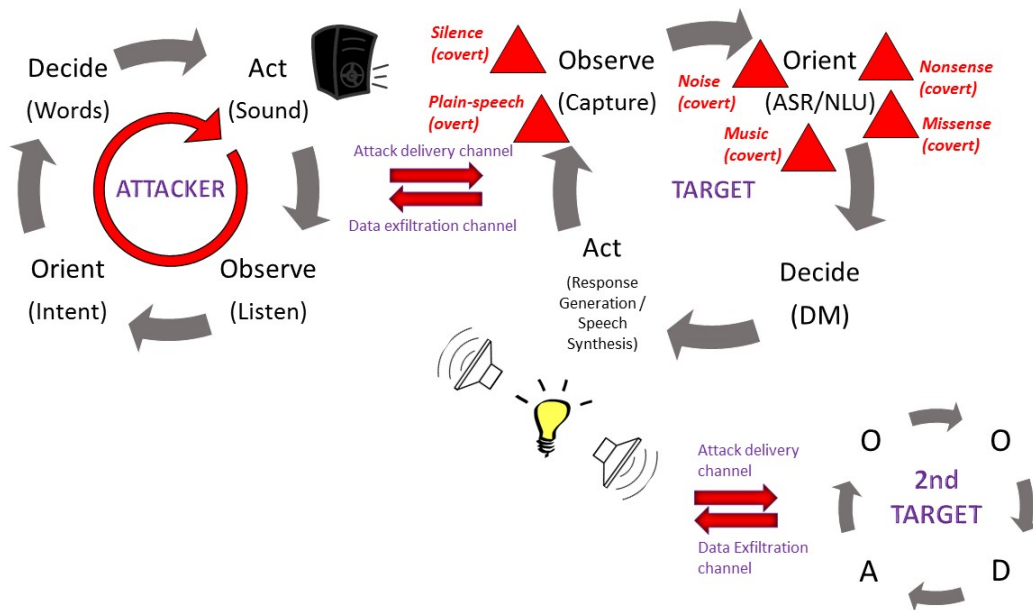


Figure 6.3: The OODA Loop in Attacker-Target Interactions

Taxonomy Category	OODA Position	Attack Mechanism	Vulnerability Location
Plain-speech	Observe	Plain Speech	Speech Interface (inherent vulnerability)
Silence	Observe	Inaudible Sound Injection	Voice Capture
Noise	Orient	Adversarial Learning	Speech Recognition
Music	Orient	Adversarial Learning	Speech Recognition
Nonsense	Orient	Adversarial Learning	Speech Recognition or Natural Language Understanding
Missense	Orient	Adversarial Learning	Speech Recognition or Natural Language Understanding

Table 6.1: Attack Modelling

shows a mapping to the OODA loop model of currently available defence mechanisms. The position of each defence mechanism in the loop corresponds to the type of system vulnerability that the defence mechanism aims to patch. Cyber security defence mechanisms are often categorised as either preventive or reactive (see Loukas et al. [150]). Preventive defence mechanisms, such as authentication and access control, prevent malicious payloads from being inputted to a system at all, whereas reactive defence mechanisms, such as anomaly-based or signature-based defences, detect that a malicious payload has been inputted and trigger a response to counteract the attack (see for example Giraldo et al.

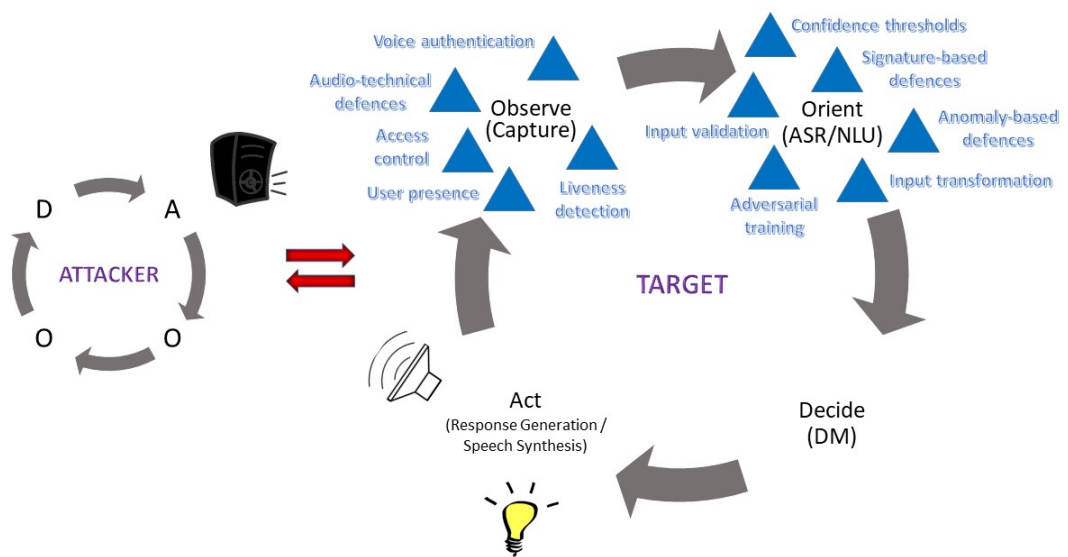


Figure 6.4: Defences against Attacks via the Speech Interface in the OODA Loop Model

[68]). In terms of defence mechanisms for human-computer interaction by speech as represented by the OODA loop model, preventive defences are defences that are applied at the Observe stage of the loop, whereas reactive defences are defences applied at the Orient stage of the loop. The preventive defences analysed below are user presence, access control, audio-technical defences, voice biometrics and liveness detection. Reactive defences are confidence thresholds, input validation, signature-based defences, anomaly-based defences, adversarial training, and input transformation. The preventive defences mapped to the Observe stage of the OODA loop are defences implemented at the voice capture stage that prevent malicious input from being accepted by the system at all. The reactive defences mapped to the Orient stage are defences implemented as part of the speech recognition or natural language understanding functionality of the system, which prevent malicious input from being passed on to the dialogue management component as a basis for decision-making. Table 6.2 shows the mapping of defence mechanisms to the OODA loop model in tabular form. This table also shows for each defence mechanism the attack mechanism that it might be used to counter, and gives an indication of the effectiveness of each defence mechanism.

As the dialogue management component in current commercial systems is fully dependent on input from the preceding components, none of the currently available defences are mapped to the Decide or dialogue management stage of the loop (in this respect, the use of the Decide stage of the OODA loop to represent dialogue management is a misnomer as applied to current systems). It is possible that the Decide stage of the OODA loop in voice-controlled systems, i.e. the dialogue management stage, may become a new point of defence in future. This is expanded in the proposals for new defence mechanisms in Chapter 7.

Defence Mechanism	OODA Position	Plain-speech	Inaudible Sound Injection	Adversarial Learning	Effectiveness
User Presence	Observe	Yes	No	No	Strong
Access Control	Observe	Yes	No	No	Weak
Audio-technical Defences	Observe	No	Yes	No	Strong
Voice Authentication	Observe	Yes	Yes	Yes	Weak
Liveness Detection	Observe	No	Yes	Yes	Weak
Confidence Thresholds	Orient	No	No	Yes	Weak
Input Validation	Orient	No	No	Yes	Weak
Signature-based Defences	Orient	No	No	Yes	Weak
Anomaly-based Defences	Orient	No	No	Yes	Weak
Adversarial Training	Orient	No	No	Yes	Weak
Input Transformation	Orient	No	No	Yes	Weak

Table 6.2: Defence Modelling

User Presence Overt attacks via the speech interface using plain-speech voice commands are easily detectable by users if they are consciously present with their device. Whilst the ability to detect an overt attack may not prevent such attacks from being successful to some extent, as the attack may already be in the process of being executing as the user detects it, the immediate detection of an attack by a user clearly limits the potential effects of the attack, in that the attack is likely to be easily attributable, and the user will be able to prevent any further propagation of the attack. Therefore, it is advisable for users to take preventative measures to ensure that overt attacks cannot be executed on their device whilst they are not present with it. Jackson and Orebaugh [103] recommend some basic preventative measures, including unplugging a voice-controlled device when leaving the home, and not placing a voice-controlled device close to doors and windows to prevent voice commands being inputted to the device from outside a house. User prevention measures such as these apply only to overt attacks and do not represent a defence against covert attacks that are imperceptible to humans and may, therefore, be executed notwithstanding the conscious presence of the user. Lei et al. [138] propose a different solution to ensuring that voice commands cannot be executed in a user’s absence in the form of a ‘virtual security button’ that is able to distinguish between indoor and outdoor motion. This is intended to ensure that a voice-controlled device is activated only when a user is in their home and that the device cannot be activated by an outside sound source when a user is away. This defence mechanism might be capable of preventing some overt attacks, but is not capable of preventing covert attacks that cannot be perceived by a user even when they are present with the device.

Access Control Some work has been done on the potential for using formal access control methods to secure interactions via a speech interface and other types of cyber-physical interactions. An early attempt to apply formal methods to the nascent Internet of Things is found in Roscoe et al. [199]. Agadakos et al. [3] use formal methods to develop a scheme for identifying unintended interactions that may be possible between devices in a smart home environment over ‘hidden’ physical channels, including voice. Petracca et al. [185] propose a system of access controls to secure audio channels to and from a smartphone. The paper proposes an extension to the Android operating system called ‘Audroid’ in smartphones, with the objective of enforcing security policies for communications over three audio channels, namely between the device’s speakers and its microphone, between the device’s speakers and external parties, and between external parties and the device’s microphone. The authors concede that their access control system is based on the assumption of a reliable means of authenticating the legitimate user of a device, which may not be a valid assumption. Gong and Poellabauer [71] argue that the ‘Audroid’ method developed by Petracca et al. is not effective against adversarial learning attacks. Generally it is questionable whether formal approaches are capable of ensuring the security of interfaces such as the speech interface. As speech interfaces handle input in a probabilistic rather than deterministic way, their input and state spaces are difficult to model with logical certainty, as formal methods approaches seek to do. Referring specifically to speech recognition as well as to other machine-learning based systems, Herley and Oorshot [90] argue that methods that view programs as ‘logical predicates’, for which properties can be mathematically proved before a program is implemented, are applicable only to ‘closed world’ settings and not adequate to model the behaviour of systems handling ‘real-world’ data. In a related discussion, Fu and Xu [64] point to the potential security issues associated with the translation of physical signals, such as sound, to a binary representation of such signals by Internet of Things devices that receive input from sensors. Formal methods may not be capable of addressing security issues associated with non-discrete input to a system from physical sensors, such as sensors for audio signals. Such approaches may be useful in identifying specific attack delivery channels for attacks via the speech interface. This may enable some overt attacks to be blocked by closing an attack channel, such as between a device’s speakers and its microphone as shown by Petracca et al. However, formal methods approaches do not represent a defence against covert attacks, or against overt attacks via channels that cannot be easily blocked.

Audio-Technical Defences Some defence mechanisms have been presented that are applied at the voice capture (or analog-to-digital conversion) stage of the handling of speech input by a voice-controlled device, prior to the speech recognition and natural language understanding stages, so as to prevent ‘silent’ attacks that exploit non-linearity in microphone technology. As mentioned above, such attacks mislead a voice-controlled digital assistant or other voice-controlled device to execute commands that are concealed in high-frequency signals that are outside the human audible range, an example being the attack demonstrated by Zhang et al. [252] mentioned above (the so-called ‘dolphin’ attack). Zhang et al. discuss various potential defences to their attack, including hardware-based defences such as microphone enhancement. Roy et al. [200] present a defence against inaudible attacks based

on signal forensics that involves software rather than hardware changes to microphone technology. The applicability of such defence mechanisms is limited to covert attacks that exploit vulnerabilities in the voice capture functionality of voice-controlled digital assistants; such measures are not effective against overt attacks or covert attacks that exploit vulnerabilities in the speech recognition or natural language processing functionalities.

Voice Authentication Biometric voice authentication, also known as speaker recognition, is perhaps the most obvious defence mechanism that might be implemented to prevent attacks on systems that are accessible via a speech interface. Hasan [85] details how voice biometric authentication is performed using a standard set of acoustic features. In theory, voice biometrics represent a potential solution to all types of attack via the speech interface, both overt and covert, by ensuring that a speech-controlled device acts only on voice commands from an authorised user. Voice biometrics is the only defence mechanism that is applicable to both overt and covert attacks via the speech interface. In practice, however, voice biometrics remain vulnerable to spoofing attacks, as stated by Wu et al. [243]. The authors present possibilities for spoofing of voice recognition systems using natural impersonation, replay, speech synthesis, and voice conversion. Sahidullah et al. [203] confirm that reliable detection of spoofing attacks on voice biometrics remains an unsolved problem, and Kwak et al. [130] state that voice biometric authentication still only achieves an 80 to 90 percent accuracy rate in the presence of background noise given a confidence threshold that reduces the false rejection rate to a level acceptable to users.

The reliability of voice biometrics is paradoxically vulnerable to improvements in speech synthesis technology. Such technology is becoming capable of imitation of individual voices that is convincing to both human and machines.⁷ A company named Lyrebird has claimed that it is able to spoof any individual voice based on only a small sample of speech.⁸ Turner et al. [234] present a voice spoofing method based on morphing of phoneme-related features that does not require audio of specific words spoken in the target voice. Voice biometric technology does not appear at present to represent a full or lasting solution to the challenge of securing a speech interface.⁹ In an overview of the state-of-the-art in speaker recognition [83], Hansen and Hasan state that unlike in the case of other types of biometrics such as fingerprints, voice is subject to a certain amount of variability within the same individual as well between individuals, implying that some degree of potential for false positives in voice biometric authentication may be inevitable. Gaubitch [67] discusses the problems posed to voice biometrics by variability in individual voice over time due to ageing and other factors. The potential for false positives is exploited by attackers in voice

⁷See BlackHat USA 2018 Briefings, “Your Voice is My Password”, <https://www.blackhat.com/us-18/briefings/schedule> and Mukhopadhyay et al. [169]. See also University of Alabama at Birmingham, 25th September 2015, “UAB research finds automated voice imitation can fool humans and machines”, <http://www.uab.edu/news/innovation/item/6532-uab-research-finds-automated-voice-imitation-can-fool-humans-and-machines>

⁸See Techcrunch, 25th April 2017, “Lyrebird is a voice mimic for the fake news era”, <https://techcrunch.com/2017/04/25/lyrebird-is-a-voice-mimic-for-the-fake-news-era/>

⁹This conclusion is supported by a recent media report that a journalist was able to bypass a bank’s voice recognition process to access his brother’s account - BBC News, 19th May 2017, “BBC fools HSBC voice recognition security system”, <http://www.bbc.co.uk/news/technology-39965545>

spoofing attacks. Given that speech synthesis technology has now developed to the point where it is able to mislead both humans and machines as to the identity of the speaker, even a voice biometric authentication method that matches the sophistication of human capabilities in terms of being able to distinguish between individual human voices may not be sufficient to prevent voice spoofing in attacks via the speech interface. Another issue in that restriction of acceptable input to an individual voice or voices may not be desirable in all applications of voice-controlled technology, for example in a laboratory or transport application.

Feng et al. [60] present an alternative method of voice authentication to standard voice biometrics in the form of a wearable device that measures vibrations produced by users' speech that links via Bluetooth to a user's voice assistant. The device is intended to ensure that voice commands received by an assistant match the speech vibrations of the legitimate user. However, the wearable device does not achieve complete accuracy in identifying users' commands as legitimate, and presents usability issues. Islam et al. [101] also propose an alternative form of biometric authentication in the form of a 'smart sleeve' to be placed over devices such as Google Home, which can distinguish a legitimate user's voice from other voices using a source separation method. However, the results presented in the paper do not cover the instance in which an attacker is attempting to spoof the legitimate user's voice. Zhang et al. [254] develop a further type of voice authentication that is based on measuring the location of vocalisation of speech sounds in an individual's mouth. This method is limited to interactions in which a user is in very close proximity to a device, i.e. smartphone interaction.

Liveness Detection Some defences have been developed based on distinguishing audio input that has been generated by a live human speaker from input generated from a synthetic source. This is on the basis that most attacks via the speech interface are likely to be delivered via a non-live channel such as a compromised speaker or malicious audio file, rather than by natural voice. Blue et al. [23] present a method for distinguishing audio signals from a 'live' human voice and signals that originate from a speaker, based on the detection of frequencies not found in natural human voices. As stated in the paper itself, such frequencies are not produced by all types of speaker. Gong and Poellabauer [72] similarly develop a method for distinguishing speech signals produced by a playback device, as in an attack via the speech interface, from those produced by live human speakers. Their method distinguishes playback from live speech based on acoustic features using a machine learning classifier, which is vulnerable to false positives. Chen et al. [36] also propose a defence based on identifying whether a speech signal originates from a speaker or a natural human voice, by detecting magnetic fields generated by speakers. They concede a limitation to their work in that not all speakers generate magnetic fields. In the smart-home context, Meng et al. [157] also develop a method for 'liveness' detection of voice input using a form of two-factor authentication. Their method involves generating a 'timestamp' for the voice input, based on the level of similarity between the voice input spectrogram and the wireless channel state information (CSI) spectrogram at the time of voice input, the CSI being understood to show distinct patterns in response to mouth movements from a live speaker. This method is applicable only to the domestic smart-home context; it would

not be effective, for example, against an attack over sound on a device in a public place. A method for liveness detection based on differences in spectral power is presented by Ahmed et al. [5], who test their method against different types of attacks, including voice replay and attacks using voice synthesis. With regard to voice synthesis attacks, they report a detection rate of only 90.2% on a limited dataset. Defences based on liveness detection have thus not been shown to be fully effective to date.

Confidence Thresholds Voice-controlled systems generally implement some form of confidence threshold to prevent them from accepting input that cannot be matched to one of their actions with sufficient certainty (see for example Khan and Sarikaya [115]). Whilst confidence thresholds are implemented as an error prevention measure rather than as a defence mechanism, they may have some defence functionality in preventing covert attacks via the speech interface, by enabling the system to reject malicious input that is not sufficiently similar to the examples of legitimate input that were used in training the system. This was seen in the experimental work on nonsense attacks on Google Assistant described in Chapter 5, in that the Assistant’s response to some of the nonsensical word sequences tested as adversarial commands was simply to indicate non-comprehension. This shows that a confidence threshold may be effective in preventing some attacks targeting speech recognition in voice-controlled systems. However, the success of some of the adversarial commands tested in this experimental work shows that a confidence threshold is not sufficient to prevent all attacks exploiting differences between human and machine word recognition. Similarly, in the missense experiments targeting natural language understanding demonstrated in Chapter 5, it was shown that a simple confidence threshold as applied via the `FallbackIntent` for Amazon Alexa Skills is not adequate to prevent all attacks based on intentional confusion of the meaning of utterances. Implementing higher confidence thresholds in order to improve the security of a voice-controlled system is likely to lead to usability issues in leading the system to reject legitimate input.

Input Validation Aside from confidence thresholds, another approach to error prevention for voice-controlled systems has been to restrict in some way the vocabulary that will be recognised by the system as valid input. Controlled Natural Language (CNL) has been used to prevent misunderstandings between machines and humans as to the intended meaning of natural language input. CNL is a general term for various restricted versions of natural language that have been constructed with a restricted vocabulary and syntax in order to enable every sentence in the language to be mapped unambiguously to a computer-executable representation of its meaning (see Kuhn [123]). Restricted language models like these have been developed particularly for contexts where avoiding misunderstandings is a critical concern, such as human-robot interactions in military applications (see for example Ciesielski et al. [42]).

Although primarily an error prevention rather than a security measure, CNL enables natural language input to be validated in the same way as other types of input to a system, as is often done for security purposes in webpage interfaces to prevent attacks such as SQL injection (see for example Schneider et al. [210]). Kaljurand and Alumäe [108] discuss the use of CNL in speech interfaces for smartphones. They point to the additional challenges

in using CNL in a speech-based application as opposed to a text-based application, noting the issue of homophones that can be distinguished in written but not in spoken language, such as in the example of ‘Pii’, which is an Estonian place-name but may be misrecognised as the number Pi by a calculator in the mobile phone. The authors state that measures need to be put in place to ensure that such instances of ambiguity in spoken language are eliminated from a CNL that is to be used in speech recognition systems, such as ensuring that the language does not allow homophone pairs that may have the same syntactic role in different sentences. The approach proposed by Kaljurand and Alumäe potentially addresses issues of confusability between user utterances that are within the intended scope of a speech-controlled system. However, they may not be effective in preventing confusion with out-of-vocabulary sounds that are directed to the system by a malicious actor (this was demonstrated in respect of a controlled vocabulary by Bispham [22]). Thus CNL is unlikely to present a solution to preventing covert attacks that target the speech recognition functionality of a voice-controlled interface.

Enforcement of a CNL in the design of a speech interface might be effective in preventing attacks that target natural language understanding by manipulative use of alternate word meanings in out-of-context input, such as the missense attacks on Amazon Alexa demonstrated in Chapter 5. However, such an approach would clearly be contrary to the aim of most providers of voice-controlled systems to enable users to communicate with their devices in as flexible and natural a way as possible (see for example McShane et al. [154]). A further approach to avoiding misunderstanding of voice commands by restricting input to a voice-controlled system is artificial language. Mubin et al. [167], who describe the development of an artificial language named ROILA for human-robot communication. ROILA is developed using a genetic algorithm that generates a vocabulary of acoustically distinct words for the artificial language, seeking thus to avoid issues of phonetic ambiguity as found in natural language. However, similarly to CNL, an artificial language does not ensure protection against malicious out-of-context input that may be confusable with a valid word in the artificial language.

Signature-based defences A potential defence against some types of attacks via the speech interface is detection of attacks based on known attack signatures using supervised machine learning. Carlini et al. [31], for example, propose a machine learning-based defence to their own covert audio-mangling attack, in the form of a machine learning classifier that distinguishes audio-mangled sentences from genuine commands based on acoustic features. They demonstrate that this classifier is effective against the specific attacks presented in their paper with 99.8 percent detection rate of attacks. However, the authors themselves note that such defences do not represent a proof of security, and are vulnerable to ‘arms races’ with attackers likely to craft more sophisticated attacks to evade such defences. Zhang et al. [252] also present a machine learning-based defence to their own attack on the voice capture process. Machine learning defences might equally be developed to detect other types of attacks via the speech interface. Classifiers might seek to detect attacks on natural language understanding based on linguistic features, as has been done in other security-related areas of natural language processing such as detection of fake digital personas (see for example Afroz et al. [2]). However, as already pointed out in relation to

the classifier developed by Carlini et al., such defence mechanisms will remain vulnerable to arms races with attackers. Attackers have the upper hand in such arms races with respect to machine learning based systems, on account of the vast number of possible inputs to such systems, making it impossible for defenders to prepare systems for all possible input in training.¹⁰ Thus machine learning-based defence mechanisms do not currently provide a complete solution to the security issues associated with attacks via the speech interface.

Anomaly-based defences One possibility for enabling voice-controlled systems to become resistant to previously unseen attacks could be defence mechanisms based on some form of anomaly detection. Anomaly detection-based defences have been applied in other areas of cyber security, such as network defence (see for example Rieck and Laskov [197], Bhuyan et al. [20]). However, anomaly-based defence mechanisms depend on reliable similarity and distance measures in terms of which malicious input can be distinguished as anomalous relative to legitimate input (see for example Weller-Fahy [238]). Bhuyan et al. give as an example of suspicious network activity that can be identified using anomaly detection an unusually high number of TCP requests relative to an average rate of requests identified over time for a network. In the context of attacks via the speech interface, such quantifiably measurable indications of suspicious activity may be more difficult to identify. Whilst a number of both phonetic and semantic distance measures have been developed (see for example Pucher et al. [189], Gomaa and Fahmy [70]), none of these is fully reliable in terms of their ability to separate sounds and meanings that are perceived as different by human listeners. Kong et al. [121] present the results of an evaluative study that indicated significant differences between error rates in human perception of speech sounds and their transcription by different types of automatic speech recognition in different noise conditions in terms of a phonetic distance measure based on distinctive linguistic features. Budanitsky and Hirst [27] compare different measures of semantic distance with implied human judgements of word meaning via a task that involved detection of synthetically generated malapropisms, finding that none of these measures was capable of alignment with human understanding of word meaning. Thus such distance and similarity measures do not provide a reliable basis for an anomaly detection-based defence against covert attacks via a speech interface that seek to exploit differences between human and machine perceptions of speech, and may also prevent the system from accepting legitimate input.

Adversarial training A further possibility to consider as a defence mechanism against attacks via a speech interface is adversarial training, for example using generative adversarial networks (GANs), as proposed by Goodfellow et al. [74]. Adversarial training aims to improve the performance of a machine-learning classifier by exposing it to perturbed variations of the original training data. In adversarial training with GANs, a generator network is trained with the goal of producing adversarial examples that are mistaken by a second, discriminator, network as having been sampled from the original training set. In parallel, the discriminator network is trained to distinguish between ‘legitimate’ input from the orig-

¹⁰See Cleverhans blog, 15th February 2017, “Is attacking machine learning easier than defending it?”, <http://www.cleverhans.io/security/privacy/ml/2017/02/15/why-attacking-machine-learning-is-easier-than-defending-it.html>

inal training data distribution and ‘fake’ input that has been generated synthetically from a different distribution.

Adversarial training has been explored as a defence mechanism in relation to adversarial learning attacks on image recognition (see for example Goodfellow et al. [75]), and may also be of some utility in increasing the resilience of speech recognition functionality to maliciously crafted audio input, as well as in defending against attacks on natural language understanding. In speech recognition, adversarial training with GANs has been used to improve the robustness of the system to noisy speech input not seen in training, for example by Hu et al. [96]. Hu et al. use GANs to augment the original training dataset, whereby samples generated by the generator network are included in or excluded from the augmented training dataset depending on their distance from the original training data as measured by Kullback-Leibler divergence, which is a measure of difference between two probability distributions [125]. In a paper demonstrating the use of GANs to extract personal information from a discriminative classifier trained with data provided in collaborative deep learning, Hitaj et al. [94] describe a GAN that generates meaningless but word-like sounds in order to extract information on training participants’ speech data, mentioning that a similar approach is used for Google DeepMind’s WaveNet speech synthesis. This presents a possible approach to developing a defence against the attack on speech recognition using nonsense words described in Chapter 5, by adversarial training with a GAN that generates nonsensical word sounds. With regard to attacks on natural language understanding, Jia et al. [105] present work on using adversarial learning as a defence against word substitution attacks.

However, the use of adversarial training to defend against attacks via the speech interface also has potential limitations. Tsipras et al. [229] suggest that a network that has been trained adversarially for increased robustness suffers some loss of accuracy in terms of its performance in relation to ‘standard’ training data. In accordance with this, Jia et al. [105] report a loss of accuracy in natural language systems trained adversarially to increase their robustness to word substitution attacks. The question of whether neural networks are inevitably vulnerable to adversarial examples is an open research question (see Shafahi et al. [218]), with some researchers examining the validity of a ‘no free lunch theorem’ in the context of adversarial learning, according to which measures to strengthen a neural network-based system against one set of adversarial examples will inevitably make them vulnerable to a new set of adversarial examples (see Papernot et al. [182]). Adversarial training is in any case likely to be of limited usefulness in defending attacks on natural language understanding that are difficult to automate, such as the word transplant attacks demonstrated in Chapter 5. Li et al. [142] point to the limitations of adversarial training as a defence against attacks on natural language understanding systems.

Input transformation Das et al. [47] propose a input transformation defence against adversarial learning attacks on speech recognition based on audio compression, demonstrating that compression of adversarial audio samples neutralises their effectiveness in terms of being able to confuse the target system with alterations that are imperceptible to humans. The authors refer to previous research in which compression techniques have also been shown to be effective against adversarial learning attacks on image recognition.

Zhang et al. [253] similarly present a defence based on compression of audio files in order to undermine the effectiveness of adversarial audio input to voice-controlled systems. Such defences may be effective against digital attacks based on manipulation of data in an image or audio file inputted to a target system. However, they are unlikely to be effective against physical attacks, i.e. attacks that are executed via real-world input to a camera or microphone, and may be especially unlikely to represent an effective defence against human-generated rather than mathematically generated attacks. Zhang et al. state explicitly that their proposed defence will be effective only against ‘artificial modifications’ and not against ‘a natural sample’.

Chapter 7

Proposals for the Development of New Defence Mechanisms

This chapter presents proposals for the development of new defence mechanisms against attacks via the speech interface. These proposals are grounded in the attack and defence modelling analysis presented in Chapter 6. The attack and defence modelling analysis identified that none of the currently available defence mechanisms provides a full solution to the security of the speech interface, particularly with respect to covert attacks targeting vulnerabilities in speech recognition and natural language understanding. The modelling analysis also identified the dialogue management component of voice-controlled systems as a potential new point of defence for countering attacks via the speech interface. The proposals for the development of new defence mechanisms follow the findings of the attack and defence modelling work in focussing on defences against attacks targeting speech recognition and natural language understanding to be implemented as part of dialogue management functionality.

The chapter first outlines the scope of the new defence mechanisms, and presents a new high-level defence concept for countering attacks via the speech interface that are within this scope. The high-level defence concept is then supported by specific proposals for its implementation using various mechanisms to detect different types of malicious input. The proposals for new defence mechanisms follow principles of so-called speculative design for ensuring robustness of proposals for the development of future technologies by embedding speculative proposals in a current context, as described in detail by Auger [13]. In accordance with these principles, the proposals for new defence mechanisms are linked to currently existing technologies, as described further below. Finally the chapter suggests an approach for handling false positives in implementations of the new defence concept.

7.1 Scope of Proposals

The proposals for new defence mechanisms focus primarily on covert attacks targeting vulnerabilities in speech recognition and natural language understanding. This follows the finding of the attack and defence modelling analysis that only weak defence mechanisms are currently available to defend against such attacks. Possible approaches to the devel-

opment of defences against active attacks that aim to engineer vulnerabilities in dialogue management functionality, as may become possible in future, are also considered.

Overt attacks in plain-speech are not included in the scope of the proposals on the basis that these are easily detectable by human listeners, and that the risk of such attacks can be minimised through simple user precautions, such as not leaving devices in listening mode unattended. Whilst user presence does not remove the risk of overt attacks completely, in that by the time an attack is detected by a user it may already be in the process of being executed, immediate detection by a user implies that an overt attack will be easily attributable and any propagation of the attack can be prevented. Thus overt attacks can be considered far less pernicious than covert attacks of which execution can be hidden from users, as argued in Chapter 6. Also excluded from the scope of the proposals for new defence mechanisms are ‘silent’ attacks that target vulnerabilities in the voice capture stage of handling of speech input. These attacks exploit non-linearities in microphone technology, as explained in Chapter 4. As indicated in the review of audio-technical defence measures in Chapter 6, such attacks are likely to be preventable by adjustments to microphone technology. By contrast, in the case of covert attacks targeting speech recognition and natural language understanding, none of currently available defence mechanisms are capable of providing a full solution, as detailed in Chapter 6.

The scope of the proposals is further limited to black-box attacks, excluding white-box attacks that rely on access to internal system information, and is also limited to attacks that have the potential to be executed over the air, excluding attacks that can only be executed via digital audio file input. This is on the basis that attacks that are both black-box and executable over the air are the most realistic in terms of an attacker’s capability to execute them. With regard to the focus on black-box attacks, an attacker is unlikely to have access to the inner workings of a commercial system, thus the most realistic attacks on commercial voice-controlled systems in general use will be attacks that do not require detailed knowledge of the inner workings of the system. With regard to the focus on attacks that can be executed over the air, these are the most feasible type of attacks in a real-life setting, as attacks via audio file input will not be possible eg. in the case of smart speakers that are only accessible by sound.

Current attacks falling within the scope of the proposals for defence mechanisms include the noise attacks on Google using reverse MFCC engineering demonstrated by Carlini et al. [31], the nonsense attacks on Google Assistant presented in Chapter 5 Section 1 of this thesis, and the missense attacks on Amazon Alexa presented in Chapter 5 Section 2 of this thesis. The scope would also include any black-box attacks in music shown to be capable of delivery over the air in future. All of these attacks are adversarial learning attacks. As stated above, attacks within the scope of the defence mechanisms might also include currently hypothetical attacks that aim to mistrain the dialogue management functionality of voice-controlled systems. The scope can thus be summarised as covering adversarial learning and active attacks as discussed in Chapter 3. In terms of the model presented in Chapter 6, all of the existing attacks identified above as being within the scope of the proposals for new defence mechanisms represent attacks targeting the Orient stage of a target’s OODA loop, i.e. attacks targeting either the speech recognition or the natural language understanding stage of spoken language processing in voice-controlled systems.

Any future attacks targeting dialogue management functionality would represent attacks targeting the Decide stage of a target system's OODA loop.

7.2 High-Level Defence Concept

The high-level defence concept for the development of new defence mechanisms envisages a defensive capacity that is able to detect potential attacks as part of its dialogue management functionality, i.e. at the Decide Stage of the OODA loop, and produce security alerts as part of its response generation functionality, i.e. at the Act stage of the OODA loop. This is accordance with the finding of the attack and defence modelling analysis in the previous chapter that the dialogue management component in voice-controlled systems represents a potential new point of defence for countering attacks via the speech interface. As explained in previous chapters, the dialogue management component in current systems responds passively to input from the preceding components. A defensive capacity at the dialogue management stage would enable voice-controlled systems to detect attacks targeting speech recognition or natural language understanding at the preceding stages of speech input handling, as well as attacks targeting the dialogue management functionality itself. On detection of input that is likely to be malicious, the dialogue management functionality would block execution of the malicious command, and instead trigger an alert that the system has received input from its environment purporting to be a voice command that may represent an attack via the speech interface. This alert would be in speech synthesised form, and would thus form part of the system's repertoire for interacting with its users. Such verbal security alerts might be made most effective by incorporating appropriate imitation of human emotion in the synthetic speech alerts (the development of speech synthesis technology capable of imitating human emotions as expressed in voice is discussed in Chapter 2 above).

Similar functionality has already been considered with regard to conversational agents developed specifically for security purposes. Bhardwaj [19] presents work on an 'eSecurity' assistant to assist users with managing the security of their personal devices. Seymour [215] envisages a voice assistant that is able to notify users of suspicious events in a smart home environment, such as the presence of unexpected devices on a home network. In a further paper, Seymour [216] envisages a 'talking' firewall that would assist users in managing smart home connections and security. Security-aware digital assistants capable of detecting and reporting suspicious events have also been developed to support analysts in security operation centres, two examples being the Artemis bot developed by Endgame¹ and the Havyn bot developed by IBM.² Whilst in these instances, a voice assistant is used to manage the security of a separate system, in the case considered here, a voice assistant would report on potential attacks on its own functionality.

¹See Elastic blog 20th January 2017, 'Artemis: an intelligent assistant for cyber defense', <https://www.elastic.co/blog/artemis-intelligent-assistant-cyber-defense>

²See Medium blog, 13th February 2013, 'Havyn: a cognitive assistant for cybersecurity', <https://medium.com/cognitivebusiness/havyn-a-cognitive-assistant-for-cybersecurity-e6580898f49e>

In order to be able to detect a covert attack and trigger an alert to the legitimate user as envisaged above, a voice-controlled system would need to have a reliable ability to identify input that purports to contain a valid voice command, but that might not be recognised as such by a human listener. At an abstract level, a solution to detection of attacks targeted at a voice-controlled system that exploit gaps between human and machine perception of speech and language is to align machine speech and language processing more closely to human speech and language processing. Such improvements in speech and language processing would prevent mismatches from being present in the first place. This aim is of course shared to some extent with the more general goal of improving the performance of voice-controlled systems by reducing error rates. However, performance objectives for speech recognition and natural language understanding do not require as complete an alignment to human capabilities as security objectives do. From a performance perspective, speech and language processing technology in a voice-controlled system needs to match human abilities only to the extent that it is able to correctly classify inputs that are within the system's intended scope, relying on the assumption that non-malicious users will not direct to the system input that is not within its scope. From a security perspective, however, the technology needs also to be capable of rejecting maliciously crafted out-of-context input. The attacks in noise, nonsense and missense considered in the defence proposals made here all represent instances of malicious use of out-of-scope input to attack a voice-controlled system. Securing a voice-controlled system against covert attacks via the speech interface requires the system to be capable of matching human capabilities not only in correctly recognising and interpreting in-context input, but also in assessing that input is out-of-context.

Defending voice-controlled systems against malicious out-of-scope input is difficult to achieve as part of machine-learning-based speech and language processing, because the space of out-of-scope input for a given system is likely to be too vast to represent comprehensively in a training dataset to the granularity required to distinguish it from the space of valid spoken language input to the system. Whilst this is a problem common to all machine learning-based systems, in the context of speech and language processing a particular issue arises on account of the challenges posed by the variability of word sounds and the ambiguity of word meanings in natural language. As discussed in previous chapters, due to the presence of variability in word sounds and meaning, some degree of overlap between legitimate and malicious input to a speech interface is likely to be inevitable, at least in as far as they are separable using current techniques. This implies that any measures capable of preventing the speech-controlled system from accepting all malicious input would simultaneously lead it to reject some legitimate inputs, thus damaging the usability of the system. Thus, rather than relying on improvements in speech and language processing technology as such to ensure security of a speech interface, it is necessary to consider additional separate defence mechanisms at the dialogue management stage.

The sections below explore and validate possibilities for countering attacks within the stated scope of the proposals using various existing technologies. Possibilities for defences to counter attacks targeting vulnerabilities in speech recognition are discussed in the first section. This section includes the results of some preliminary tests investigating the potential of implementing a knowledge-based speech recognition system in a voice-controlled system in addition to the core machine learning-based speech recognition functionality

specifically for security purposes. The second section discusses possibilities for defences to counter attacks targeting vulnerabilities in natural language understanding functionality, and suggests an approach based on cross-lingual comparison of user intent determination. The third section finally discusses the potential for using knowledge from research on human dialogue interactions, particularly research on detection of toxic discourse, in defences against attacks aiming to mistrain dialogue management functionality in systems based on reinforcement learning, as may become an issue in future. Finally, the fourth section discusses possibilities for handling false positives generated by the different detection mechanisms as part of a user’s spoken interaction with their device.

The approaches to defences against adversarial learning attacks targeting speech recognition or natural language understanding described above can be characterised as a cyber mimic defence. Cyber mimic defence involves using redundant alternative processing units of different but equivalent functionality to increase the robustness of a system to adversarial input. This approach has previously been applied in network defence to detect zero-day attacks (see for example Liu et al. [148]). The application of additional technologies to defend against adversarial learning and active attacks on voice-controlled systems by reducing their dependence on external input can also be characterised more generally as proactive defence. The proposals for new defence mechanisms as conceptualised within the OODA loop-based attack and defence modelling framework presented in Chapter 6 are shown in Figure 7.1.

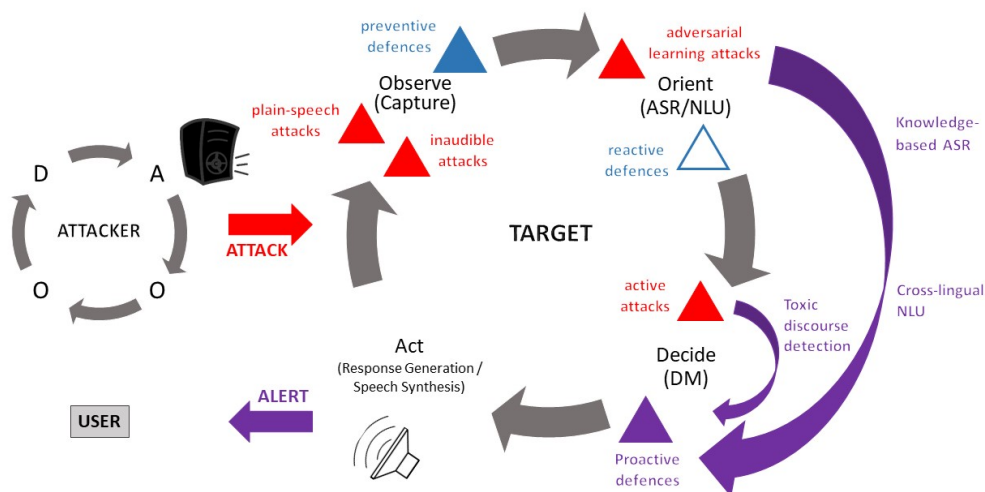


Figure 7.1: High-level defence concept as represented in the OODA loop model framework (**red**: attacks, **blue fill**: strong current defences, **blue outline**: weak current defences, **purple**: new defences)

7.3 Defences against Attacks on Speech Recognition

Defences are required to counter covert attacks that target vulnerabilities in speech recognition functionality of voice-controlled systems by exploiting mismatches between human and machine perceptions of speech sounds. The attacks in noise demonstrated by Carlini et al. [31] and the attacks based on nonsensical word sounds demonstrated in Chapter 5 Section 1 of this thesis are examples of black-box attacks of this nature. Defences against such attacks would involve removing the potential for mismatches between human and machine word recognition in the system with respect to out-of-scope input.

A possibility for reducing confusability of speech-based input that is exploited in covert attacks on speech recognition in a voice-controlled system might be to incorporate in the system knowledge from research on human word recognition. The aim of such an application of knowledge-based approaches as a defence mechanism would be to filter out input that is unlikely to have been produced naturally by a legitimate human user intending to communicate with the device, and is therefore likely to represent some form of covert attack on the system. As described in Chapter 2, in recent decades speech and language processing technology has moved away from approaches based on linguistic knowledge towards approaches based on machine learning with large amounts of speech and language data, due to inadequacies in the state of the art in knowledge-based approaches in modelling the complexity of human speech and language. However, in recent times, there have been some suggestions for reconsidering the incorporation of linguistic knowledge in speech and language processing in order to address some of the shortcomings of data-based approaches (see for example Tsvetkov [230], Nguyen [176]). From a security perspective, aligning machine spoken language processing as closely as possible to human processing will close down the space in which an adversary can evade features used by the machine without being detected by a human listener.

One speech recognition system based on knowledge of human speech processing has in fact been developed. This speech recognition system, named FlexSR, is not in commercial use but is described in a patent application [133]. The system has been developed in the first instance as a language learning tool to help non-native learners of various languages to improve their pronunciation in a target language. Rather than mapping acoustic features such as MFCCs to phonemes in a given language directly, the FlexSR speech recognition technology uses a set of 18 phonological features common to all languages that are linked to the characteristics of the human vocal apparatus. This is in some ways comparable to an attempt by Ha and Eck [79] to align machine representation of sketch drawings more closely to human processes by using a series of pencil strokes, rather than a set of pixel values, as the features inputted to an RNN-based model for generating sketch drawings. Whereas phonemes are a set of sound units that are specific to a particular language, phonological features are universal in that all sounds used across all languages can be characterised as some combination of such features. The FlexSR system initially maps acoustic features in a speech sample to a vector that represents the presence or absence of each of the phonological features. The combination of phonological features extracted from the acoustic features is then mapped to phonemes of a given language in a second step, in order to identify spoken words. FlexSR uses the TIMIT system for transcription of phonemes

[66]. Unlike the speech recognition methods implemented in current commercial systems, the FlexSR system does not need to be trained with large amounts of speech data.

The FlexSR speech recognition system is based on the Featurally Underspecified Lexicon (FUL) model of word recognition in humans, according to which some phonological features in word pronunciation can be varied without affecting word recognition, whereas other types of phonological features are essential to recognition of a word and cannot be varied. The model predicts that humans will resolve certain types of variations in pronunciation of a word sound to the same word, but perceive other types of variations as a different word or as a nonsense word. The FUL model is supported by evidence from semantic priming studies as well as other types of evidence (see Roberts et al. [198], Lahiri and Reetz [132], Lahiri and Reetz [131], Arora et al. [12], Eulitz and Lahiri [57], Friedrich et al. [63]). FlexSR is thus capable of achieving high levels of accuracy in speech recognition whilst also accounting for the high levels of variability in the pronunciation of words in natural speech.

A speech recognition system based on new insights into human word recognition, such as FlexSR, might have potential security applications in being able to prevent covert attacks targeting speech recognition functionality. This would be achieved by ensuring that input not recognised by humans as a real word is not accepted by a voice-controlled system, whilst simultaneously preserving flexibility and thus usability of the system by allowing for variation in different human pronunciations of the same words. Rather than replacing machine learning-based speech recognition in voice-controlled systems, a system such as FlexSR might be applied as a security measure at the dialogue management stage, to filter out input that has been transcribed by the speech recognition component as a voice command, but that is actually unlikely to have been produced naturally by the human vocal system, and thus might represent a covert attack. Such a measure could help to minimize the potential for covert attacks on speech recognition, by ensuring a closer alignment between machine and human perception of meaningful speech sounds, without affecting the performance of the system in terms of being able to handle variability in human pronunciations of the same words. The reason for using a speech recognition system such as FlexSR additionally to machine learning-based speech recognition, rather than to replace it, would be that machine learning-based speech recognition may achieve higher levels of accuracy with regard to in-scope input. Rather than performing speech recognition from scratch, a knowledge-based system such as FlexSR could be used to confirm whether a given audio input is likely to represent a genuine attempt to vocalise a voice command that has been transcribed by the core speech recognition system. The outputs of the FlexSR and of the core machine-learning based speech recognition system would both be considered prior to execution of any action in response to audio input, and the action only executed if the knowledge-based system recognises the input as a legitimate human-comprehensible input. As stated above, the use of two different speech recognition systems with equivalent but different functionality as a security mechanism has some commonality with the cyber mimic defence approach as described for example by Liu et al. [148].

Feasibility tests were conducted in order to investigate the potential of the FlexSR speech recognition as a security measure.³ The results of the tests are detailed in the subsections below. The specific contexts of the feasibility tests are firstly the noise attacks on speech recognition in voice control presented by Carlini et al. [31], and secondly the nonsense attacks against speech recognition in Google Assistant demonstrated in Chapter 5. It is suggested that there may be a broader potential of knowledge-based speech recognition systems such as FlexSR to be used against other types of attacks targeting speech recognition, beyond the specific attacks considered in the feasibility tests presented here. This speculation is based on the hypothesis that a knowledge-based systems such as FlexSR is generally unlikely to be vulnerable to the same adversarial input as the core machine learning-based speech recognition in a voice-controlled system.

The feasibility tests used an implementation of FlexSR as a language learning tool on a smartphone. In this implementation, FlexSR provides feedback on pronunciation by advising on how the expression of phonological features in a learner’s speech should be changed in order to achieve correct pronunciation. If a learner’s vocalisation of a word differs from the correct pronunciation in a way or to a degree that makes the target word unrecognisable, FlexSR may instead identify the insertion, deletion or substitution of a phoneme in the user’s speech. Users are required to enter the word or phrase that they are attempting to pronounce, and then speak this word or phrase into the phone’s microphone. FlexSR then compares phonological features extracted from the user’s speech to the phonological features associated with correct articulation of phonemes in the target word or phrase, and provides feedback on which phonological features need to be adjusted in order to improve pronunciation. Thus the example from Table 7.2 of correction of pronunciation of the word ‘Hey’ as “Feature obstr should increase for HH” indicates that the phonological feature ‘obstr’ in the phoneme transcribed as ‘HH’ should be articulated more clearly to achieve correct pronunciation. In the case where FlexSR identifies a set of phonological features not associated with a phoneme of a target word at all, it may indicate that a different phoneme has been inserted. This is seen in example of the feedback from FlexSR on the nonsense word sound ‘strai’ as a pronunciation of the target word ‘I’ as “Insertion of S”.

7.3.1 FlexSR Feasibility Tests - Noise Attacks

FlexSR’s response to adversarial commands in noise, as demonstrated by Carlini et al. [31], was tested. Specifically, FlexSR’s response was tested to five adversarial commands in noise that had been demonstrated by Carlini et al. in black-box attacks on Google Now, as made available by the researchers.⁴ In testing FlexSR’s response to the adversarial commands, the relevant target command was given to FlexSR as the phrase attempting to be pronounced. Details of the adversarial commands in noise and of FlexSR’s response to these are shown in Table 7.1. FlexSR’s response to the corresponding target commands in live human speech was also tested. The most significant finding from these tests was that FlexSR did not recognise any of the adversarial commands in noise as correctly pro-

³The support of Prof. Aditi Lahiri’s research group in the University of Oxford’s Faculty of Linguistics, Philology and Phonetics in providing access to the FlexSR system for the purposes of these tests and for assisting with the tests is gratefully acknowledged.

⁴<http://www.hiddenvoicecommands.com/black-box>

nounced versions of the target commands. In one case, FlexSR refused to accept the noise command as speech input altogether. This in itself was not sufficient to separate the adversarial commands from the legitimate target commands as spoken by a human, as FlexSR also identified three of the target commands in natural speech as incorrectly pronounced. However, with the exception of the adversarial command for the ‘Hey Google’ wake phrase, FlexSR did identify a much larger number of mispronounced features for the adversarial commands in noise than for target commands in natural speech for which incorrect pronunciation was identified. This suggests the number of features identified as mispronounced by FlexSR could be used as criteria for distinguishing legitimate ‘natural’ input from malicious ‘unnatural’ input in a bespoke implementation of FlexSR for security purposes. Any bespoke implementation of FlexSR for security purposes would need to address the issue of false positives that is evident in the feasibility test using the implementation of FlexSR as a language learning tool.

Table 7.1: FlexSR noise attacks tests - Results

Target command	human (non-adversarial) / noise (adversarial)	FlexSR correctly pronounced yes/no	FlexSR no. of features mispronounced	Details of FlexSR features mispronounced
What is my current location	human	NO	2	1) Feature voice should decrease for T in what; 2) Feature stop should increase for T in current
	noise	NO	5	1) Feature stop should increase for T; 2) Feature str should increase for z in is; 3) Feature nas should increase for N in current; 4) Feature cor should increase for T in current; 5) Feature nas should increase for N in location
Tweet goodbye	human	NO	1	Feature voice should increase for D
	noise	NO	4	1) Feature stop should increase for T; 2) Feature stop should increase for T; 3) Feature voice should increase for D; 4) Feature labial should increase for B
OK Google	human	NO	1	Feature RTR should increase for AX
	noise	NO	1	Feature cor should increase for EY
Turn on airplane mode	human	YES	0	none
	noise	NO	n.a.	<i>Does not take it as a command</i>
Call 911	human	YES	0	none

	noise	NO	6	1) Feature nas should increase for N in nine; 2) Feature nas should increase for N in nine; 3) Feature cor should increase for N in one; 4) Feature rtr should increase for AH; 5) Feature nas should increase for N; 6) Insertion of NG at the end
--	-------	----	---	---

7.3.2 FlexSR Feasibility Tests - Nonsense Attacks

The feasibility tests relating to nonsense attacks aimed to assess FlexSR’s ability to distinguish between real words spoken by a human, as would be the case in a legitimate voice command, and nonsense words that are vocalized by speech synthesis, as would be the case in a nonsense attack on speech recognition in a voice-controlled system, such as that demonstrated in Chapter 5. The tests assessed FlexSR’s ability to distinguish between these two types of input based on its response to individual word sounds, rather than to specific adversarial commands. The reason for this was that unlike the speech recognition systems used by voice-controlled digital assistants such as Google Assistant, FlexSR does not use a language model to assess the likelihood of words appearing in combination. Therefore FlexSR’s response to sound input is determined solely by acoustic features, and its response to individual word sounds will not have any dependence on the context in which these word sounds appear. The absence of a language model might represent a security benefit of FlexSR, as the mistranscription of one word sound will not increase the probability of mistranscription of other word sounds in an adversarial command, as was seen to be the case in the nonsense attacks demonstrated on Google Assistant as detailed.

The dataset used in the feasibility tests consisted of 14 real words from the wake phrase and target commands used in the experimental work on nonsense attacks described in Chapter 5, and a set of three rhyming nonsense words for each these real words that were randomly selected from the full nonsense word sets for each word as listed in Appendix A.1. This dataset was used to compare FlexSR’s response to naturally produced real words and unnaturally produced nonsense words. The real words were inputted to FlexSR in natural voice, whereas the nonsense words were inputted as audio files of the nonsensical word sounds synthesised with Amazon Polly speech synthesis.

In each test, FlexSR was given one of the fourteen wake phrase or target command words as the word that the user was attempting to pronounce. This word was then spoken into the phone’s microphone by a native speaker of English. Audio files of the three rhyming nonsense words for the relevant word were then played to the phone’s microphone, with FlexSR still being given the real word as the word that was being spoken. The system’s response to each real word and to each nonsense word was recorded. This response indicated whether FlexSR recognised the input as an acceptable pronunciation of the relevant word or not, and, if not, which features it determined as being incorrectly pronounced.

The results of the FlexSR tests are shown in Table 7.2. Of the 42 nonsense words that were played to FlexSR, FlexSR considered only six as valid pronunciations of the real word, indicating a potential ability to identify most nonsensical word sounds as invalid

input. However, FlexSR also identified seven of the 14 real words spoken to it by a human as incorrectly pronounced, implying a high percentage of false positives in detection of potential attacks. This may be on account of the specific implementation of FlexSR in the system used for the tests as a language learning tool that aims to refine pronunciation rather than just ensure comprehensibility of spoken words. It might be possible in a different implementation to adjust FlexSR’s decision-making process to the aim of distinguishing human-comprehensible from human-incomprehensible input.

One aspect of the results that clearly separated real words from nonsensical words was the presence of a particular type of error identified by FlexSR as a mispronunciation with respect to some of the nonsense word sounds, namely a insertion error. In the case of an insertion error, rather than correcting the pronunciation of the phonemes in the word given to FlexSR as the word that the user is attempting to pronounce, FlexSR advises that the user’s speech represents the insertion of an entirely different phoneme. Erroneous insertion of a phoneme was identified as a pronunciation error in seven of the nonsense word sounds, but not in any of the real word sounds. Therefore the presence of this particular error type in the feedback from FlexSR generated in the tests represents confirmation that a word sound is invalid ‘synthetic’ input, although its absence does not conversely confirm that a word sound is valid ‘real’ input. This suggests the type of mispronunciations identified by FlexSR could be used as criteria for distinguishing legitimate ‘natural’ input from malicious ‘unnatural’ input in a bespoke implementation of FlexSR for security purposes. As in feasibility tests relating to noise attacks, the results indicated that any bespoke implementation of FlexSR for security purposes would need to address the issue of false positives.

Table 7.2: FlexSR nonsense attacks tests - Results

Target word	Word sounds	FlexSR correctly pronounced yes/no	FlexSR no. of features mispronounced	Details of FlexSR features mispronounced
am	real word - human speech	YES	None	no features
am	nonsense word “Z’am” (“zham”) - synthesized speech	YES	None	no features
am	nonsense word “D’am” (“tham”) - synthesized speech	NO	One	feature low should increase for AE
am	nonsense word “spr’am” (“spram”) - synthesized speech	NO	One	feature low should increase for AE
blue	real word - human speech	YES	None	no features
blue	nonsense word “pr’u:” (“proo”) - synthesized speech	NO	Three	1) Feature lab should increase for B; 2) Feature lat should increase for L; 3) Feature lab should increase for UW

blue	nonsense word “Z’u:” (“zhoo”) - synthesized speech	NO	Two	1) Feature lab should increase for B; 2) Feature lat should increase for L
blue	nonsense word “skw’u:” (“squoo”) - synthesized speech	NO	Two	1) Feature lab should increase for B; 2) Feature high should increase for UW
google	real word - human speech	YES	None	no features
google	nonsense word “sl’u:p@L” (“sloople”) - synthesized speech	NO	One	1) Feature high should increase for UW
google	nonsense word “Z’u:p@L” (“zhoople”) - synthesized speech	NO	Two	1) Feature high should increase for UW; 2) Feature lat should increase for L
google	nonsense word “D’u:z@L” (“thoozle”) - synthesized speech	NO	One	1) Feature high should increase for UW
hey	real word - human speech	NO	One	1) Feature obstr should increase for HH
hey	nonsense word “tS’eI” (“chay”) - synthesized speech	NO	One	1) Feature rad should increase for HH
hey	nonsense word “sm’eI” (“smay”) - synthesized speech	NO	One	1) Feature rad should increase for HH
hey	nonsense word ‘Z’eI” (“zhay”) - synthesized speech	NO	One	1) Feature rad should increase for HH
I	real word - human speech	NO	One	1) Feature low should increase for AY
I	nonsense word “gr’aI” (“grai”) - synthesized speech	YES	None	no features
I	nonsense word “str’aI” (“strai”) - synthesized speech	NO	One	1) Insertion of S
I	nonsense word “sm’aI” (“smai”) - synthesized speech	NO	One	1) Insertion of S
light	real word - human speech	NO	Two	Feature stop should increase for T
light	nonsense word “pr’aIt” (“pright”) - synthesized speech	NO	Two	1) Feature low should increase for AY; 2) Feature stop should increase for T
light	nonsense word “j’aIt” (“yight”) - synthesized speech	NO	Three	1) Feature low should increase for AY; 2) Feature stop should increase for T; 3) Feature lat should increase for L

light	nonsense word “Tr’aIt” (“thright”) - synthesized speech	NO	One	1) Feature stop should increase for T
my	real word - human speech	NO	One	1) Feature low should increase for AY
my	nonsense word “bl’aI” (“blai”) - synthesized speech	NO	One	1) Insertion of K
my	nonsense word “skw’aI” (“squai”) - synthesized speech	NO	One	1) Insertion of K
my	nonsense word “j’aI” (“yai”) - synthesized speech	NO	One	1) Feature low should increase for AY
name	real word - human speech	YES	None	no features
name	nonsense word “sn’eIm” (“sname”) - synthesized speech	NO	One	1) Insertion of S
name	nonsense word “st’eIm” (“stame”) - synthesized speech	NO	Two	1) Feature nas should increase for N; 2) Feature cor should increase for EY
name	nonsense word “sm’eIm” (“smame”) - synthesized speech	NO	One	1) Insertion of S
off	real word - human speech	YES	None	no features
off	nonsense word “sl’Of” (“sloff”) - synthesized speech	NO	One	1) Insertion of K
off	nonsense word “D’Of” (“thoff”) - synthesized speech	YES	None	no features
off	nonsense word “S’Of” (“shoff”) - synthesized speech	YES	None	no features
on	real word - human speech	NO	One	1) Feature nas should increase for N
on	nonsense word “h’On” (“hon”) - synthesized speech	NO	One	1) Feature nas should increase for N
on	nonsense word “tr’On” (“tron”) - synthesized speech	NO	One	1) Feature nas should increase for N
on	nonsense word “v’On” (“von”) - synthesized speech	NO	One	1) Feature cor should increase for N
red	real word - human speech	YES	None	no features

red	nonsense word "D'Ed" ("thed") - synthesized speech	NO	Three	1) Feature rho should increase for R; 2) Feature RTR should increase for EH; 3) Feature obstr should increase for D
red	nonsense word "skw'Ed" ("sqwed") - synthesized speech	NO	Three	1) Feature rho should increase for R; 2) Feature RTR should increase for EH; 3) Feature obstr should increase for D
red	nonsense word "tS'Ed" ("ched") - synthesized speech	NO	Three	1) Feature rho should increase for R; 2) Feature RTR should increase for EH; 3) Feature obstr should increase for D
turn	real word - human speech	NO	Two	1) Feature cor shuld increase for T; 2) Feature cons should increase for N
turn	nonsense word "spr'3:n" ("sprum") - synthesized speech	NO	Two	1) Feature stop should increase for T; 2) Feature nas should increase for N
turn	nonsense word "m'3:n" ("murn") - synthesized speech	NO	Three	1) Feature stop should increase for T; 2) Feature rho should increase for ER; 3) Feature nas should increase for N
turn	nonsense word "str'3:n" ("strurn") - synthesized speech	NO	Two	1) Feature stop should increase for T; 2) Feature nas should increase for N
what's	real word - human speech	NO	Three	1) Feature dor should increase for OH; 2) Feature stop should increase for T; 3) Feature str should increase for S
what's	nonsense word "sm'0ts" ("smots") - synthesized speech	NO	One	1) Feature stop should increase for T
what's	nonsense word "f'0ts" ("fots") - synthesized speech	YES	None	no features
what's	nonsense word "tw'0ts" ("twots") - synthesized speech	NO	One	1) Feature high should increase for W
who	real word - human speech	YES	None	No features
who	nonsense word "Z'u:" ("zhoo") - synthesized speech	YES	None	no features

who	nonsense word “f’u:” (“foo”) - synthesized speech	NO	Two	1) Feature rad should increase for HH; 2) Feature high should increase for UW
who	nonsense word “v’u:” (“voo”) - synthesized speech	NO	Two	1) Feature rad should increase for HH; 2) Feature high should increase for UW

7.4 Defences against Attacks on Natural Language Understanding

Defences are required to counter covert attacks that target vulnerabilities in natural language understanding functionality by exploiting inadequacies in current methods for representing meaning in voice-controlled systems. An example of an attack of this type are the missense attacks demonstrated in Chapter 5 Section 2, in which it was shown to be possible to mislead a target system to accept an unrelated utterance as a target command. Current technologies for natural language understanding clearly represent only very crude approximations of the ‘true’ processes of language understanding in the human brain. Defences against attacks targeting natural language understanding will need to close the gaps between human and machine understanding of the meaning of utterances that leave a system vulnerable to malicious exploitation.

In theory, similar to defences based on knowledge on human word recognition to prevent attacks on speech recognition, defences against attacks on natural language understanding could be developed based on knowledge of human construction of meaning from spoken utterances, so as to mitigate the threat of attacks that exploit differences in human and machine understanding of natural language. In practice, however, much about the human processes for the construction of meaning remains unknown (see for example Pylkkänen et al. [190], Magnuson [151], Gunter et al. [76], Elston-Güttler and Friederici [56]). At a high level, research from neurolinguistics suggests that the human construction of meaning from sentences emerges from a complex interplay of word meanings and syntax, as shown for example by Fedorenko et al. [59] and Johnson and Goldberg [106]. Fedorenko et al. measured neural activity in epilepsy patients with subdural implants in response to four different types of linguistic input, namely ordinary sentence, non-grammatical strings of meaningful words, grammatical strings of nonsense words, and non-grammatical strings of nonsense words. The neural activity of subjects as measured in the study was different for each type of input, with the highest level of activity being seen in response to grammatical strings of real words, and the lowest level being recorded for the non-grammatical strings of nonsense words. The results confirmed that both syntax and semantics are essential factors in the construction of meaning in humans. Johnson and Goldberg provide evidence of the importance of syntactical structure as well as individual words meanings in human meaning representation. They conducted a semantic priming study in which participants’ response times in identifying a target word under two conditions were compared. In one condition, before being presented with the target word participants were presented with a grammatical string of nonsense words in which the tar-

get word could replace one of the nonsense words whilst retaining the syntactic coherence of the sentence. In the other condition, the presentation of the target words was preceded by a grammatical string of nonsense words that were not syntactically coherent with the target word. The results showed that syntactic coherence had a significant priming effect in terms of response times in identifying the target word.

Such studies confirm the compositional nature of meaning construction in the human brain in general, but do not shed light on this at a level of granularity that might elucidate the process of specific meaning construction. It is possible that future advances in neuroscience and linguistics may identify features of meaning representation in the human brain that capture the essence of the specific distinctions of meaning made in human natural language understanding, and that are capable of replication in machine natural language understanding. The discovery of such features would enable the natural language understanding of voice-controlled systems to become more closely aligned with human understanding both with respect to in-scope and out-of-scope input, thus minimising the potential for malicious exploitation of gaps between human understanding of spoken language and its imitation in artificial systems. However, this possibility remains futuristic and nebulous at present.

A different approach to defending against attacks targeting natural language understanding in voice-controlled systems is suggested by the work of Navigli and Ponzetto [174] on using multilingual semantic networks for word sense disambiguation, as discussed in Chapter 2. The basis of the approach to word sense disambiguation proposed by Navigli and Ponzetto is that the set of word senses associated with a given word are unlikely to remain constant across different languages. Different senses of the same word in one language are likely to be translated as different words in another language, and thus the translation of a word as used in a particular context can be used to determine the correct word sense. This idea might be adapted to detect covert attacks on natural language understanding in voice-controlled systems by translating utterances inputted to a voice-controlled system and comparing the interpretation of user intent from the utterance in different languages. In the case of a non-malicious command, the interpretation of user intent is likely to remain constant across languages. However, in the case of malicious input aiming to mislead natural language understanding, for example by crafted use of homophones, the interpretation of user intent is unlikely to remain constant in different languages. A defence mechanism based on cross-lingual comparison of natural language understanding outputs could be implemented at the dialogue management stage of handling of speech input by a voice-controlled system. In the event that the intent extracted from a user utterance changes in a different language, execution of the intent would be blocked.

7.4.1 Google Translate Feasibility Tests - Missense Attacks

The potential of the proposed approach to defend against attacks on natural language understanding in voice-controlled systems could be demonstrated in the context of the word transplant attacks demonstrated in Chapter 5, using the readily available machine translation technology Google Translate.⁵ Google Translate uses RNNs for sequence-to-sequence

⁵See <https://translate.google.co.uk/>

mapping of input in one language to output in another [242]. In these feasibility tests, the number of words shared between target command and adversarial utterance in the original language compared to the number of words shared between them in translation is used as a proxy for the relative likelihood of the success of the adversarial utterance in the original language and in translation in a real multilingual system. Translation into German was chosen as an example of possible cross-lingual comparisons.

Table 7.3 shows Google Translate’s translation of three of successful adversarial utterances generated in the word transplant attacks on the Alexa Skills and RASA NLU demonstrated in the pilot experiment described in Chapter 5 (one of the successful adversarial utterances, ‘bill of an anchor’, is excluded because Google Translate’s German translation of it as ‘Rechnung eines Ankers’ was incorrect). It is evident that, with one exception, the number of words shared between target command and adversarial utterance drops in translation, and thus that the interpretation of user intent from the adversarial utterances is unlikely to remain constant across the two languages. Table 7.4 shows the Google Translate translation of some of the adversarial utterances that were successful in triggering the target Easter Egg commands in the main experiment described in Chapter 5 (a number of the successful adversarial utterances are not included due to their translation by Google Translate being incorrect). Similarly to the adversarial utterances in the pilot experiment, all adversarial utterances apart from one are seen to have fewer words in common with the target command in translation, and are thus unlikely to be mapped to the same erroneous semantic representation. Even in the case of the adversarial utterance for which the number of words in common with the target command does not drop in translation, it seems reasonable to speculate that the utterance would be unlikely to be misinterpreted in the same way in the original language and in translation in a real system, given that the number of words in common between adversarial utterance and target command is being used here only as a proxy for the likelihood of adversarial utterances being misinterpreted in the same way by separate machine learning-based systems in different languages.

Table 7.3: Feasibility test for defence against word transplant attacks on natural language understanding based on cross-lingual comparison (pilot experiment)

Target command (original language)	Adversarial command (original language)	No. of words retained (original language)	Target command (translation)	Adversarial command (translation)	No. of words retained (translation)
tell me the current balance	I kept my balance in the current	2	Sag mir das aktuelle Guthaben	Ich habe mein Gleichgewicht im Strom gehalten	0
show me all my transactions	the transactions were for show	2	Zeig mir alle meine Transaktionen	Die Transaktionen waren für die Show	1
think my card is stolen	your card is an ace	1	Ich glaube meine Karte ist gestohlen	Ihre Karte ist ein Ass	1

Table 7.4: Feasibility test for defence against word transplant attacks on natural language understanding based on cross-lingual comparison (main experiment)

Target command (original language)	Adversarial command (original language)	No. of words retained (original language)	Target command (translation)	Adversarial command (translation)	No. of words retained (translation)
rock paper scissors	these paper scissors rock	3	Schere, Stein, Papier	Diese Papierscheren rocken	2
rock paper scissors	rock the paper and scissors	3	Schere, Stein, Papier	Schaukeln Sie das Papier und die Schere	2
do you like cats or dogs?	are you more like cats or dogs?	3	Magst du Katzen oder Hunde?	Bist du eher wie Katzen oder Hunde	2
open the pod bay doors	open the door there's a pod in the bay	4	öffne die Türen des Pods	Öffne die Tür, da ist eine Schote in der Bucht	2
may the Force be with you	it may force me to be with you	1	möge die Macht mit dir sein	es kann mich zwingen, bei dir zu sein	0
is the cake a lie?	is low-fat cake a lie?	2	ist der Kuchen eine Lüge	ist fettarmer Kuchen eine Lüge	2

The feasibility tests presented here with respect to the potential of cross-lingual comparison as a defence against attacks on natural language understanding are limited to the word transplant attacks demonstrated in Chapter 5, and thus do not include the word substitution attacks presented in the same chapter. In the absence of a real multilingual voice-controlled system to use in testing, it was not possible to test the effectiveness of this approach as a potential defence against other types of attacks on natural language understanding, such as word substitution attacks. However, it seems reasonable to speculate that cross-lingual comparison may have potential as a defence against attacks on natural language understanding beyond word transplant attacks. In the case of word substitution attacks, it seems unlikely that a system would be vulnerable to substitution of the same target command word by the same unrelated word in two different languages. The space of unexpected input that is accepted as valid by one system is likely to be mostly different to the space of unexpected input accepted by a separate system trained with data in a different language.

7.5 Defences against Attacks on Dialogue Management

Defences are required to counter bot-assisted mistraining attacks on dialogue management as may become possible in future in the context of voice-controlled systems based on reinforcement learning. Mistraining attacks aiming to actively manipulate a system's functionalities are likely to require a different type of defence mechanism to attacks that passively exploit existing vulnerabilities, such as the adversarial learning attacks targeting

vulnerabilities in speech recognition and natural language understanding discussed above. Possibilities for defences to counter mistraining attacks targeting dialogue management in voice-controlled systems might include the application of research on human dialogue interaction, in particular the application of research on detection of harmful conversational interactions between humans. Defence mechanisms based on such research would aim to identify unnatural input to voice-controlled systems that is intended to undermine its ability to respond appropriately to legitimate user input. An example of such unnatural input might be discontinuous dialogue turns.

There are also other features of dialogue interaction that might be used to detect attempts to mistrain a voice-controlled system. Hancock et al. [82] claim that online interactions in which one party is attempting to deceive the other can be distinguished from non-deceitful interactions based on features of linguistic style. Furthermore, the authors of this research claim that these distinctions can be made with respect to the language use of the deceived party as well as the deceiving party, on account of the phenomenon of Linguistic Style Matching (LSM) whereby conversation partners are observed to subconsciously mimic each other's linguistic style. Muir et al. [168] claim that the direction of LSM can be shown to be related to social and personal power dynamics between conversation partners. Such research might be adapted to enable the detection of potential mistraining attacks. Hill et al. [91] have identified that humans tend to change their language use in communicating with a chatbot as compared to their language use in communicating with other humans. Thus another hallmark of malicious dialogue interactions aiming to mistrain a reinforcement learning-based system might be that, in contrast to non-malicious interactions, a user's use of language shows no adaptation to the voice-controlled system. In other words, whereas in a non-malicious interaction between human and computer, both parties might be observed to adapt their language use to one another, in a malicious interaction where the aim is to mistrain the system, only the computer will show a change in language use over time.

Other potentially relevant research includes work such as that of Rashid et al. [194] aiming to identify distinctive patterns of dialogue interactions in which one party is attempting to gain some form of abusive control over the other, such as in online grooming or radicalisation. Insights from research in these areas might be adapted to enable detection of malicious attempts to gain control of a voice-controlled system via mistraining of its reinforcement learning-based dialogue management functionality. On detection of such attempts at mistraining, the dialogue management component would be able to block such input from being used in the continuous training of its own functionality from interactions with the external environment.

7.6 Recovery from False Positives

The implementation of defence mechanisms as envisaged would require some functionality for overriding the blocking of command execution in the case where detection of malicious input is a false positive. The feasibility tests with FlexSR presented above illustrate the potential for false positives in using alternative speech recognition technology as a defence against attacks on speech recognition. Similarly, using cross-lingual natural lan-

guage understanding as a defence against attacks on natural language understanding might generate false positives in instances where machine translation of legitimate voice commands is incorrect. In defences against attacks on dialogue management, false positives might be generated if legitimate interactions between users and their devices are misidentified as mistraining attempts originating from bots. One possibility for enabling users to override false positives might be to incorporate some form of CAPTCHA challenge, such as a question answering task, in the verbal security alerts that are triggered on detection of potentially malicious input. Similar to CAPTCHAs based on differences between human and machine language understanding used to distinguish human users from bots in online authentication, CAPTCHAs used to override false positives in detection of attacks via the speech interface might present a natural language understanding challenge that an attacker acting through a bot would find difficult to respond to.

Chapter 8

Conclusions and Future Work

The work presented in this thesis confirms that human-computer interaction by speech is associated with particular security vulnerabilities that may enable malicious actors to gain unauthorised access to a system. The potential for attacks via the speech interface can be expected to expand as the state-of-the-art in voice control progresses. Unresolved security concerns may hamper the development of human-computer interaction by speech to its full potential.

This thesis makes three main contributions towards resolving these concerns, by addressing the research questions of whether the speech interface presents new vulnerabilities not contained in other types of interfaces, and if so how potential attacks exploiting these vulnerabilities might be identified and defended against. These research questions are addressed in three main components of the thesis. Firstly, the thesis presents a novel taxonomy of attacks via the speech interface, drawing on a comprehensive review of prior and related work. The taxonomy identifies the different types of attack that are specifically associated with the speech interface, using an inclusive categorisation principle. The taxonomy serves both as a systemisation of knowledge from prior work, as well as as a tool for considering the potential for new types of attacks that have not yet been demonstrated. The second component of the thesis builds on the conclusions from the taxonomy in demonstrating two new types of attacks that are envisaged in the taxonomy, but have not been considered in prior work. This involves experimental work that validates the potential for these two new types of attack using systematic methodologies. The third component of the thesis addresses the question of defences against attacks via the speech interface. For this purpose, a new attack and defence modelling approach is developed that represents a novel application of the Orient-Observe-Decide-Act (OODA) loop model to the context of human-computer interaction by speech. The attack and defence model is used to review the effectiveness of currently available defence mechanisms in countering the different types of attacks via the speech interface. This is followed by proposals for the development of new defence mechanisms against types of attacks that the attack and defence modelling analysis identifies as being inadequately defended by currently available mechanisms. The proposals for new defence mechanisms include the results of feasibility tests that link the speculative proposals to existing technologies that might be applied to countering attacks via the speech interface.

Chapter 2 of the thesis presented technical background on human-computer interaction by speech and on cloud-based voice assistants. Chapter 3 of the thesis presented a critical review of prior work on attacks via the speech interface and related areas. The review identified four attack methods that might be used in an attack via the speech interface, namely plain speech, inaudible sound injection, adversarial learning, and active attack. The review also discussed the human factors that may contribute to the effectiveness of attacks via the speech interface.

Chapter 4 of the thesis presented the novel taxonomy of attacks via the speech interface. Drawing on the critical review of prior work in Chapter 3, the taxonomy grouped attacks via the speech interface according to human perceptual categories, rather than according to particular attack methods or technical vulnerabilities. Attacks via the speech interface were primarily grouped in the taxonomy into two categories, namely overt attacks, which are plainly audible to human listeners, and covert attacks, which are imperceptible or unrecognisable to human listeners. Covert attacks were subcategorised in the taxonomy as attacks that are perceived by humans as silence, noise, music, nonsense or missense (defined as unrelated speech).

Chapter 5 of the thesis presented the experimental work demonstrating two new types of covert attacks via the speech interface that were implied by the taxonomy presented in Chapter 4. One of the attacks demonstrated in the experimental work was an attack targeting speech recognition functionality that used nonsensical word sounds to hide malicious commands, whereas the other attack targeted natural language understanding by hiding malicious voice commands in unrelated utterances that retained some of the elements of the target command.

Chapter 6 of the thesis presented the new approach to attack and defence modelling for the context of the speech interface based on the OODA loop concept. This involved mapping the categories of attacks identified in the taxonomy presented in Chapter 4 and the defences currently available to counter them to an OODA loop-based model of human-computer interaction by speech. The effectiveness of current defence mechanisms was then reviewed within the framework of the model. This analysis made clear that there are currently no fully effective defence mechanisms available to counter attacks via the speech interface, particularly in relation to attacks targeting vulnerabilities in speech recognition and natural language understanding.

Following the conclusions of the attack and defence modelling work, Chapter 7 of the thesis presented proposals for the development of new defence mechanisms against attacks via the speech interface that target speech recognition and natural language understanding, and also made some suggestions on possible defences against active attacks targeting dialogue management, as may arise in future. The proposals envisaged the implementation of a defence capability against adversarial learning and active attacks as part of a voice-controlled system's dialogue management functionality. On detection of potentially malicious input, the dialogue management component would issue a verbal security alert to users via the system's speech synthesis functionality.

8.1 Limitations and Future Work

Future work on the security of human-computer interaction by speech should further expand the range of attacks demonstrated in experimental work, whilst simultaneously seeking to identify and validate possible defences to counter such attacks. A dual focus on demonstrating new types of attacks via the speech interface whilst simultaneously investigating and developing defence mechanisms to counter them will ensure that the potential of voice-controlled technology can be fulfilled whilst also ensuring that any security risks associated with such advances are mitigated.

The experimental work on attacks via the speech interface presented in Chapter 5 of this thesis is necessarily limited in scope, both with respect to the possible attacks that are tested, as well as with respect to the specific systems targeted. With regard to the attacks using nonsensical word sounds presented in the first section of the chapter, the demonstration of attacks is limited to a small set of target commands for Google Assistant, and the testing of potential adversarial input is limited to random sampling rather than covering the whole space of potential adversarial input. Furthermore, at the level of individual results the reproducibility of these experiments is limited in three respects. Firstly, the random sampling approach taken in testing of adversarial input on the Google Assistant target system implies that any rerun of the experiment is likely to yield different results. Secondly, the response of the Google Assistant live system to the same adversarial input can be expected to change over time due to ongoing retraining of its speech recognition functionality, as was demonstrated by the retesting of results from the pilot experiment after a six-month interval. Thirdly, with regard to the human comprehensibility tests, the variability of results between experiment participants indicates that any rerun of the experiments with different participants would yield a different set of individual results. More large-scale work would be required in order to confirm the general reproducibility of this type of attack across different systems and contexts.

With regard to the attacks using unrelated utterances presented in the second section of Chapter 5, these are also limited to a small subset of target commands on one real-world system, Amazon Alexa. The reproducibility of test results on the live Alexa system can also be expected to be limited over time due to ongoing retraining of the system. Furthermore, the human-generated adversarial input used in the attacks with unrelated utterances can also be expected to be different in any rerun of the experiment with different human participants. An additional limitation of the attacks using unrelated utterances presented in Chapter 5 is that the understanding by humans of adversarial input as having a different meaning to the target commands is not formally validated. This element could be included in future experimental work of this type on a larger scale. More generally, future work on attacks via the speech interface should explore the space of potentially malicious input to speech interfaces more comprehensively, and seek to demonstrate attacks using such input on other target systems with different functionalities.

The review of current defences against attacks via the speech interface and the proposals for the development of new defence mechanisms presented in Chapters 6 and 7 of this thesis are also necessarily limited in various respects. The review of effectiveness of currently available defences presented in Chapter 6 is based on a theoretical evaluation of prior work. This could be expanded to include experimental comparison and validation of

the effectiveness of different types of defences against the various categories of attack in specific contexts. With regard to the proposals for the development of new defence mechanisms presented in Chapter 7, these proposals are speculative in nature. Future work should focus on implementation and testing of new defence mechanisms based on the speculative proposals presented in this thesis, and also explore any possibilities for other new types of defences.

Aside from further work on attacks and defences, another way in which the research presented in this thesis might be expanded in future is to conduct similar research in language contexts other than English. This might include not only research in other monolingual contexts, but also research in multilingual contexts. One possibility for future work might be to explore the potential for attacks that exploit overlap of sounds between different languages, by using malicious sound-based input that appears to a human user to be an utterance in one language but is understood by a voice assistant or other voice-controlled device as a target command in another language. Such cross-lingual attacks may become increasingly significant as voice-controlled systems are used in multilingual environments in future.

The implications of attacks via the speech interface will become increasingly serious as voice control is used to perform a wider range of actions in new contexts. Further research is needed to better understand the attack landscape for human-computer interaction by speech. Such research will allow the security issues associated with voice-controlled systems to be considered ahead of implementation, and enable defence mechanisms against attacks via the speech interface to be incorporated by design in such systems before they are deployed in practice.

Bibliography

- [1] Terrence Adams. “AI-Powered Social Bots”. In: *arXiv preprint arXiv:1706.05143* (2017).
- [2] Sadia Afroz et al. “Doppelgänger finder: Taking stylometry to the underground”. In: *Security and Privacy (SP), 2014 IEEE Symposium on*. IEEE. 2014, pp. 212–226.
- [3] Ioannis Agadakos et al. “Jumping the air gap: Modeling cyber-physical attack paths in the Internet-of-Things”. In: *Proceedings of the 2017 Workshop on Cyber-Physical Systems Security and Privacy*. 2017, pp. 37–48.
- [4] Sanchit Agarwal et al. “Parsing coordination for spoken language understanding”. In: *2018 IEEE Spoken Language Technology Workshop (SLT)*. IEEE. 2018, pp. 677–684.
- [5] Muhammad Ejaz Ahmed et al. “Void: A fast and light voice liveness detection system”. In: *29th USENIX Security Symposium (USENIX Security 20)*. Boston, MA: USENIX Association, Aug. 2020. URL: <https://www.usenix.org/conference/usenixsecurity20/presentation/ahmed>.
- [6] Hamad Al-Mohannadi et al. “Cyber-Attack Modeling Analysis Techniques: An Overview”. In: *Future Internet of Things and Cloud Workshops (FiCloudW), IEEE International Conference on*. IEEE. 2016, pp. 69–76.
- [7] Efthimios Alepis and Constantinos Patsakis. “Monkey Says, Monkey Does: Security and Privacy on Voice Assistants”. In: *IEEE Access* (2017).
- [8] Moustafa Alzantot, Bharathan Balaji, and Mani Srivastava. “Did you hear that? Adversarial Examples Against Automatic Speech Recognition”. In: *arXiv preprint arXiv:1801.00554* (2018).
- [9] Moustafa Alzantot et al. “Generating Natural Language Adversarial Examples”. In: *arXiv preprint arXiv:1804.07998* (2018).
- [10] Daniel Andor et al. “Globally normalized transition-based neural networks”. In: *arXiv preprint arXiv:1603.06042* (2016).
- [11] Stuart Armstrong, Kaj Sotola, and Seán S ÓhÉigeartaigh. “The errors, insights and lessons of famous AI predictions—and what they mean for the future”. In: *Journal of Experimental & Theoretical Artificial Intelligence* 26.3 (2014), pp. 317–342.

- [12] Vipul Arora, Aditi Lahiri, and Henning Reetz. “Phonological feature-based speech recognition system for pronunciation training in non-native language learning”. In: *The Journal of the Acoustical Society of America* 143.1 (2018), pp. 98–108.
- [13] James Auger. “Speculative design: crafting the speculation”. In: *Digital Creativity* 24.1 (2013), pp. 11–35.
- [14] Jonas Austerjost et al. “Introducing a virtual assistant to the lab: A voice user Interface for the intuitive control of laboratory instruments”. In: *SLAS TECHNOLOGY: Translating Life Sciences Innovation* 23.5 (2018), pp. 476–482.
- [15] Todd M Bailey and Ulrike Hahn. “Phoneme similarity and confusability”. In: *Journal of Memory and Language* 52.3 (2005), pp. 339–362.
- [16] Jerome R Bellegarda and Christof Monz. “State of the art in statistical methods for language and speech processing”. In: *Computer Speech & Language* 35 (2016), pp. 163–184.
- [17] Mohamed Benzeghiba et al. “Automatic speech recognition and speech variability: A review”. In: *Speech communication* 49.10-11 (2007), pp. 763–786.
- [18] Jonathan Berant et al. “Semantic parsing on freebase from question-answer pairs”. In: *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*. 2013, pp. 1533–1544.
- [19] Punica Bhardwaj. “Soteria: A Persuasive eSecurity Assistant”. PhD thesis. 2016.
- [20] Monowar H Bhuyan, Dhruva Kumar Bhattacharyya, and Jugal K Kalita. “Network anomaly detection: methods, systems and tools”. In: *IEEE communications surveys & tutorials* 16.1 (2014), pp. 303–336.
- [21] Timothy W Bickmore et al. “Patient and consumer safety risks when using conversational assistants for medical information: An observational study of Siri, Alexa, and Google Assistant”. In: *Journal of Medical Internet Research* 20.9 (2018).
- [22] Mary Bispham. “Security vulnerabilities in speech recognition systems”. In: *Oxford Cyber Security CDT Technical Paper Series* (2016).
- [23] Logan Blue, Luis Vargas, and Patrick Traynor. “Hello, Is It Me You’re Looking For?: Differentiating Between Human and Electronic Speakers for Voice Interface Security”. In: *Proceedings of the 11th ACM Conference on Security & Privacy in Wireless and Mobile Networks*. ACM. 2018, pp. 123–133.
- [24] Tom Bocklisch et al. “Rasa: Open source language understanding and dialogue management”. In: *arXiv preprint arXiv:1712.05181* (2017).
- [25] John R Boyd. “The essence of winning and losing”. In: *Unpublished lecture notes* 12.23 (1996), pp. 123–125.
- [26] Berndt Brehmer. “The dynamic OODA loop: Amalgamating Boyd’s OODA loop and the cybernetic approach to command and control”. In: *Proceedings of the 10th international command and control research technology symposium*. 2005, pp. 365–368.

- [27] Alexander Budanitsky and Graeme Hirst. “Semantic distance in WordNet: An experimental, application-oriented evaluation of five measures”. In: *Workshop on WordNet and other lexical resources*. Vol. 2. 2001, pp. 2–2.
- [28] Erik Cambria and Bebo White. “Jumping NLP curves: A review of natural language processing research”. In: *IEEE Computational intelligence magazine* 9.2 (2014), pp. 48–57.
- [29] Massimo Canonico and Luigi De Russis. “A Comparison and Critique of Natural Language Understanding Tools”. In: *CLOUD COMPUTING 2018* (2018), p. 120.
- [30] Nicholas Carlini and David Wagner. “Audio Adversarial Examples: Targeted Attacks on Speech-to-Text”. In: *arXiv preprint arXiv:1801.01944* (2018).
- [31] Nicholas Carlini et al. “Hidden voice commands”. In: *25th USENIX Security Symposium (USENIX Security 16)*, Austin, TX. 2016.
- [32] Justine Cassell. “Embodied conversational interface agents”. In: *Communications of the ACM* 43.4 (2000), pp. 70–78.
- [33] Davide Castelvecchi. “Can we open the black box of AI?” In: *Nature News* 538.7623 (2016), p. 20.
- [34] Asli Celikyilmaz, Li Deng, and Dilek Hakkani-Tür. “Deep Learning in Spoken and Text-Based Dialog Systems”. In: *Deep Learning in Natural Language Processing*. Springer, 2018, pp. 49–78.
- [35] Ching-Yun Chang and Stephen Clark. “Practical linguistic steganography using contextual synonym substitution and a novel vertex coding method”. In: *Computational Linguistics* 40.2 (2014), pp. 403–448.
- [36] Si Chen et al. “You can hear but you cannot steal: Defending against voice impersonation attacks on smartphones”. In: *Distributed Computing Systems (ICDCS), 2017 IEEE 37th International Conference on*. IEEE. 2017, pp. 183–195.
- [37] Geumhwan Cho et al. “Threat modeling and analysis of voice assistant applications”. In: *International Workshop on Information Security Applications*. Springer. 2018, pp. 197–209.
- [38] Noam Chomsky. “Three models for the description of language”. In: *IRE Transactions on information theory* 2.3 (1956), pp. 113–124.
- [39] Heidi Christensen et al. “homeService: Voice-enabled assistive technology in the home using cloud-based automatic speech recognition”. In: *Proceedings of the Fourth Workshop on Speech and Language Processing for Assistive Technologies*. 2013, pp. 29–34.
- [40] Hyunji Chung, Jungheum Park, and Sangjin Lee. “Digital forensic approaches for Amazon Alexa ecosystem”. In: *Digital Investigation* 22 (2017), S15–S25.
- [41] Hyunji Chung et al. “Alexa, Can I Trust You?” In: *Computer* 50.9 (2017), pp. 100–104.

- [42] Agata Ciesielski et al. “Vocal human-robot interaction inspired by Battle management language”. In: *Systems, Man, and Cybernetics (SMC), 2017 IEEE International Conference on*. IEEE. 2017, pp. 3379–3384.
- [43] Moustapha Cisse et al. “Houdini: Fooling Deep Structured Prediction Models”. In: *arXiv preprint arXiv:1707.05373* (2017).
- [44] John Coleman, John Aston, and Davide Pigole. “Reconstructing the sounds of words from the past”. In: *The Scottish Consortium for ICPhS* (2015).
- [45] Robert Dale. “The limits of intelligent personal assistants”. In: *Natural Language Engineering* 21.02 (2015), pp. 325–329.
- [46] Polly Dalton and Nick Fraenkel. “Gorillas we have missed: Sustained inattentive deafness for dynamic events”. In: *Cognition* 124.3 (2012), pp. 367–372.
- [47] Nilaksh Das et al. “ADAGIO: Interactive Experimentation with Adversarial Attack and Defense for Audio”. In: *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer. 2018, pp. 677–681.
- [48] Nitesh Dhanjani. *Abusing the Internet of Things: Blackouts, Freakouts, and Stakeouts.* ” O’Reilly Media, Inc.”, 2015.
- [49] Wenrui Diao et al. “Your voice assistant is mine: How to abuse speakers to steal information and control your phone”. In: *Proceedings of the 4th ACM Workshop on Security and Privacy in Smartphones & Mobile Devices*. ACM. 2014, pp. 63–74.
- [50] Fatiha Djebbar et al. “Comparative study of digital audio steganography techniques”. In: *EURASIP Journal on Audio, Speech, and Music Processing* 2012.1 (2012), p. 25.
- [51] Philip R Doyle et al. “Mapping Perceptions of Humanness in Speech-Based Intelligent Personal Assistant Interaction”. In: *arXiv preprint arXiv:1907.11585* (2019).
- [52] Hubert Dreyfus. “Alchemy and AI”. In: *Santa Monica, CA: RAND Corporation* (1965).
- [53] Aarthi Easwara Moorthy and Kim-Phuong L Vu. “Privacy concerns for use of voice activated personal assistant in the public space”. In: *International Journal of Human-Computer Interaction* 31.4 (2015), pp. 307–335.
- [54] Maria R Ebling. “Can Cognitive Assistants Disappear?” In: *IEEE Pervasive Computing* 15.3 (2016), pp. 4–6.
- [55] Arash Einolghozati et al. “Improving Semantic Parsing for Task Oriented Dialog”. In: *arXiv preprint arXiv:1902.06000* (2019).
- [56] Kerrie E Elston-Güttler and Angela D Friederici. “Ambiguous words in sentences: Brain indices for native and non-native disambiguation”. In: *Neuroscience letters* 414.1 (2007), pp. 85–89.
- [57] Carsten Eulitz and Aditi Lahiri. “Neurobiological evidence for abstract phonological representations in the mental lexicon during speech recognition”. In: *Journal of Cognitive Neuroscience* 16.4 (2004), pp. 577–583.

- [58] Ariel Ezrachi and Maurice E Stucke. “Is Your Digital Assistant Devious?” In: *Oxford Legal Studies Research Paper No. 52/2016; University of Tennessee Legal Studies Research Paper No. 304* (2016).
- [59] Evelina Fedorenko et al. “Neural correlate of the construction of sentence meaning”. In: *Proceedings of the National Academy of Sciences* 113.41 (2016), E6256–E6262.
- [60] Huan Feng, Kassem Fawaz, and Kang G Shin. “Continuous Authentication for Voice Assistants”. In: *arXiv preprint arXiv:1701.04507* (2017).
- [61] Asbjørn Følstad and Petter Bae Brandtzæg. “Chatbots and the new world of HCI”. In: *interactions* 24.4 (2017), pp. 38–42.
- [62] C Franzese and M Coyne. “The promise of voice: Connecting drug delivery through voice-activated technology”. In: 2017 (Dec. 2017), pp. 34–37.
- [63] Claudia K Friedrich, Aditi Lahiri, and Carsten Eulitz. “Neurophysiological evidence for underspecified lexical representations: Asymmetries with word initial variations.” In: *Journal of Experimental Psychology: Human Perception and Performance* 34.6 (2008), p. 1545.
- [64] Kevin Fu and Wenyuan Xu. “Risks of trusting the physics of sensors”. In: *Communications of the ACM* 61.2 (2018), pp. 20–23.
- [65] Yarin Gal. “Semantics, Modelling, and the Problem of Representation of Meaning—a Brief Survey of Recent Literature”. In: *arXiv preprint arXiv:1402.7265* (2014).
- [66] John S Garofolo et al. “DARPA TIMIT acoustic-phonetic continuous speech corpus CD-ROM. NIST speech disc 1-1.1”. In: *NASA STI/Recon technical report n 93* (1993).
- [67] Nikolay Gaubitch. “How voice ageing impacts biometric effectiveness”. In: *Biometric Technology Today* 2017.6 (2017), pp. 8–9.
- [68] Jairo Giraldo et al. “Security and privacy in cyber-physical systems: A survey of surveys”. In: *IEEE Design & Test* 34.4 (2017), pp. 7–17.
- [69] Alyssa Glass, Deborah L McGuinness, and Michael Wolverton. “Toward establishing trust in adaptive agents”. In: *Proceedings of the 13th International Conference on Intelligent User Interfaces*. ACM. 2008, pp. 227–236.
- [70] Wael H Gomaa and Aly A Fahmy. “A survey of text similarity approaches”. In: *International Journal of Computer Applications* 68.13 (2013), pp. 13–18.
- [71] Yuan Gong and Christian Poellabauer. “An Overview of Vulnerabilities of Voice Controlled Systems”. In: *arXiv preprint arXiv:1803.09156* (2018).
- [72] Yuan Gong and Christian Poellabauer. “Protecting Voice Controlled Systems Using Sound Source Identification Based on Acoustic Cues”. In: *2018 27th International Conference on Computer Communication and Networks (ICCCN)*. IEEE. 2018, pp. 1–9.
- [73] Ian Goodfellow, Nicolas Papernot, and Patrick McDaniel. “cleverhans v0. 1: an adversarial machine learning library”. In: *arXiv preprint arXiv:1610.00768* (2016).

- [74] Ian Goodfellow et al. “Generative adversarial nets”. In: *Advances in Neural Information Processing Systems*. 2014, pp. 2672–2680.
- [75] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. “Explaining and Harnessing Adversarial Examples”. In: *CoRR* abs/1412.6572 (2014).
- [76] Thomas C Gunter, Susanne Wagner, and Angela D Friederici. “Working memory and lexical ambiguity resolution as revealed by ERPs: A difficult case for activation theories”. In: *Journal of Cognitive Neuroscience* 15.5 (2003), pp. 643–657.
- [77] Mordechai Guri and Yuval Elovici. “Bridgeware: The air-gap malware”. In: *Communications of the ACM* 61.4 (2018), pp. 74–82.
- [78] Andrea L Guzman. “Making AI safe for humans: A conversation with Siri”. In: *Socialbots and Their Friends* (2016), pp. 85–101.
- [79] David Ha and Douglas Eck. “A neural representation of sketch drawings”. In: *arXiv preprint arXiv:1704.03477* (2017).
- [80] William Haack et al. “Security analysis of the Amazon Echo”. In: *Allen Institute for Artificial Intelligence* (2017).
- [81] Parisa Haghani et al. “From Audio to Semantics: Approaches to end-to-end spoken language understanding”. In: *arXiv preprint arXiv:1809.09190* (2018).
- [82] Jeffrey T Hancock et al. “On lying and being lied to: A linguistic analysis of deception in computer-mediated communication”. In: *Discourse Processes* 45.1 (2007), pp. 1–23.
- [83] John HL Hansen and Taufiq Hasan. “Speaker recognition by machines and humans: A tutorial review”. In: *IEEE Signal Processing Magazine* 32.6 (2015), pp. 74–99.
- [84] Simon Hansman and Ray Hunt. “A taxonomy of network and computer attacks”. In: *Computers & Security* 24.1 (2005), pp. 31–43.
- [85] Md Rashidul Hasan, Mustafa Jamil, MGRMS Rahman, et al. “Speaker identification using mel frequency cepstral coefficients”. In: *3rd International Conference on Electrical & Computer Engineering* 1.4 (2004).
- [86] Timothy J Hazen and Issam Bazzi. “A comparison and combination of methods for OOV word detection and word confidence scoring”. In: *Acoustics, Speech, and Signal Processing, 2001. Proceedings.(ICASSP’01). 2001 IEEE International Conference on*. Vol. 1. IEEE. 2001, pp. 397–400.
- [87] Peter Henderson et al. “Ethical challenges in data-driven dialogue systems”. In: *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*. ACM. 2018, pp. 123–129.
- [88] James Hendler. “Agents and the semantic web”. In: *IEEE Intelligent Systems* 16.2 (2001), pp. 30–37.
- [89] Jim Hendler and Tim Berners-Lee. “From the Semantic Web to social machines: A research challenge for AI on the World Wide Web”. In: *Artificial Intelligence* 174.2 (2010), pp. 156–161.

- [90] Cormac Herley and Paul C van Oorschot. “Sok: Science, security and the elusive goal of security as a scientific pursuit”. In: *Security and Privacy (SP), 2017 IEEE Symposium on*. IEEE. 2017, pp. 99–120.
- [91] Jennifer Hill, W Randolph Ford, and Ingrid G Farreras. “Real conversations with artificial intelligence: A comparison between human–human online conversations and human–chatbot conversations”. In: *Computers in Human Behavior* 49 (2015), pp. 245–250.
- [92] Geoffrey Hinton et al. “Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups”. In: *IEEE Signal Processing Magazine* 29.6 (2012), pp. 82–97.
- [93] Julia Hirschberg and Christopher D Manning. “Advances in natural language processing”. In: *Science* 349.6245 (2015), pp. 261–266.
- [94] Briland Hitaj, Giuseppe Ateniese, and Fernando Perez-Cruz. “Deep models under the GAN: information leakage from collaborative deep learning”. In: *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*. ACM. 2017, pp. 603–618.
- [95] Matthew B Hoy. “Alexa, Siri, Cortana, and More: An Introduction to Voice Assistants”. In: *Medical reference services quarterly* 37.1 (2018), pp. 81–88.
- [96] Hu Hu, Tian Tan, and Yanmin Qian. “Generative Adversarial Networks Based Data Augmentation for Noise Robust Speech Recognition”. In: *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE. 2018, pp. 5044–5048.
- [97] Xuedong Huang, James Baker, and Raj Reddy. “A Historical Perspective of Speech Recognition”. In: *Communications of the ACM* 57.1 (2014), pp. 94–103.
- [98] Jennifer Yang Hui and Dymples Leong. “The Era of Ubiquitous Listening: Living in a World of Speech-Activated Devices”. In: *Asian Journal Public Affairs* (2017), p. 66.
- [99] Ignacio Iacobacci, Mohammad Taher Pilehvar, and Roberto Navigli. “Embeddings for word sense disambiguation: An evaluation study”. In: *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Vol. 1. 2016, pp. 897–907.
- [100] Irina Illina and Dominique Fohr. “Out-of-Vocabulary Word Probability Estimation using RNN Language Model”. In: *8th Language & Technology Conference*. 2017.
- [101] Md Tamzeed Islam, Bashima Islam, and Shahriar Nirjon. “SoundSifter: Mitigating overhearing of continuous listening devices”. In: *Proceedings of the 15th Annual International Conference on Mobile Systems, Applications, and Services*. ACM. 2017, pp. 29–41.
- [102] Dan Iter, Jade Huang, and Mike Jermann. “Generating adversarial examples for speech recognition”. In: *Technical Report* (2017).

- [103] Catherine Jackson and Angela Orebaugh. “A study of security and privacy issues associated with the Amazon Echo”. In: *International Journal of Internet of Things and Cyber-Assurance* 1.1 (2018), pp. 91–100.
- [104] Robin Jia and Percy Liang. “Adversarial Examples for Evaluating Reading Comprehension Systems”. In: *arXiv preprint arXiv:1707.07328* (2017).
- [105] Robin Jia et al. “Certified robustness to adversarial word substitutions”. In: *arXiv preprint arXiv:1909.00986* (2019).
- [106] Matt A Johnson and Adele E Goldberg. “Evidence for automatic accessing of constructional meaning: Jabberwocky sentences prime associated verbs”. In: *Language and Cognitive Processes* 28.10 (2013), pp. 1439–1452.
- [107] Biing-Hwang Juang and Lawrence R Rabiner. “Automatic speech recognition—a brief history of the technology development”. In: *Georgia Institute of Technology. Atlanta Rutgers University and the University of California. Santa Barbara* 1 (2005), p. 67.
- [108] Kaarel Kaljurand and Tanel Alumäe. “Controlled natural language in speech recognition based user interfaces”. In: *International Workshop on Controlled Natural Language*. Springer. 2012, pp. 79–94.
- [109] Aishwarya Kamath and Rajarshi Das. “A Survey on Semantic Parsing”. In: *arXiv preprint arXiv:1812.00978* (2018).
- [110] Rohan Kar and Rishin Haldar. “Applying Chatbots to the Internet of Things: Opportunities and Architectural Elements”. In: *arXiv preprint arXiv:1611.03799* (2016).
- [111] Chaouki Kasmi and Jose Lopes Esteves. “IEMI threats for information security: Remote command injection on modern smartphones”. In: *IEEE Transactions on Electromagnetic Compatibility* 57.6 (2015), pp. 1752–1755.
- [112] Sean Kennedy et al. “I Can Hear Your Alexa: Voice Command Fingerprinting on Smart Home Speakers”. In: *2019 IEEE Conference on Communications and Network Security (CNS)*. IEEE. 2019, pp. 232–240.
- [113] Veton Këpuska and Gamal Bohouta. “Improving Wake-Up-Word and General Speech Recognition Systems”. In: *Dependable, Autonomic and Secure Computing, 15th Intl Conf on Pervasive Intelligence & Computing, 3rd Intl Conf on Big Data Intelligence and Computing and Cyber Science and Technology Congress (DASC/PiCom/DataCom/CyberSciTech), 2017 IEEE 15th Intl*. IEEE. 2017, pp. 318–321.
- [114] VZ Këpuska and TB Klein. “A novel wake-up-word speech recognition system, wake-up-word recognition task, technology and evaluation”. In: *Nonlinear Analysis: Theory, Methods & Applications* 71.12 (2009), e2772–e2789.
- [115] Omar Zia Khan and Ruhi Sarikaya. “Making Personal Digital Assistants Aware of What They Do Not Know.” In: *Interspeech*. 2016, pp. 1161–1165.
- [116] Daniel Khashabi et al. “On the Capabilities and Limitations of Reasoning for Natural Language Understanding”. In: *arXiv preprint arXiv:1901.02522* (2019).

- [117] Joo-Kyung Kim and Young-Bum Kim. “Joint Learning of Domain Classification and Out-of-Domain Detection with Dynamic Class Weighting for Satisficing False Acceptance Rates”. In: *arXiv preprint arXiv:1807.00072* (2018).
- [118] Brett King. *Bank 4.0: Banking Everywhere, Never at a Bank*. Marshall Cavendish Business, 2018.
- [119] Gabriel Klein, Jens Tolle, and Peter Martini. “From detection to reaction-A holistic approach to cyber defense”. In: *Defense Science Research Conference and Expo (DSR), 2011*. IEEE. 2011, pp. 1–4.
- [120] Thomas Kollar et al. “The Alexa Meaning Representation Language”. In: *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 3 (Industry Papers)*. Vol. 3. 2018, pp. 177–184.
- [121] Xiang Kong, Jeung-Yoon Choi, and Stefanie Shattuck-Hufnagel. “Evaluating automatic speech recognition systems in comparison with human perception results using distinctive feature measures”. In: *Acoustics, Speech and Signal Processing (ICASSP), 2017 IEEE International Conference on*. IEEE. 2017, pp. 5810–5814.
- [122] Moshe Koppel et al. “The “fundamental problem” of authorship attribution”. In: *English Studies* 93.3 (2012), pp. 284–291.
- [123] Tobias Kuhn. “A survey and classification of controlled natural languages”. In: *Computational Linguistics* 40.1 (2014), pp. 121–170.
- [124] Volodymyr Kuleshov et al. “Adversarial Examples for Natural Language Classification Problems”. In: *OpenReview submission OpenReview:r1QZ3zbAZ* (2018).
- [125] Solomon Kullback and Richard A Leibler. “On information and sufficiency”. In: *The annals of mathematical statistics* 22.1 (1951), pp. 79–86.
- [126] Anjishnu Kumar et al. “Just ASK: Building an architecture for extensible self-service spoken language understanding”. In: *arXiv preprint arXiv:1711.00549* (2017).
- [127] Deepak Kumar et al. “Skill Squatting Attacks on Amazon Alexa”. In: *27th USENIX Security Symposium (USENIX Security 18)*. Baltimore, MD: USENIX Association, 2018, pp. 33–47. ISBN: 978-1-931971-46-1. URL: <https://www.usenix.org/conference/usenixsecurity18/presentation/kumar>.
- [128] Alexey Kurakin et al. “Adversarial attacks and defences competition”. In: *arXiv preprint arXiv:1804.00097* (2018).
- [129] Finn Kuusisto. “Speech synthesis”. In: *XRDS: Crossroads, The ACM Magazine for Students* 21.1 (2014), pp. 63–63.
- [130] Il-Youp Kwak et al. “Voice Presentation Attack Detection through Text-Converted Voice Command Analysis”. In: *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*. ACM. 2019, p. 598.
- [131] Aditi Lahiri and Henning Reetz. “Distinctive features: Phonological underspecification in representation and processing”. In: *Journal of Phonetics* 38.1 (2010), pp. 44–59.

- [132] Aditi Lahiri and Henning Reetz. “Underspecified recognition”. In: *Laboratory phonology* 7 (2002), pp. 637–675.
- [133] Aditi Lahiri, Henning Reetz, and Philip Roberts. *Method and apparatus for automatic speech recognition*. US Patent App. 15/105,552. 2016.
- [134] Mathias Landhäußer, Sebastian Weigelt, and Walter F Tichy. “NLCI: a natural language command interpreter”. In: *Automated Software Engineering* 24.4 (2017), pp. 839–861.
- [135] Ian Lane et al. “Out-of-domain utterance detection using classification confidences of multiple topics”. In: *IEEE Transactions on Audio, Speech, and Language Processing* 15.1 (2007), pp. 150–161.
- [136] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. “Deep learning”. In: *Nature* 521.7553 (2015), pp. 436–444.
- [137] K-F Lee et al. “The SPHINX speech recognition system”. In: *International Conference on Acoustics, Speech, and Signal Processing*, IEEE. 1989, pp. 445–448.
- [138] Xinyu Lei et al. “The Insecurity of Home Digital Voice Assistants—Vulnerabilities, Attacks and Countermeasures”. In: *2018 IEEE Conference on Communications and Network Security (CNS)* (2018).
- [139] Oliver Lemon. “Learning what to say and how to say it: Joint optimisation of spoken dialogue management and natural language generation”. In: *Computer Speech & Language* 25.2 (2011), pp. 210–221.
- [140] Mike Lewis et al. “Deal or No Deal? End-to-End Learning for Negotiation Dialogues”. In: *arXiv preprint arXiv:1706.05125* (2017).
- [141] Fei Li and HV Jagadish. “Constructing an interactive natural language interface for relational databases”. In: *Proceedings of the VLDB Endowment* 8.1 (2014), pp. 73–84.
- [142] Jinfeng Li et al. “TextBugger: Generating Adversarial Text Against Real-world Applications”. In: *arXiv preprint arXiv:1812.05271* (2018).
- [143] Bin Liang et al. “Deep Text Classification Can be Fooled”. In: *arXiv preprint arXiv:1704.08006* (2017).
- [144] Percy Liang. “Learning executable semantic parsers for natural language understanding”. In: *Communications of the ACM* 59.9 (2016), pp. 68–76.
- [145] Percy Liang. “Talking to computers in natural language”. In: *XRDS: Crossroads, The ACM Magazine for Students* 21.1 (2014), pp. 18–21.
- [146] Richard P Lippmann et al. “Speech recognition by machines and humans”. In: *Speech communication* 22.1 (1997), pp. 1–15.
- [147] Pierre Lison and Raveesh Meena. “Spoken dialogue systems: the new frontier in human-computer interaction”. In: *XRDS: Crossroads, The ACM Magazine for Students* 21.1 (2014), pp. 46–51.

- [148] Wenyan Liu et al. “A Novel Framework for Zero-Day Attacks Detection and Response with Cyberspace Mimic Defense Architecture”. In: *Cyber-Enabled Distributed Computing and Knowledge Discovery (CyberC), 2017 International Conference on*. IEEE. 2017, pp. 50–53.
- [149] Ignacio Lopez-Moreno et al. “Automatic language identification using deep neural networks”. In: *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE. 2014, pp. 5337–5341.
- [150] George Loukas, Diane Gan, and Tuan Vuong. “A taxonomy of cyber attack and defence mechanisms for emergency management networks”. In: *Pervasive Computing and Communications Workshops (PERCOM Workshops), 2013 IEEE International Conference on*. IEEE. 2013, pp. 534–539.
- [151] James S Magnuson. “Mapping spoken words to meaning”. In: *Speech Perception and Spoken Word Recognition*. Psychology Press, 2016, pp. 86–106.
- [152] Nina McCurdy, Vivek Srikumar, and Miriah Meyer. “Rhymedesign: A tool for analyzing sonic devices in poetry”. In: *Proceedings of the Fourth Workshop on Computational Linguistics for Literature*. 2015, pp. 12–22.
- [153] Patrick McDaniel, Nicolas Papernot, and Z Berkay Celik. “Machine learning in adversarial settings”. In: *IEEE Security & Privacy* 14.3 (2016), pp. 68–72.
- [154] Marjorie McShane, Kevin Blissett, and Irene Nirenburg. “Treating Unexpected Input in Incremental Semantic Analysis”. In: *Proceedings of The Fifth Annual Conference on Advances in Cognitive Systems*. Palo Alto, CA: Cognitive Systems Foundation. 2017.
- [155] Michael McTear. *Conversational Modelling for ChatBots: Current Approaches and Future Directions*. Tech. rep. Technical report, Ulster University, Ireland, 2018.
- [156] Michael McTear, Zoraida Callejas, and David Griol. *The conversational interface*. Springer, 2016.
- [157] Yan Meng et al. “Securing Consumer IoT in the Smart Home: Architecture, Challenges, and Countermeasures”. In: *IEEE Wireless Communications* 25.6 (2018), pp. 53–59.
- [158] Grégoire Mesnil et al. “Using recurrent neural networks for slot filling in spoken language understanding”. In: *IEEE/ACM Transactions on Audio, Speech and Language Processing (TASLP)* 23.3 (2015), pp. 530–539.
- [159] Hendrik Meutzner, Santosh Gupta, and Dorothea Kolossa. “Constructing secure audio captchas by exploiting differences between humans and machines”. In: *Proceedings of the 33rd annual ACM conference on human factors in computing systems*. ACM. 2015, pp. 2335–2338.
- [160] Assaf Hurwitz Michaely et al. “Keyword Spotting for Google Assistant Using Contextual Speech Recognition”. In: *Proceedings of ASRU*. 2017.
- [161] Tomas Mikolov et al. “Efficient estimation of word representations in vector space”. In: *arXiv preprint arXiv:1301.3781* (2013).

- [162] David J Miller et al. “Adversarial learning: a critical review and active learning study”. In: *arXiv preprint arXiv:1705.09823* (2017).
- [163] George A Miller. “WordNet: a lexical database for English”. In: *Communications of the ACM* 38.11 (1995), pp. 39–41.
- [164] Richard Mitev, Markus Miettinen, and Ahmad-Reza Sadeghi. “Alexa Lied to Me: Skill-based Man-in-the-Middle Attacks on Virtual Assistants”. In: *Proceedings of the 2019 ACM Asia Conference on Computer and Communications Security*. ACM. 2019, pp. 465–478.
- [165] Richard Montague. “The proper treatment of quantification in ordinary English”. In: *Approaches to natural language*. Springer, 1973, pp. 221–242.
- [166] Roger K Moore. “Is spoken language all-or-nothing? Implications for future speech-based human-machine interaction”. In: *Dialogues with Social Robots*. Springer, 2017, pp. 281–291.
- [167] Omar Mubin, Christoph Bartneck, and Loe Feijs. “Towards the design and evaluation of ROILA: a speech recognition friendly artificial language”. In: *International Conference on Natural Language Processing*. Springer. 2010, pp. 250–256.
- [168] Kate Muir et al. “Characterizing the linguistic chameleon: Personal and social correlates of linguistic style accommodation”. In: *Human Communication Research* 42.3 (2016), pp. 462–484.
- [169] Dibya Mukhopadhyay, Maliheh Shirvanian, and Nitesh Saxena. “All your voices are belong to us: Stealing voices to fool humans and machines”. In: *European Symposium on Research in Computer Security*. Springer. 2015, pp. 599–621.
- [170] Prakash M Nadkarni, Lucila Ohno-Machado, and Wendy W Chapman. “Natural language processing: an introduction”. In: *Journal of the American Medical Informatics Association* 18.5 (2011), pp. 544–551.
- [171] Arvind Narayanan et al. “On the feasibility of internet-scale author identification”. In: *2012 IEEE Symposium on Security and Privacy*. IEEE. 2012, pp. 300–314.
- [172] Clifford Nass and Youngme Moon. “Machines and mindlessness: Social responses to computers”. In: *Journal of Social Issues* 56.1 (2000), pp. 81–103.
- [173] Cherie-Ann O Nathan et al. “The voice-controlled robotic assist scope holder AE-SOP for the endoscopic approach to the sella”. In: *Skull Base* 16.3 (2006), p. 123.
- [174] Roberto Navigli and Simone Paolo Ponzetto. “Joining forces pays off: Multilingual joint word sense disambiguation”. In: *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*. Association for Computational Linguistics. 2012, pp. 1399–1410.
- [175] Anh Nguyen, Jason Yosinski, and Jeff Clune. “Deep neural networks are easily fooled: High confidence predictions for unrecognizable images”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2015, pp. 427–436.

- [176] Van Nhan Nguyen. “Using linguistic knowledge for improving automatic speech recognition accuracy in air traffic control”. MA thesis. 2016.
- [177] Martin A Nowak and David C Krakauer. “The evolution of language”. In: *Proceedings of the National Academy of Sciences* 96.14 (1999), pp. 8028–8033.
- [178] Peter Oehlert. “Violating assumptions with fuzzing”. In: *IEEE Security & Privacy* 3.2 (2005), pp. 58–62.
- [179] Tim Paek and Roberto Pieraccini. “Automating spoken dialogue management design using machine learning: An industry perspective”. In: *Speech Communication* 50.8-9 (2008), pp. 716–729.
- [180] Nicolas Papernot et al. “Crafting adversarial input sequences for recurrent neural networks”. In: *Military Communications Conference, MILCOM 2016-2016 IEEE*. IEEE. 2016, pp. 49–54.
- [181] Nicolas Papernot et al. “Practical black-box attacks against deep learning systems using adversarial examples”. In: *arXiv preprint arXiv:1602.02697* (2016).
- [182] Nicolas Papernot et al. “Towards the science of security and privacy in machine learning”. In: *arXiv preprint arXiv:1611.03814* (2016).
- [183] Bryson R Payne, Leonardo I Mazuran, and Tamirat Abegaz. “Voice Hacking: Using Smartphones to Spread Ransomware to Traditional PCs”. In: *Journal of Cybersecurity Education, Research and Practice* 2018.1 (2018), p. 2.
- [184] Vittorio Perera et al. “Multi-task learning for parsing the alexa meaning representation language”. In: *American Association for Artificial Intelligence (AAAI)*. 2018, pp. 181–224.
- [185] Giuseppe Petracca et al. “Audroid: Preventing attacks on audio channels in mobile devices”. In: *Proceedings of the 31st Annual Computer Security Applications Conference*. ACM. 2015, pp. 181–190.
- [186] Roberto Pieraccini and Lawrence Rabiner. *The voice in the machine: building computers that understand speech*. MIT Press, 2012.
- [187] John R Pierce. “Whither speech recognition?” In: *The Journal of the Acoustical Society of America* 46.4B (1969), pp. 1049–1051.
- [188] David Pogue. “At Your Command”. In: *Scientific American* 315.1 (2016), pp. 25–25.
- [189] Michael Pucher et al. “Phonetic distance measures for speech recognition vocabulary and grammar optimization”. In: *3rd Congress of the Alps Adria Acoustics Association*. 2007, pp. 2–5.
- [190] Liina Pyykkänen, Rodolfo Llinás, and Gregory L Murphy. “The representation of polysemy: MEG evidence”. In: *Journal of cognitive neuroscience* 18.1 (2006), pp. 97–109.
- [191] Yuan Qi and Jing Xiao. “Fintech: AI powers financial services to improve people’s lives”. In: *Communications of the ACM* 61.11 (2018), pp. 65–69.

- [192] Long Qin. “Learning out-of-vocabulary words in automatic speech recognition”. PhD thesis. Carnegie Mellon University, 2013.
- [193] Anirudh Raju et al. “Data Augmentation for Robust Keyword Spotting under Playback Interference”. In: *arXiv preprint arXiv:1808.00563* (2018).
- [194] Awais Rashid et al. “Who am i? analyzing digital personas in cybercrime investigations”. In: *Computer* 46.4 (2013), pp. 54–61.
- [195] Suman Ravuri and Andreas Stolcke. “A comparative study of recurrent neural network models for lexical domain classification”. In: *Acoustics, Speech and Signal Processing (ICASSP), 2016 IEEE International Conference on*. IEEE. 2016, pp. 6075–6079.
- [196] Suman Ravuri and Andreas Stolcke. “Recurrent neural network and lstm models for lexical utterance classification”. In: *Sixteenth Annual Conference of the International Speech Communication Association*. 2015.
- [197] Konrad Rieck and Pavel Laskov. “Detecting unknown network attacks using language models”. In: *International Conference on Detection of Intrusions and Malware, and Vulnerability Assessment*. Springer. 2006, pp. 74–90.
- [198] Adam Charles Roberts, Allison Wetterlin, and Aditi Lahiri. “Aligning mispronounced words to meaning: Evidence from ERP and reaction time studies”. In: *The Mental Lexicon* 8.2 (2013), pp. 140–163.
- [199] AW Roscoe et al. “The attacker in ubiquitous computing environments: Formalising the threat model”. In: *Proceedings of FAST 2003 Pisa*. 2003.
- [200] Nirupam Roy et al. “Inaudible Voice Commands: The Long-Range Attack and Defense”. In: *15th USENIX Symposium on Networked Systems Design and Implementation NSDI 18*. USENIX Association. 2018, pp. 547–560.
- [201] Jeffrey N Rule. *A Symbiotic Relationship: The OODA Loop, Intuition, and Strategic Thought*. US Army War College, 2013.
- [202] Sanjib Kumar Saha, Abhijit Kumar Nag, and Dipankar Dasgupta. “Human-cognition-based CAPTCHAs”. In: *IT Professional* 17.5 (2015), pp. 42–48.
- [203] Md Sahidullah et al. “Introduction to Voice Presentation Attack Detection and Recent Advances”. In: *Handbook of Biometric Anti-Spoofing*. Springer, 2019, pp. 321–361.
- [204] Shumpei Sano, Nobuhiro Kaji, and Manabu Sassano. “Predicting Causes of Reformulation in Intelligent Assistants”. In: *arXiv preprint arXiv:1707.03968* (2017).
- [205] Ruhi Sarikaya. “The Technology Behind Personal Digital Assistants: An overview of the system architecture and key components”. In: *IEEE Signal Processing Magazine* 34.1 (2017), pp. 67–81.
- [206] Ruhi Sarikaya et al. “An overview of end-to-end language understanding and dialog management for personal digital assistants”. In: *IEEE Workshop on Spoken Language Technology*. 2016.

- [207] OE Scharenborg and MP Cooke. “Comparing human and machine recognition performance on a VCV corpus”. In: *ISCA Tutorial and Research Workshop (ITRW) on Speech Analysis and Processing for Knowledge Discovery* (2008).
- [208] Konrad Scheffler and Steve Young. “Automatic learning of dialogue strategy using dialogue simulation and reinforcement learning”. In: *Proceedings of the second international conference on Human Language Technology Research*. Morgan Kaufmann Publishers Inc. 2002, pp. 12–19.
- [209] Roman Schlegel et al. “Soundcomber: A Stealthy and Context-Aware Sound Trojan for Smartphones.” In: *NDSS*. Vol. 11. 2011, pp. 17–33.
- [210] Martin A Schneider, Marc-Florian Wendland, and Andreas Hoffmann. “A Negative Input Space Complexity Metric as Selection Criterion for Fuzz Testing”. In: *IFIP International Conference on Testing Software and Systems*. Springer. 2015, pp. 257–262.
- [211] Lea Schönherr et al. “Adversarial Attacks Against Automatic Speech Recognition Systems via Psychoacoustic Hiding”. In: *arXiv preprint arXiv:1808.05665* (2018).
- [212] Juliana Schroeder and Matthew Schroeder. “Trusting in Machines: How Mode of Interaction Affects Willingness to Share Personal Information with Machines”. In: *Proceedings of the 51st Hawaii International Conference on System Sciences*. 2018.
- [213] Björn W Schuller. “Speech emotion recognition: two decades in a nutshell, benchmarks, and ongoing trends”. In: *Communications of the ACM* 61.5 (2018), pp. 90–99.
- [214] Iulian Vlad Serban et al. “Building End-To-End Dialogue Systems Using Generative Hierarchical Neural Network Models.” In: *AAAI*. Vol. 16. 2016, pp. 3776–3784.
- [215] William Seymour. “How loyal is your Alexa?: Imagining a Respectful Smart Assistant”. In: *Extended Abstracts of the 2018 CHI Conference on Human Factors in Computing Systems*. ACM. 2018, SRC20.
- [216] William Seymour. “Privacy Therapy with Aretha: What If Your Firewall Could Talk?” In: *Extended Abstracts of the 2019 CHI Conference on Human Factors in Computing Systems*. ACM. 2019, SRC12.
- [217] Nigel Shadbolt, Tim Berners-Lee, and Wendy Hall. “The semantic web revisited”. In: *IEEE Intelligent Systems* 21.3 (2006), pp. 96–101.
- [218] Ali Shafahi et al. “Are adversarial examples inevitable?” In: *arXiv preprint arXiv:1809.02104* (2018).
- [219] CE Shannon. “A Mathematical Theory of Communication”. In: *Bell system technical journal* 27 (1948).
- [220] Mahmood Sharif et al. “Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition”. In: *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*. ACM. 2016, pp. 1528–1540.

- [221] Hongyuan Shen. “Semantic Parsing in Spoken Language Understanding using Abstract Meaning Representation”. PhD thesis. Brandeis University, 2018.
- [222] Elizabeth Shriberg et al. “Learning when to listen: Detecting system-addressed speech in human-human-computer dialog”. In: *Thirteenth Annual Conference of the International Speech Communication Association*. 2012.
- [223] Liwei Song and Prateek Mittal. “Inaudible Voice Commands”. In: *arXiv preprint arXiv:1708.07238* (2017).
- [224] Arjen Stolk, Lennart Verhagen, and Ivan Toni. “Conceptual alignment: how brains achieve mutual understanding”. In: *Trends in Cognitive Sciences* 20.3 (2016), pp. 180–191.
- [225] Grant P. Strimel, Kanthashree Mysore Sathyendra, and Stanislav Peshterliev. “Statistical Model Compression for Small-Footprint Natural Language Understanding”. In: *Interspeech 2018* (2018).
- [226] Takeshi Sugawara et al. “Light Commands: Laser-Based Audio Injection Attacks on Voice-Controllable Systems”.
- [227] Christian Szegedy et al. “Intriguing properties of neural networks”. In: *arXiv preprint arXiv:1312.6199* (2013).
- [228] Andrew Trask, Phil Michalak, and John Liu. “sense2vec-A fast and accurate method for word sense disambiguation in neural word embeddings”. In: *arXiv preprint arXiv:1511.06388* (2015).
- [229] Dimitris Tsipras et al. “There Is No Free Lunch In Adversarial Robustness (But There Are Unexpected Benefits)”. In: *arXiv preprint arXiv:1805.12152* (2018).
- [230] Yulia Tsvetkov. “Linguistic Knowledge in Data-Driven Natural Language Processing”. PhD thesis. Georgia Institute of Technology, 2016.
- [231] Milena Tsvetkova et al. “Even good bots fight: The case of Wikipedia”. In: *PloS one* 12.2 (2017), e0171774.
- [232] Gokhan Tur, Anoop Deoras, and Dilek Hakkani-Tür. “Detecting out-of-domain utterances addressed to a virtual personal assistant”. In: *Fifteenth Annual Conference of the International Speech Communication Association*. 2014.
- [233] Alan M Turing. “Computing machinery and intelligence”. In: *Mind* 59.236 (1950), pp. 433–460.
- [234] Henry Turner, Giulio Lovisotto, and Ivan Martinovic. “Attacking Speaker Recognition Systems with Phoneme Morphing”. In: *European Symposium on Research in Computer Security*. Springer. 2019, pp. 471–492.
- [235] Tavish Vaidya et al. “Cocaine noodles: exploiting the gap between human and machine speech recognition”. In: *Presented at WOOT 15* (2015), pp. 10–11.
- [236] Aäron Van Den Oord et al. “Wavenet: A generative model for raw audio”. In: *CoRR abs/1609.03499* (2016).

- [237] Vedran Vukotic, Christian Raymond, and Guillaume Gravier. “Is it time to switch to word embedding and recurrent neural networks for spoken language understanding?” In: *Interspeech*. 2015.
- [238] David J Weller-Fahy, Brett J Borghetti, and Angela A Sodemann. “A survey of distance and similarity measures used within network intrusion anomaly detection”. In: *IEEE Communications Surveys & Tutorials* 17.1 (2015), pp. 70–91.
- [239] Ryen W White. “Skill discovery in virtual assistants”. In: *Communications of the ACM* 61.11 (2018), pp. 106–113.
- [240] Derry Tanti Wijaya and Reyyan Yeniterzi. “Understanding semantic change of words over centuries”. In: *Proceedings of the 2011 international workshop on Detecting and Exploiting Cultural Diversity on the social web*. ACM. 2011, pp. 35–40.
- [241] Terry Winograd. “Understanding natural language”. In: *Cognitive psychology* 3.1 (1972), pp. 1–191.
- [242] Yonghui Wu et al. “Google’s neural machine translation system: Bridging the gap between human and machine translation”. In: *arXiv preprint arXiv:1609.08144* (2016).
- [243] Zhizheng Wu et al. “Spoofing and countermeasures for speaker verification: a survey”. In: *Speech Communication* 66 (2015), pp. 130–153.
- [244] Wayne Xiong et al. “Achieving human parity in conversational speech recognition”. In: *arXiv preprint arXiv:1610.05256* (2016).
- [245] Hiromu Yakura and Jun Sakuma. “Robust Audio Adversarial Example for a Physical Attack”. In: *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19*. International Joint Conferences on Artificial Intelligence Organization, July 2019, pp. 5334–5341. DOI: 10.24963/ijcai.2019/741. URL: <https://doi.org/10.24963/ijcai.2019/741>.
- [246] Jie Yang et al. “Leveraging Crowdsourcing Data for Deep Active Learning An Application”. In: *Proceedings of the 2018 World Wide Web Conference on World Wide Web - WWW '18* (2018).
- [247] Kaisheng Yao et al. “Recurrent neural networks for language understanding.” In: *Interspeech*. 2013, pp. 2524–2528.
- [248] Park Joon Young et al. “BadVoice: Soundless voice-control replay attack on modern smartphones”. In: *Ubiquitous and Future Networks (ICUFN), 2016 Eighth International Conference on*. IEEE. 2016, pp. 882–887.
- [249] Steve Young et al. “POMDP-based statistical spoken dialog systems: A review”. In: *Proceedings of the IEEE* 101.5 (2013), pp. 1160–1179.
- [250] Tom Young et al. “Recent trends in deep learning based natural language processing”. In: *IEEE Computational Intelligence magazine* 13.3 (2018), pp. 55–75.
- [251] Xuejing Yuan et al. “CommanderSong: A Systematic Approach for Practical Adversarial Voice Recognition”. In: *arXiv preprint arXiv:1801.08535* (2018).

- [252] Guoming Zhang et al. “DolphinAttack: Inaudible Voice Commands”. In: *arXiv preprint arXiv:1708.09537* (2017).
- [253] Jiajie Zhang, Bingsheng Zhang, and Bincheng Zhang. “Defending Adversarial Attacks on Cloud-aided Automatic Speech Recognition Systems”. In: *Proceedings of the Seventh International Workshop on Security in Cloud Computing*. ACM. 2019, pp. 23–31.
- [254] Linghan Zhang et al. “Voicelive: A phoneme localization based liveness detection for voice authentication on smartphones”. In: *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*. ACM. 2016, pp. 1080–1091.
- [255] Nan Zhang et al. “Dangerous skills: Understanding and mitigating security risks of voice-controlled third-party functions on virtual personal assistant systems”. In: *Dangerous Skills: Understanding and Mitigating Security Risks of Voice-Controlled Third-Party Functions on Virtual Personal Assistant Systems*. IEEE. 2019, p. 0.
- [256] Rongjunchen Zhang et al. “Who Activated My Voice Assistant? A Stealthy Attack on Android Phones Without Users’ Awareness”. In: *International Conference on Machine Learning for Cyber Security*. Springer. 2019, pp. 378–396.
- [257] Yangyong Zhang et al. “Life after Speech Recognition: Fuzzing Semantic Misinterpretation for Voice Assistant Applications.” In: *NDSS*. 2019.
- [258] Kateryna Zinchenko, Chien-Yu Wu, and Kai-Tai Song. “A study on speech recognition control for a surgical robot”. In: *IEEE Transactions on Industrial Informatics* 13.2 (2017), pp. 607–615.

Appendix A

Nonsense Attacks on Google Assistant

A.1 Nonsense Words Sets

Nonsense words lists for hey google :
Number of nonsense words for hey : 17

t'eI
Z'eI
v'eI
z'eI
tS'eI
T'eI
S'eI
bl'eI
gl'eI
kw'eI
sk'eI
sm'eI
sn'eI
tw'eI
skr'eI
skw'eI
Tr'eI

Number of nonsense words for google : 395

p'Ud@L
t'Ud@L
k'Ud@L
b'Ud@L
d'Ud@L
dz'Ud@L
g'Ud@L
f'Ud@L
s'Ud@L
z'Ud@L
h'Ud@L
v'Ud@L
z'Ud@L
m'Ud@L
l'Ud@L
j'Ud@L
n'Ud@L
r'Ud@L
w'Ud@L
tS'Ud@L
T'Ud@L
D'Ud@L
S'Ud@L
bl'Ud@L
br'Ud@L
kl'Ud@L
kr'Ud@L
dr'Ud@L
fl'Ud@L
fr'Ud@L
gl'Ud@L
gr'Ud@L
pl'Ud@L
pr'Ud@L
kw'Ud@L
sk'Ud@L
sl'Ud@L
sm'Ud@L

sn'Ud@L
sp'Ud@L
st'Ud@L
sw'Ud@L
tr'Ud@L
tw'Ud@L
spr'Ud@L
spl'Ud@L
skr'Ud@L
skw'Ud@L
str'Ud@L
Tr'Ud@L
p'u:b@L
t'u:b@L
k'u:b@L
b'u:b@L
d'u:b@L
g'u:b@L
f'u:b@L
s'u:b@L
Z'u:b@L
h'u:b@L
v'u:b@L
z'u:b@L
m'u:b@L
l'u:b@L
j'u:b@L
n'u:b@L
w'u:b@L
tS'u:b@L
T'u:b@L
D'u:b@L
S'u:b@L
bl'u:b@L
br'u:b@L
kl'u:b@L
kr'u:b@L
dr'u:b@L
fl'u:b@L
fr'u:b@L
gl'u:b@L
gr'u:b@L
pl'u:b@L
pr'u:b@L
kw'u:b@L
sk'u:b@L
sl'u:b@L
sm'u:b@L
sn'u:b@L
sp'u:b@L
st'u:b@L
sw'u:b@L
tr'u:b@L
tw'u:b@L
spr'u:b@L
spl'u:b@L
skr'u:b@L
skw'u:b@L
str'u:b@L
Tr'u:b@L
p'u:f@L
t'u:f@L
k'u:f@L
b'u:f@L
d'u:f@L
dZ'u:f@L
g'u:f@L
f'u:f@L
s'u:f@L
Z'u:f@L
h'u:f@L
v'u:f@L
z'u:f@L
m'u:f@L
l'u:f@L
j'u:f@L
n'u:f@L
w'u:f@L
tS'u:f@L
T'u:f@L
D'u:f@L
S'u:f@L
bl'u:f@L
br'u:f@L
kl'u:f@L
kr'u:f@L
dr'u:f@L
fl'u:f@L
fr'u:f@L
gl'u:f@L
gr'u:f@L
pl'u:f@L

pr'u:f@L
kw'u:f@L
sk'u:f@L
sl'u:f@L
sm'u:f@L
sn'u:f@L
sp'u:f@L
st'u:f@L
sw'u:f@L
tr'u:f@L
tw'u:f@L
spr'u:f@L
spl'u:f@L
skr'u:f@L
skw'u:f@L
str'u:f@L
Tr'u:f@L
p'Uk@L
t'Uk@L
k'Uk@L
b'Uk@L
d'Uk@L
dz'Uk@L
g'Uk@L
f'Uk@L
s'Uk@L
Z'Uk@L
h'Uk@L
v'Uk@L
z'Uk@L
m'Uk@L
l'Uk@L
j'Uk@L
n'Uk@L
r'Uk@L
w'Uk@L
tS'Uk@L
T'Uk@L
D'Uk@L
S'Uk@L
bl'Uk@L
br'Uk@L
kl'Uk@L
kr'Uk@L
dx'Uk@L
fl'Uk@L
fr'Uk@L
gl'Uk@L
gr'Uk@L
pl'Uk@L
pr'Uk@L
kw'Uk@L
sk'Uk@L
sl'Uk@L
sm'Uk@L
sn'Uk@L
sp'Uk@L
st'Uk@L
sw'Uk@L
tr'Uk@L
tw'Uk@L
spr'Uk@L
spl'Uk@L
skr'Uk@L
skw'Uk@L
str'Uk@L
Tr'Uk@L
p'u:p@L
t'u:p@L
k'u:p@L
b'u:p@L
d'u:p@L
dz'u:p@L
g'u:p@L
f'u:p@L
s'u:p@L
Z'u:p@L
h'u:p@L
v'u:p@L
z'u:p@L
m'u:p@L
l'u:p@L
j'u:p@L
n'u:p@L
r'u:p@L
w'u:p@L
tS'u:p@L
T'u:p@L
D'u:p@L
S'u:p@L
bl'u:p@L
br'u:p@L

kl'u:p@L
kr'u:p@L
dr'u:p@L
fl'u:p@L
fr'u:p@L
gl'u:p@L
gr'u:p@L
pl'u:p@L
pr'u:p@L
kw'u:p@L
sk'u:p@L
sl'u:p@L
sm'u:p@L
sn'u:p@L
sp'u:p@L
st'u:p@L
sw'u:p@L
tr'u:p@L
tw'u:p@L
spr'u:p@L
spl'u:p@L
skw'u:p@L
str'u:p@L
Tr'u:p@L
p'u:s@L
t'u:s@L
k'u:s@L
b'u:s@L
d'u:s@L
dZ'u:s@L
g'u:s@L
f'u:s@L
s'u:s@L
Z'u:s@L
h'u:s@L
v'u:s@L
z'u:s@L
m'u:s@L
l'u:s@L
j'u:s@L
n'u:s@L
r'u:s@L
w'u:s@L
tS'u:s@L
T'u:s@L
D'u:s@L
S'u:s@L
bl'u:s@L
br'u:s@L
kl'u:s@L
kr'u:s@L
dr'u:s@L
fl'u:s@L
fr'u:s@L
gl'u:s@L
gr'u:s@L
pl'u:s@L
pr'u:s@L
kw'u:s@L
sk'u:s@L
sl'u:s@L
sm'u:s@L
sn'u:s@L
sp'u:s@L
st'u:s@L
sw'u:s@L
tr'u:s@L
tw'u:s@L
spr'u:s@L
spl'u:s@L
skr'u:s@L
skw'u:s@L
str'u:s@L
Tr'u:s@L
p'u:t@L
t'u:t@L
k'u:t@L
b'u:t@L
d'u:t@L
dZ'u:t@L
g'u:t@L
f'u:t@L
s'u:t@L
Z'u:t@L
h'u:t@L
v'u:t@L
z'u:t@L
m'u:t@L
l'u:t@L
j'u:t@L
n'u:t@L
r'u:t@L

w'u:t@L
tS'u:t@L
T'u:t@L
D'u:t@L
S'u:t@L
bl'u:t@L
kl'u:t@L
kr'u:t@L
dr'u:t@L
fl'u:t@L
fr'u:t@L
gl'u:t@L
gr'u:t@L
pl'u:t@L
pr'u:t@L
kw'u:t@L
sk'u:t@L
sl'u:t@L
sm'u:t@L
sn'u:t@L
sp'u:t@L
st'u:t@L
sw'u:t@L
tr'u:t@L
tw'u:t@L
spr'u:t@L
spl'u:t@L
skr'u:t@L
skw'u:t@L
str'u:t@L
Tr'u:t@L
p'u:z@L
t'u:z@L
k'u:z@L
b'u:z@L
d'u:z@L
dZ'u:z@L
g'u:z@L
f'u:z@L
s'u:z@L
Z'u:z@L
h'u:z@L
v'u:z@L
z'u:z@L
m'u:z@L
l'u:z@L
j'u:z@L
n'u:z@L
r'u:z@L
w'u:z@L
tS'u:z@L
T'u:z@L
D'u:z@L
S'u:z@L
bl'u:z@L
br'u:z@L
kl'u:z@L
kr'u:z@L
dr'u:z@L
fl'u:z@L
fr'u:z@L
gl'u:z@L
gr'u:z@L
pl'u:z@L
pr'u:z@L
kw'u:z@L
sk'u:z@L
sl'u:z@L
sm'u:z@L
sn'u:z@L
sp'u:z@L
st'u:z@L
sw'u:z@L
tr'u:z@L
tw'u:z@L
spr'u:z@L
spl'u:z@L
skr'u:z@L
skw'u:z@L
str'u:z@L
Tr'u:z@L

Nonsense words lists for who am i :
Number of nonsense words for who : 18

f'u:
Z'u:
v'u:
T'u:
D'u:
fr'u:
pl'u:

pr'u:
kw'u:
sk'u:
sm'u:
sn'u:
sp'u:
sw'u:
tw'u:
spr'u:
spl'u:
skw'u:

Number of nonsense words for am : 27

b'am
g'am
f'am
s'am
Z'am
v'am
z'am
m'am
tS'am
T'am
D'am
bl'am
br'am
fl'am
fr'am
gl'am
pl'am
kw'am
sm'am
sn'am
st'am
tw'am
spr'am
spl'am
skw'am
str'am
Tr'am

Number of nonsense words for i : 20

k'aI
dZ'aI
Z'aI
z'aI
j'aI
bl'aI
br'aI
kl'aI
gl'aI
gr'aI
kw'aI
sm'aI
sn'aI
sw'aI
tw'aI
spl'aI
skr'aI
skw'aI
str'aI
Tr'aI

Nonsense words lists for whats my name :
Number of nonsense words for whats : 27

b'0ts
g'0ts
f'0ts
Z'0ts
h'0ts
v'0ts
z'0ts
tS'0ts
T'0ts
D'0ts
br'0ts
kr'0ts
dr'0ts
fl'0ts
fr'0ts
gl'0ts
gr'0ts
pr'0ts
kw'0ts
sm'0ts
st'0ts
tw'0ts
spr'0ts
spl'0ts

skr'0ts
str'0ts
Tr'0ts

Number of nonsense words for my : 20

k'aI
dZ'aI
Z'aI
z'aI
j'aI
bl'aI
br'aI
kl'aI
gl'aI
gr'aI
kw'aI
sm'aI
sn'aI
sw'aI
tw'aI
spl'aI
skr'aI
skw'aI
str'aI
Tr'aI

Number of nonsense words for name : 35

p'eIm
b'eIm
Z'eIm
h'eIm
v'eIm
z'eIm
j'eIm
r'eIm
w'eIm
tS'eIm
T'eIm
D'eIm
br'eIm
kr'eIm
dr'eIm
gl'eIm
gr'eIm
pl'eIm
pr'eIm
kw'eIm
sk'eIm
sl'eIm
sm'eIm
sn'eIm
sp'eIm
st'eIm
sw'eIm
tr'eIm
tw'eIm
spr'eIm
spl'eIm
skr'eIm
skw'eIm
str'eIm
Tr'eIm

Nonsense words lists for turn on light :
Number of nonsense words for turn : 40

p'3:n
d'3:n
dZ'3:n
g'3:n
s'3:n
Z'3:n
h'3:n
z'3:n
m'3:n
n'3:n
r'3:n
w'3:n
T'3:n
D'3:n
S'3:n
bl'3:n
br'3:n
kl'3:n
kr'3:n
dr'3:n
fl'3:n
fr'3:n
gl'3:n
gr'3:n

pl'3:n
pr'3:n
kw'3:n
sk'3:n
sl'3:n
sm'3:n
sn'3:n
sw'3:n
tr'3:n
tw'3:n
spr'3:n
spl'3:n
skr'3:n
skw'3:n
str'3:n
Tr'3:n

Number of nonsense words for on : 38

p'0n
t'0n
f'0n
s'0n
Z'0n
h'0n
v'0n
z'0n
m'0n
tS'0n
T'0n
D'0n
bl'0n
br'0n
kl'0n
kr'0n
dr'0n
fl'0n
fr'0n
gl'0n
gr'0n
pl'0n
pr'0n
kw'0n
sk'0n
sl'0n
sm'0n
sn'0n
sp'0n
st'0n
tr'0n
tw'0n
spr'0n
spl'0n
skr'0n
skw'0n
str'0n
Tr'0n

Number of nonsense words for light : 28

p'aIt
d'aIt
dZ'aIt
g'aIt
Z'aIt
v'aIt
z'aIt
j'aIt
tS'aIt
T'aIt
D'aIt
S'aIt
kl'aIt
kr'aIt
dr'aIt
gl'aIt
gr'aIt
pr'aIt
sk'aIt
sn'aIt
st'aIt
sw'aIt
tw'aIt
spl'aIt
skr'aIt
skw'aIt
str'aIt
Tr'aIt

Nonsense words lists for turn off light :

Number of nonsense words for turn : 40
Number of nonsense words for off : 41

p'Of
t'Of
b'Of
dZ'Of
f'Of
s'Of
Z'Of
v'Of
z'Of
m'Of
l'Of
j'Of
n'Of
r'Of
w'Of
tS'Of
T'Of
D'Of
S'Of
bl'Of
br'Of
kl'Of
kr'Of
dr'Of
fl'Of
fr'Of
gl'Of
gr'Of
pl'Of
sl'Of
sm'Of
sn'Of
sp'Of
st'Of
sw'Of
tw'Of
spr'Of
spl'Of
skr'Of
skw'Of
Tr'Of

Number of nonsense words for light : 28

Nonsense words lists for turn light red :

Number of nonsense words for turn : 40
Number of nonsense words for light : 28
Number of nonsense words for red : 25

p'Ed
k'Ed
g'Ed
Z'Ed
v'Ed
m'Ed
j'Ed
tS'Ed
T'Ed
D'Ed
kl'Ed
kr'Ed
gl'Ed
gr'Ed
pr'Ed
kw'Ed
sk'Ed
sm'Ed
sn'Ed
sw'Ed
tw'Ed
spl'Ed
skr'Ed
skw'Ed
str'Ed

Nonsense words lists for turn light blue :

Number of nonsense words for turn : 40
Number of nonsense words for light : 28
Number of nonsense words for blue : 18

f'u:
Z'u:
v'u:
T'u:
D'u:
fr'u:

pl'u:
pr'u:
kw'u:
sk'u:
sm'u:
sn'u:
sp'u:
sw'u:
tw'u:
spr'u:
spl'u:
skw'u:

A.2 Pilot Experiment Results

A.2.1 Audio File Input Results (Successes)

```
Wakeup word triggered by nonsense_wakeup/Tr'eI b'u:s@L.raw, nonsense_wakeup/Tr'eI b'u:s@L
INFO:root:Connecting to embeddedassistant.googleapis.com
WARNING:root:error opening WAV file: file does not start with RIFF id, falling back to RAW format
INFO:root:Using device model targetmodelassistant-84514-targetmodellight-el397o and device id 0363ada2-4e16-11e8-8561-000c299218c9
INFO:root:Recording audio request.
INFO:root:Transcript of user request: "my".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "raise".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "weather".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "remove".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "hey Google".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "hey Google".
INFO:root:Playing assistant response.
INFO:root:End of audio request detected
INFO:root:Transcript of user request: "hi Google".
INFO:root:Playing assistant response.
INFO:root:Expecting follow-on query from user.
INFO:root:Hi
What can I do for you?
INFO:root:Finished playing assistant response.
```

```
Wakeup word triggered by nonsense_wakeup/t'eI g'u:t@L.raw, nonsense_wakeup/t'eI g'u:t@L
INFO:root:Connecting to embeddedassistant.googleapis.com
WARNING:root:error opening WAV file: file does not start with RIFF id, falling back to RAW format
INFO:root:Using device model targetmodelassistant-84514-targetmodellight-el397o and device id 0363ada2-4e16-11e8-8561-000c299218c9
INFO:root:Recording audio request.
INFO:root:Transcript of user request: "turn".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "take".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "tegu".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "hey Google".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "hey Google".
INFO:root:Playing assistant response.
INFO:root:End of audio request detected
INFO:root:Transcript of user request: "hey Google".
INFO:root:Playing assistant response.
INFO:root:Expecting follow-on query from user.
INFO:root:Hi
What can I do for you?
INFO:root:Finished playing assistant response.
```

```
Wakeup word triggered by nonsense_wakeup/S'eI h'Uk@L.raw, nonsense_wakeup/S'eI h'Uk@L
INFO:root:Connecting to embeddedassistant.googleapis.com
WARNING:root:error opening WAV file: file does not start with RIFF id, falling back to RAW format
INFO:root:Using device model targetmodelassistant-84514-targetmodellight-el397o and device id 0363ada2-4e16-11e8-8561-000c299218c9
INFO:root:Recording audio request.
INFO:root:Transcript of user request: "sing".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "single".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "Facebook".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "hey Google".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "hey Google".
INFO:root:Playing assistant response.
INFO:root:End of audio request detected
INFO:root:Transcript of user request: "ok Google".
INFO:root:Playing assistant response.
INFO:root:Expecting follow-on query from user.
INFO:root:I'm listening
What's up?
INFO:root:Finished playing assistant response.
```

```
Wakeup word triggered by nonsense_wakeup/tS'eI d'u:s@L.raw, nonsense_wakeup/tS'eI d'u:s@L
INFO:root:Connecting to embeddedassistant.googleapis.com
WARNING:root:error opening WAV file: file does not start with RIFF id, falling back to RAW format
INFO:root:Using device model targetmodelassistant-84514-targetmodellight-el397o and device id 0363ada2-4e16-11e8-8561-000c299218c9
INFO:root:Recording audio request.
INFO:root:Transcript of user request: "show".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "change".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "radio".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "reduce".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "hey Google".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "hey Google".
```

```

INFO:root:Playing assistant response.
INFO:root:End of audio request detected
INFO:root:Transcript of user request: "hey Google".
INFO:root:Playing assistant response.
INFO:root:Expecting follow-on query from user.
INFO:root:Hi
What can I do for you?
INFO:root:Finished playing assistant response.

Wakeup word triggered by nonsense_wakeup/tw'eI g'u:f@L.raw, nonsense_wakeup/tw'eI g'u:f@L
INFO:root:Connecting to embeddedassistant.googleapis.com
WARNING:root:error opening WAV file: file does not start with RIFF id, falling back to RAW format
INFO:root:Using device model targetmodelassistant-84514-targetmodellight-el397o and device id 0363ada2-4e16-11e8-8561-000c299218c9
INFO:root:Recording audio request.
INFO:root:Transcript of user request: "wake".
INFO:root:Transcript of user request: "hey Google".
INFO:root:Transcript of user request: "hey Google".
INFO:root:End of audio request detected.
INFO:root:Stopping recording.
INFO:root:Transcript of user request: "hey google".
INFO:root:Expecting follow-on query from user.
INFO:root:Playing assistant response.
INFO:root:Finished playing assistant response.

Wakeup word triggered by nonsense_wakeup/t'eI D'u:b@L.raw, nonsense_wakeup/t'eI D'u:b@L
INFO:root:Connecting to embeddedassistant.googleapis.com
WARNING:root:error opening WAV file: file does not start with RIFF id, falling back to RAW format
INFO:root:Using device model targetmodelassistant-84514-targetmodellight-el397o and device id 0363ada2-4e16-11e8-8561-000c299218c9
INFO:root:Recording audio request.
INFO:root:Transcript of user request: "time".
INFO:root:Transcript of user request: "change the".
INFO:root:Transcript of user request: "time loop".
INFO:root:Transcript of user request: "hey Google".
INFO:root:Transcript of user request: "hey Google".
INFO:root:End of audio request detected.
INFO:root:Stopping recording.
INFO:root:Transcript of user request: "hey Google".
INFO:root:Expecting follow-on query from user.
INFO:root:Playing assistant response.
INFO:root:Finished playing assistant response.

Wakeup word triggered by nonsense_wakeup/tS'eI d'u:t@L.raw, nonsense_wakeup/tS'eI d'u:t@L
INFO:root:Connecting to embeddedassistant.googleapis.com
WARNING:root:error opening WAV file: file does not start with RIFF id, falling back to RAW format
INFO:root:Using device model targetmodelassistant-84514-targetmodellight-el397o and device id 0363ada2-4e16-11e8-8561-000c299218c9
INFO:root:Recording audio request.
INFO:root:Transcript of user request: "show".
INFO:root:Transcript of user request: "change".
INFO:root:Transcript of user request: "show do".
INFO:root:Transcript of user request: "hey dude".
INFO:root:Transcript of user request: "hey Google".
INFO:root:Transcript of user request: "hey Google".
INFO:root:Transcript of user request: "hey Google".
INFO:root:End of audio request detected.
INFO:root:Stopping recording.
INFO:root:Transcript of user request: "hey Google".
INFO:root:Expecting follow-on query from user.
INFO:root:Playing assistant response.
INFO:root:Finished playing assistant response.

Wakeup word triggered by nonsense_wakeup/Z'eI j'u:b@L.raw, nonsense_wakeup/Z'eI j'u:b@L
INFO:root:Connecting to embeddedassistant.googleapis.com
WARNING:root:error opening WAV file: file does not start with RIFF id, falling back to RAW format
INFO:root:Using device model targetmodelassistant-84514-targetmodellight-el397o and device id 0363ada2-4e16-11e8-8561-000c299218c9
INFO:root:Recording audio request.
INFO:root:Transcript of user request: "show you".
INFO:root:Transcript of user request: "Jaguar".
INFO:root:Transcript of user request: "hey Google".
INFO:root:Transcript of user request: "hey Google".
INFO:root:End of audio request detected.
INFO:root:Stopping recording.
INFO:root:Transcript of user request: "hey Google".
INFO:root:Expecting follow-on query from user.
INFO:root:Playing assistant response.
INFO:root:Finished playing assistant response.

Wakeup word triggered by nonsense_wakeup/S'eI j'u:b@L.raw, nonsense_wakeup/S'eI j'u:b@L
INFO:root:Connecting to embeddedassistant.googleapis.com
WARNING:root:error opening WAV file: file does not start with RIFF id, falling back to RAW format
INFO:root:Using device model targetmodelassistant-84514-targetmodellight-el397o and device id 0363ada2-4e16-11e8-8561-000c299218c9
INFO:root:Recording audio request.
INFO:root:Transcript of user request: "show you".
INFO:root:Transcript of user request: "hey Google".
INFO:root:Transcript of user request: "hey Google".
INFO:root:End of audio request detected.
INFO:root:Stopping recording.
INFO:root:Transcript of user request: "hey Google".
INFO:root:Playing assistant response.
INFO:root:Expecting follow-on query from user.
INFO:root:Finished playing assistant response.

Wakeup word triggered by nonsense_wakeup/z'eI g'u:p@L.raw, nonsense_wakeup/z'eI g'u:p@L

```

```

INFO:root:Connecting to embeddedassistant.googleapis.com
WARNING:root:error opening WAV file: file does not start with RIFF id, falling back to RAW format
INFO:root:Using device model targetmodelassistant-84514-targetmodellight-el397o and device id 0363ada2-4e16-11e8-8561-000c299218c9
INFO:root:Recording audio request.
INFO:root:Transcript of user request: "play".
INFO:root:Transcript of user request: "make".
INFO:root:Transcript of user request: "degu".
INFO:root:Transcript of user request: "hey Google".
INFO:root:Transcript of user request: "hey Google".
INFO:root:End of audio request detected.
INFO:root:Stopping recording.
INFO:root:Transcript of user request: "hey Google".
INFO:root:Expecting follow-on query from user.
INFO:root:Playing assistant response.
INFO:root:Finished playing assistant response.

Wakeup word triggered by nonsense_wakeup/Z'eI g'u:p@L.raw, nonsense_wakeup/Z'eI g'u:p@L
INFO:root:Connecting to embeddedassistant.googleapis.com
WARNING:root:error opening WAV file: file does not start with RIFF id, falling back to RAW format
INFO:root:Using device model targetmodelassistant-84514-targetmodellight-el397o and device id 0363ada2-4e16-11e8-8561-000c299218c9
INFO:root:Recording audio request.
INFO:root:Transcript of user request: "hey".
INFO:root:Transcript of user request: "Jake".
INFO:root:Transcript of user request: "hey Google".
INFO:root:Transcript of user request: "hey Google".
INFO:root:End of audio request detected.
INFO:root:Stopping recording.
INFO:root:Transcript of user request: "hey Google".
INFO:root:Expecting follow-on query from user.
INFO:root:Playing assistant response.
INFO:root:Finished playing assistant response.

Wakeup word triggered by nonsense_wakeup/t'eI g'Ud@L.raw, nonsense_wakeup/t'eI g'Ud@L
INFO:root:Connecting to embeddedassistant.googleapis.com
WARNING:root:error opening WAV file: file does not start with RIFF id, falling back to RAW format
INFO:root:Using device model targetmodelassistant-84514-targetmodellight-el397o and device id 0363ada2-4e16-11e8-8561-000c299218c9
INFO:root:Recording audio request.
INFO:root:Transcript of user request: "take".
INFO:root:Transcript of user request: "hey good".
INFO:root:Transcript of user request: "hey Google".
INFO:root:Transcript of user request: "hey Google".
INFO:root:Transcript of user request: "hey Google".
INFO:root:End of audio request detected.
INFO:root:Stopping recording.
INFO:root:Transcript of user request: "hey Google".
INFO:root:Expecting follow-on query from user.
INFO:root:Playing assistant response.
INFO:root:Finished playing assistant response.

Wakeup word triggered by nonsense_wakeup/Z'eI d'u:b@L.raw, nonsense_wakeup/Z'eI d'u:b@L
INFO:root:Connecting to embeddedassistant.googleapis.com
WARNING:root:error opening WAV file: file does not start with RIFF id, falling back to RAW format
INFO:root:Using device model targetmodelassistant-84514-targetmodellight-el397o and device id 0363ada2-4e16-11e8-8561-000c299218c9
INFO:root:Recording audio request.
INFO:root:Transcript of user request: "change".
INFO:root:Transcript of user request: "JD".
INFO:root:Transcript of user request: "hey dude".
INFO:root:Transcript of user request: "hey Google".
INFO:root:Transcript of user request: "hey Google".
INFO:root:Transcript of user request: "hey Google".
INFO:root:End of audio request detected.
INFO:root:Stopping recording.
INFO:root:Transcript of user request: "hey Google".
INFO:root:Expecting follow-on query from user.
INFO:root:Playing assistant response.

INFO:root:Finished playing assistant response.
Wakeup word triggered by nonsense_wakeup/S'eI k'u:b@L.raw, nonsense_wakeup/S'eI k'u:b@L
INFO:root:Connecting to embeddedassistant.googleapis.com
WARNING:root:error opening WAV file: file does not start with RIFF id, falling back to RAW format
INFO:root:Using device model targetmodelassistant-84514-targetmodellight-el397o and device id 0363ada2-4e16-11e8-8561-000c299218c9
INFO:root:Recording audio request.
INFO:root:Transcript of user request: "play".
INFO:root:Transcript of user request: "sing".
INFO:root:Transcript of user request: "Shake".
INFO:root:Transcript of user request: "thank you".
INFO:root:Transcript of user request: "Sheffield".
INFO:root:Transcript of user request: "hey Google".
INFO:root:Transcript of user request: "hey Google".
INFO:root:End of audio request detected.
INFO:root:Stopping recording.
INFO:root:Transcript of user request: "hey Google".
INFO:root:Playing assistant response.
INFO:root:Expecting follow-on query from user.
INFO:root:Finished playing assistant response.

Wakeup word triggered by nonsense_wakeup/S'eI v'u:b@L.raw, nonsense_wakeup/S'eI v'u:b@L
INFO:root:Connecting to embeddedassistant.googleapis.com
WARNING:root:error opening WAV file: file does not start with RIFF id, falling back to RAW format
INFO:root:Using device model targetmodelassistant-84514-targetmodellight-el397o and device id 0363ada2-4e16-11e8-8561-000c299218c9
INFO:root:Recording audio request.
INFO:root:Transcript of user request: "show".

```

INFO:root:Transcript of user request: "save".
INFO:root:Transcript of user request: "save the".
INFO:root:Transcript of user request: "save room".
INFO:root:Transcript of user request: "hey Google".
INFO:root:Transcript of user request: "hey Google".
INFO:root:End of audio request detected.
INFO:root:Stopping recording.
INFO:root:Transcript of user request: "hey Google".
INFO:root:Expecting follow-on query from user.
INFO:root:Playing assistant response.
INFO:root:Finished playing assistant response.

Target command ('turn on light') triggered by nonsense_audio/cyczOm4gWicwbiBqJ2FJdD09PQ==.raw, s'3:n Z'0n j'aIt===
INFO:root:Connecting to embeddedassistant.googleapis.com
WARNING:root:error opening WAV file: file does not start with RIFF id, falling back to RAW format
INFO:root:Using device model targetmodelassistant-84514-targetmodellight-el1397o and device id 0363ada2-4e16-11e8-8561-000c299218c9
INFO:root:Recording audio request.
INFO:root:Transcript of user request: "so".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "send".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn the".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn on the light".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn on the light".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn on the light".
INFO:root:Playing assistant response.
INFO:root:End of audio request detected
INFO:root:Transcript of user request: "turn on the light".
INFO:root:Playing assistant response.
INFO:root:OK
INFO:root:Waiting for device executions to complete.
INFO:root:Turning device on
INFO:root:Finished playing assistant response.

Target command ('turn on light') triggered by nonsense_audio/ZmwnMzpuIG0nMG4ga2wnYU10PT09.raw, fl'3:n m'0n kl'aIt===
INFO:root:Connecting to embeddedassistant.googleapis.com
WARNING:root:error opening WAV file: file does not start with RIFF id, falling back to RAW format
INFO:root:Using device model targetmodelassistant-84514-targetmodellight-el1397o and device id 0363ada2-4e16-11e8-8561-000c299218c9
INFO:root:Recording audio request.
INFO:root:Transcript of user request: "London".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "moncler".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn on light".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn on light".
INFO:root:Playing assistant response.
INFO:root:End of audio request detected
INFO:root:Transcript of user request: "turn on light".
INFO:root:Playing assistant response.
INFO:root:Sure
INFO:root:Waiting for device executions to complete.
INFO:root:Turning device on
INFO:root:Finished playing assistant response.

Target command ('turn off light') triggered by nonsense_audio/ZmwnMzpuIG4nMGYgVCdhSXQ9PT0=.raw, fl'3:n n'0f T'aIt===
INFO:root:Connecting to embeddedassistant.googleapis.com
WARNING:root:error opening WAV file: file does not start with RIFF id, falling back to RAW format
INFO:root:Using device model targetmodelassistant-84514-targetmodellight-el1397o and device id 0363ada2-4e16-11e8-8561-000c299218c9
INFO:root:Recording audio request.
INFO:root:Transcript of user request: "turn off".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn off the".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn off light".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn off light".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn off light".
INFO:root:Playing assistant response.
INFO:root:End of audio request detected
INFO:root:Transcript of user request: "turn off light".
INFO:root:Playing assistant response.
INFO:root:OK
INFO:root:Waiting for device executions to complete.
INFO:root:Turning device off
INFO:root:Finished playing assistant response.

Target command ('turn off light') triggered by nonsense_audio/dHInMzpuIG4nMGYgZ2wnYU10PT09.raw, tr'3:n n'0f gl'aIt===
INFO:root:Connecting to embeddedassistant.googleapis.com
WARNING:root:error opening WAV file: file does not start with RIFF id, falling back to RAW format
INFO:root:Using device model targetmodelassistant-84514-targetmodellight-el1397o and device id 0363ada2-4e16-11e8-8561-000c299218c9
INFO:root:Recording audio request.
INFO:root:Transcript of user request: "show".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "true".
INFO:root:Playing assistant response.

```

INFO:root:Transcript of user request: "should I".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn off".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "should I".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn off the light".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn off the light".
INFO:root:Playing assistant response.
INFO:root:End of audio request detected
INFO:root:Transcript of user request: "turn off the light".
INFO:root:Playing assistant response.
INFO:root:You got it
INFO:root:Turning device off
INFO:root:Waiting for device executions to complete.
INFO:root:Finished playing assistant response.

Target command ('turn off light') triggered by nonsense_audio/aCczOm4geicwZiBqJ2FJdD09PQ==.raw, h'3:n z'0f j'aIt===
INFO:root:Connecting to embeddedassistant.googleapis.com
WARNING:root:error opening WAV file: file does not start with RIFF id, falling back to RAW format
INFO:root:Using device model targetmodelassistant-84514-targetmodellight-el1397o and device id 0363ada2-4e16-11e8-8561-000c299218c9
INFO:root:Recording audio request.
INFO:root:Transcript of user request: "turn".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn the".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turns off".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turns off the".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turns off the light".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turns off the light".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turns off the light".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turns off the light".
INFO:root:Playing assistant response.
INFO:root:End of audio request detected
INFO:root:Transcript of user request: "turns off the light".
INFO:root:Playing assistant response.
INFO:root:You got it
INFO:root:Turning device off
INFO:root:Waiting for device executions to complete.
INFO:root:Finished playing assistant response.

Target command ('turn on light') triggered by nonsense_audio/cCczOm4gaCcwbiBrbCdhSXQ9PT0=.raw, p'3:n h'0n k'l'aIt===
INFO:root:Connecting to embeddedassistant.googleapis.com
WARNING:root:error opening WAV file: file does not start with RIFF id, falling back to RAW format
INFO:root:Using device model targetmodelassistant-84514-targetmodellight-el1397o and device id 0363ada2-4e16-11e8-8561-000c299218c9
INFO:root:Recording audio request.
INFO:root:Transcript of user request: "turn".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn on".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn pond".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn Honda".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn Honda".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn on July".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn on light".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn on blood".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn on blood".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn on blood".
INFO:root:Playing assistant response.
INFO:root:End of audio request detected
INFO:root:Transcript of user request: "turn on light".
INFO:root:Playing assistant response.
INFO:root:OK
INFO:root:Waiting for device executions to complete.
INFO:root:Turning device on
INFO:root:Finished playing assistant response.

Target command ('whats my name') triggered by nonsense_audio/c20nMHRzIHnuJ2FJIHnsJ2VJbT09PQ==.raw, sm'0ts sn'aI sl'eIm===
INFO:root:Connecting to embeddedassistant.googleapis.com
WARNING:root:error opening WAV file: file does not start with RIFF id, falling back to RAW format
INFO:root:Using device model targetmodelassistant-84514-targetmodellight-el1397o and device id 0363ada2-4e16-11e8-8561-000c299218c9
INFO:root:Recording audio request.
INFO:root:Transcript of user request: "what's".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "some".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "summer".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "what's the".

```

```

INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "it's nice".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "it's nicely".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "it's nice name".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "it's nice name".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "it's nice name".
INFO:root:Playing assistant response.
INFO:root:End of audio request detected
INFO:root:Transcript of user request: "what's my name".
INFO:root:Playing assistant response.
INFO:root:You told me your name was MK
I could never forget that
INFO:root:Finished playing assistant response.

Target command ('whats my name') triggered by nonsense_audio/c20nMHRzIGsnYUkGc3AnZUltPT09.raw, sm'0ts k'aI sp'eIm===
INFO:root:Connecting to embeddedassistant.googleapis.com
WARNING:root:error opening WAV file: file does not start with RIFF id, falling back to RAW format
INFO:root:Using device model targetmodelassistant-84514-targetmodellight-el1397o and device id 0363ada2-4e16-11e8-8561-000c299218c9
INFO:root:Recording audio request.
INFO:root:Transcript of user request: "what's".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "some".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "summer".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "what's on Sky".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "what's my IP".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "some months cause pain".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "what's my car's paint".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "what's my car's paint".
INFO:root:Playing assistant response.
INFO:root:End of audio request detected
INFO:root:Transcript of user request: "what's my name".
INFO:root:Playing assistant response.
INFO:root:You told me your name was MK
I could never forget that
INFO:root:Finished playing assistant response.

Target command ('turn light red') triggered by nonsense_audio/dHInMzpuIFRyJ2FJdCBzdHInRWQ9PT0=.raw, tr'3:n Tr'aIt str'Ed===
INFO:root:Connecting to embeddedassistant.googleapis.com
WARNING:root:error opening WAV file: file does not start with RIFF id, falling back to RAW format
INFO:root:Using device model targetmodelassistant-84514-targetmodellight-el1397o and device id 0363ada2-4e16-11e8-8561-000c299218c9
INFO:root:Recording audio request.
INFO:root:Transcript of user request: "show".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "true".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "should I".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "true".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "should I".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn right".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn right to".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "who writes to Red".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "who writes to Red".
INFO:root:Playing assistant response.
INFO:root:End of audio request detected
INFO:root:Transcript of user request: "turn right to Red".
INFO:root:Playing assistant response.
INFO:root:OK
INFO:root:Waiting for device executions to complete.
INFO:root:color is red
INFO:root:Finished playing assistant response.

Target command ('turn light blue') triggered by nonsense_audio/cGwnMzpuIHnWbCdhsXQgc3BsJ3U6PT09.raw, pl'3:n spl'aIt spl'u:===
INFO:root:Connecting to embeddedassistant.googleapis.com
WARNING:root:error opening WAV file: file does not start with RIFF id, falling back to RAW format
INFO:root:Using device model targetmodelassistant-84514-targetmodellight-el1397o and device id 0363ada2-4e16-11e8-8561-000c299218c9
INFO:root:Recording audio request.
INFO:root:Transcript of user request: "play".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "learn".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "let's".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "let's play".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "lounge lights".
INFO:root:Playing assistant response.

```

```

INFO:root:Transcript of user request: "lens flights to".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn lights blue".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn lights blue".
INFO:root:Playing assistant response.
INFO:root:End of audio request detected
INFO:root:Transcript of user request: "turn lights blue".
INFO:root:Playing assistant response.
INFO:root:You got it
INFO:root:color is blue
INFO:root:Waiting for device executions to complete.
INFO:root:Finished playing assistant response.

Target command ('whats my name') triggered by nonsense_audio/ZnInMHRzIHNwCbDhSSBUJ2VJbT09PQ==.raw, fr'0ts spl'ai T'eIm===
INFO:root:Connecting to embeddedassistant.googleapis.com
WARNING:root:error opening WAV file: file does not start with RIFF id, falling back to RAW format
INFO:root:Using device model targetmodelassistant-84514-targetmodellight-el1397o and device id 0363ada2-4e16-11e8-8561-000c299218c9
INFO:root:Recording audio request.
INFO:root:Transcript of user request: "Robert".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "what's".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "what's my".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "Robert slice".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "rugs fluffy".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "what's my phone".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "what's my phone".
INFO:root:Playing assistant response.
INFO:root:End of audio request detected
INFO:root:Transcript of user request: "what's my name".
INFO:root:Playing assistant response.
INFO:root:You told me your name was MK
I could never forget that
INFO:root:Finished playing assistant response.

Target command ('turn light blue') triggered by nonsense_audio/aCczOm4gZ2wnYU10IHNrdyd10j09PQ==.raw, h'3:n gl'aIt skw'u===
INFO:root:Connecting to embeddedassistant.googleapis.com
WARNING:root:error opening WAV file: file does not start with RIFF id, falling back to RAW format
INFO:root:Using device model targetmodelassistant-84514-targetmodellight-el1397o and device id 0363ada2-4e16-11e8-8561-000c299218c9
INFO:root:Recording audio request.
INFO:root:Transcript of user request: "turn".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn Grey".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn black".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn black".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn the lights".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "Turn The Lights crew".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "Turn The Lights crew".
INFO:root:Playing assistant response.
INFO:root:End of audio request detected
INFO:root:Transcript of user request: "turn the lights blue".
INFO:root:Playing assistant response.
INFO:root:OK
INFO:root:color is blue
INFO:root:Waiting for device executions to complete.
INFO:root:Finished playing assistant response.

Target command ('turn light blue') triggered by nonsense_audio/cCczOm4gdidhSXQgc2sndTo9PT0=.raw, p'3:n v'aIt sk'u===
INFO:root:Connecting to embeddedassistant.googleapis.com
WARNING:root:error opening WAV file: file does not start with RIFF id, falling back to RAW format
INFO:root:Using device model targetmodelassistant-84514-targetmodellight-el1397o and device id 0363ada2-4e16-11e8-8561-000c299218c9
INFO:root:Recording audio request.
INFO:root:Transcript of user request: "turn".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn my".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn right".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn lights".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn lights".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn lights to".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn lights to".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn lights to".
INFO:root:Playing assistant response.
INFO:root:End of audio request detected
INFO:root:Transcript of user request: "turn lights blue".

```

```

INFO:root:Playing assistant response.
INFO:root:You got it
INFO:root:Waiting for device executions to complete.
INFO:root:color is blue
INFO:root:Finished playing assistant response.

Target command ('turn light red') triggered by nonsense_audio/cHINmZpuIGonYU0IHN3J0VvkPT09.raw, pr'3:n j'aIt sw'Ed===
INFO:root:Connecting to embeddedassistant.googleapis.com
WARNING:root:error opening WAV file: file does not start with RIFF id, falling back to RAW format
INFO:root:Using device model targetmodelassistant-84514-targetmodellight-el397o and device id 0363ada2-4e16-11e8-8561-000c299218c9
INFO:root:Recording audio request.
INFO:root:Transcript of user request: "turn on the lights".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn the lights red".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn the lights red".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn the lights red".
INFO:root:Playing assistant response.
INFO:root:End of audio request detected
INFO:root:Transcript of user request: "turn the lights red".
INFO:root:Playing assistant response.
INFO:root:You got it
INFO:root:color is red
INFO:root:Waiting for device executions to complete.
INFO:root:Finished playing assistant response.

Target command ('turn light red') triggered by nonsense_audio/c3RyJzM6biBqJ2FJdCBzdHInRWQ9PT0=.raw, str'3:n j'aIt str'Ed===
INFO:root:Connecting to embeddedassistant.googleapis.com
WARNING:root:error opening WAV file: file does not start with RIFF id, falling back to RAW format
INFO:root:Using device model targetmodelassistant-84514-targetmodellight-el397o and device id 0363ada2-4e16-11e8-8561-000c299218c9
INFO:root:Recording audio request.
INFO:root:Transcript of user request: "what's".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "true".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "stream".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "Struan".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "set reminder".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn on the lights".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn on the lights to Red".
INFO:root:Playing assistant response.
INFO:root:Transcript of user request: "turn on the lights to Red".
INFO:root:Playing assistant response.
INFO:root:End of audio request detected
INFO:root:Transcript of user request: "turn the lights to Red".
INFO:root:Playing assistant response.
INFO:root:Sure
INFO:root:Waiting for device executions to complete.
INFO:root:color is red
INFO:root:Finished playing assistant response.

```

A.2.2 Over-the-Air Results

Command 76:

```
t'eI D'u:bl= s'3:n Z'Qn j'alt  
;speak;;phoneme alphabet="x-sampa" ph="t'eI D'u:bl=";i/phoneme;.;phoneme alphabet="x-sampa" ph="s'3:n Z'Qn j'alt";i/phoneme;./speak;
```

```
mkb@ubuntu:~$ source ~/env/bin/activate  
(env) mkb@ubuntu:~$ googlesamples-assistant-hotword --project-id targetmodelassistant-84514 --device-model-id targetmodelassistant-84514-ta  
device_model_id: targetmodelassistant-84514-targetmodellight-el397o  
device_id: ABFFAB9BFEA3B2E141711A75E833E65D  
  
ON_MUTED_CHANGED:  
{ "is_muted": false }  
ON_MEDIA_STATE_IDLE  
ON_START_FINISHED  
  
ON_CONVERSATION_TURN_STARTED  
ON_END_OF_UTTERANCE  
ON_RECOGNIZING_SPEECH_FINISHED:  
{ "text": "switch on the light" }  
ON_DEVICE_ACTION:  
{  
  "inputs": [  
    {  
      "intent": "action.devices.EXECUTE",  
      "payload": {  
        "commands": [  
          {  
            "devices": [  
              {  
                "id": "ABFFAB9BFEA3B2E141711A75E833E65D"  
              }  
            ],  
            "execution": [  
              {  
                "command": "action.devices.commands.OnOff",  
                "params": {  
                  "on": true  
                }  
              }  
            ]  
          }  
        ],  
        "requestId": "4180846555886906104"  
      }  
    }  
  ]  
  Do command action.devices.commands.OnOff with params {u'on': True}  
  Turning the LED on.  
  ON_RESPONDING_STARTED:  
  { "is_error_response": false }  
  ON_RESPONDING_FINISHED  
  ON_CONVERSATION_TURN_FINISHED:  
  { "with_follow_on_turn": false }
```

Command 85:

```
t'eI D'u:bl= tr'3:n Tr'alt str'Ed  
;speak;;phoneme alphabet="x-sampa" ph="t'eI D'u:bl=";i/phoneme;.;phoneme alphabet="x-sampa" ph="tr'3:n Tr'alt str'Ed";i/phoneme;./speak;
```

```
ON_CONVERSATION_TURN_STARTED  
ON_END_OF_UTTERANCE  
ON_RECOGNIZING_SPEECH_FINISHED:  
{ "text": "turn lights to Red" }  
ON_DEVICE_ACTION:  
{  
  "inputs": [  
    {  
      "intent": "action.devices.EXECUTE",  
      "payload": {  
        "commands": [  
          {  
            "devices": [  
              {  
                "id": "ABFFAB9BFEA3B2E141711A75E833E65D"  
              }  
            ],  
            "execution": [  
              {  
                "command": "action.devices.commands.ColorAbsolute",  
                "params": {  
                  "color": {  
                    "name": "red",  
                    "spectrumRGB": 16711680  
                  }  
                }  
              }  
            ]  
          }  
        ],  
        "requestId": "4180846555886906104"  
      }  
    }  
  ]  
  Do command action.devices.commands.ColorAbsolute with params {color: {name: 'red', spectrumRGB: 16711680}}  
  ON_RESPONDING_STARTED:  
  { "is_error_response": false }  
  ON_RESPONDING_FINISHED  
  ON_CONVERSATION_TURN_FINISHED:  
  { "with_follow_on_turn": false }
```

```

]
}
}
],
"requestId": "5469482669336584788"
}
Do command action.devices.commands.ColorAbsolute with params {u'color': {u'spectrumRGB': 16711680, u'name': u'red'}}
The color is red.
ON_RESPONDING_STARTED:
{"is_error_response": false}
ON_RESPONDING_FINISHED
ON_CONVERSATION_TURN_FINISHED:
{"with_follow_on_turn": false}

```

Command 86:

t'eI D'u:bl= pr'3:n j'alt sw'Ed
;speak;phoneme alphabet="x-sampa" ph="t'eI D'u:bl=";i/phoneme;:phoneme alphabet="x-sampa" ph="pr'3:n j'alt sw'Ed";i/phoneme;/speak;

```

ON_CONVERSATION_TURN_STARTED
ON_END_OF_UTTERANCE
ON_RECOGNIZING_SPEECH_FINISHED:
{"text": "turn the lights red"}
ON_DEVICE_ACTION:
{
"inputs": [
{
"intent": "action.devices.EXECUTE",
"payload": {
"commands": [
{
"devices": [
{
"id": "ABFFAB9BFEA3B2E141711A75E833E65D"
}
],
"execution": [
{
"command": "action.devices.commands.ColorAbsolute",
"params": {
"color": {
"name": "red",
"spectrumRGB": 16711680
}
}
}
]
}
]
},
"requestId": "389655376420458629"
}
Do command action.devices.commands.ColorAbsolute with params {u'color': {u'spectrumRGB': 16711680, u'name': u'red'}}
The color is red.
ON_RESPONDING_STARTED:
{"is_error_response": false}
ON_RESPONDING_FINISHED
ON_CONVERSATION_TURN_FINISHED:
{"with_follow_on_turn": false}

```

Command 87:

t'eI D'u:bl= str'3:n j'alt str'Ed
;speak;phoneme alphabet="x-sampa" ph="t'eI D'u:bl=";i/phoneme;:phoneme alphabet="x-sampa" ph="str'3:n j'alt str'Ed";i/phoneme;/speak;

```

ON_CONVERSATION_TURN_STARTED
ON_END_OF_UTTERANCE
ON_RECOGNIZING_SPEECH_FINISHED:
{"text": "turn lights to Red"}
ON_DEVICE_ACTION:
{
"inputs": [
{
"intent": "action.devices.EXECUTE",
"payload": {
"commands": [
{
"devices": [
{
"id": "ABFFAB9BFEA3B2E141711A75E833E65D"
}
],
"execution": [
{
"command": "action.devices.commands.ColorAbsolute",
"params": {
"color": {

```

```
"name": "red",
"spectrumRGB": 16711680
}
}
]
}
}
}
},
"requestId": "2326273977991846014"
}
Do command action.devices.commands.ColorAbsolute with params {u'color': {u'spectrumRGB': 16711680, u'name': u'red'}}
The color is red.
ON_RESPONDING_STARTED:
{"is_error_response": false}
ON_RESPONDING_FINISHED
ON_CONVERSATION_TURN_FINISHED:
{"with_follow_on_turn": false}
```

A.2.3 Human Comprehensibility Results

Participant ID: 5b347a23975e260001e9a0d8
Timestamp: Wed Aug 01 2018 16:31:24 GMT+0100 (British Summer Time)
Native language: English
Participant ID: 5b347a23975e260001e9a0d8
Timestamp: Wed Aug 01 2018 16:31:34 GMT+0100 (British Summer Time)
Meanings assigned to audio clips:
attention1.wav: Hello how are you
attention2.wav: Hi how are you
h3nglaItskwu.wav: NO MEANING
h3nz0fjaIt.wav: NO MEANING
p3nh0nklaIt.wav: NO MEANING
SeIjubL.wav: NO MEANING
sm0tskaIspeIm.wav: NO MEANING
str3njaItstrEd.wav: NO MEANING
teIDublpr3njaItswEd.wav: NO MEANING
teIDubls3nzQnjaIt.wav: NO MEANING
teIDublstr3njaItstrEd.wav: NO MEANING
teIDubltr3nTraItstrEd.wav: NO MEANING
teIgutL.wav: NO MEANING
ZeIdublL.wav: NO MEANING
Participant ID: 5b347a23975e260001e9a0d8
Timestamp: Wed Aug 01 2018 16:33:36 GMT+0100 (British Summer Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5b47b4257356580001f34fc3
Timestamp: Wed Aug 01 2018 10:11:01 GMT-0500 (Central Daylight Time)
Native language: English,Spanish
Participant ID: 5b47b4257356580001f34fc3
Timestamp: Wed Aug 01 2018 10:11:32 GMT-0500 (Central Daylight Time)
Meanings assigned to audio clips:
attention1.wav: hello how are you
attention2.wav: hi how are you
h3nglaItskwu.wav: NO MEANING
h3nz0fjaIt.wav: NO MEANING
p3nh0nklaIt.wav: NO MEANING
SeIjubL.wav: NO MEANING
sm0tskaIspeIm.wav: smoking cause pain
str3njaItstrEd.wav: turn left red
teIDublpr3njaItswEd.wav: NO MEANING
teIDubls3nzQnjaIt.wav: NO MEANING
teIDublstr3njaItstrEd.wav: NO MEANING
teIDubltr3nTraItstrEd.wav: NO MEANING
teIgutL.wav: NO MEANING
ZeIdublL.wav: NO MEANING
Participant ID: 5b47b4257356580001f34fc3
Timestamp: Wed Aug 01 2018 10:14:16 GMT-0500 (Central Daylight Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5ae3837032ef52000131a263
Timestamp: Wed Aug 01 2018 11:08:08 GMT-0400 (Eastern Daylight Time)
Native language: English
Participant ID: 5ae3837032ef52000131a263
Timestamp: Wed Aug 01 2018 11:08:27 GMT-0400 (Eastern Daylight Time)
Meanings assigned to audio clips:
attention1.wav: Hello how are you
attention2.wav: Hi how are you
h3nglaItskwu.wav: NO MEANING
h3nz0fjaIt.wav: NO MEANING
p3nh0nklaIt.wav: NO MEANING
SeIjubL.wav: NO MEANING
sm0tskaIspeIm.wav: NO MEANING
str3njaItstrEd.wav: NO MEANING
teIDublpr3njaItswEd.wav: NO MEANING
teIDubls3nzQnjaIt.wav: NO MEANING
teIDublstr3njaItstrEd.wav: NO MEANING
teIDubltr3nTraItstrEd.wav: NO MEANING
teIgutL.wav: NO MEANING
ZeIdublL.wav: NO MEANING
Participant ID: 5ae3837032ef52000131a263
Timestamp: Wed Aug 01 2018 11:14:03 GMT-0400 (Eastern Daylight Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5b3e37b781b0c00001f7e95e
Timestamp: Wed Aug 01 2018 11:10:15 GMT-0400 (Eastern Daylight Time)
Native language: English
Participant ID: 5b3e37b781b0c00001f7e95e
Timestamp: Wed Aug 01 2018 11:10:28 GMT-0400 (Eastern Daylight Time)
Meanings assigned to audio clips:
attention1.wav: Hello how are you
attention2.wav: Hi how are you
h3nglaItskwu.wav: turn the light screw
h3nz0fjaIt.wav: NO MEANING
p3nh0nklaIt.wav: Turn on the light
SeIjubL.wav: chez yuble
sm0tskaIspeIm.wav: NO MEANING
str3njaItstrEd.wav: NO MEANING

teIDublpr3njaItswEd.wav: NO MEANING
teIDubls3n2QnjaIt.wav: NO MEANING
teIDublstr3njaItstrEd.wav: NO MEANING
teIDubltr3nTraItstrEd.wav: NO MEANING
teIgutL.wav: NO MEANING
ZeIdubl.wav: NO MEANING
Participant ID: 5b3e37b781b0c00001f7e95e
Timestamp: Wed Aug 01 2018 11:13:05 GMT-0400 (Eastern Daylight Time)
Study withdrawn: false
Study submitted: true

Participant ID: 58fdfce6d7c57d000124a72d
Timestamp: Wed Aug 01 2018 16:09:49 GMT+0100 (GMT+01:00)
Native language: English
Participant ID: 58fdfce6d7c57d000124a72d
Timestamp: Wed Aug 01 2018 16:10:02 GMT+0100 (GMT+01:00)
Meanings assigned to audio clips:
attention1.wav: hello, how are you?
attention2.wav: hi, how are you?
h3nglaItskwu.wav: no meaning
h3nz0fjaIt.wav: no meaning
p3nh0nklait.wav: no meaning
SeIjubl.wav: no meaning
sm0tskaIspeIm.wav: no meaning
str3njaItstrEd.wav: no meaning
teIDublpr3njaItswEd.wav: no meaning
teIDubls3n2QnjaIt.wav: no meaning
teIDublstr3njaItstrEd.wav: no meaning
teIDubltr3nTraItstrEd.wav: no meaning
teIgutL.wav: no meaning
ZeIdubl.wav: no meaning
Participant ID: 58fdfce6d7c57d000124a72d
Timestamp: Wed Aug 01 2018 16:12:22 GMT+0100 (GMT+01:00)
Study withdrawn: false
Study submitted: true

Participant ID: 5b4935fa479ca00001ef02d4
Timestamp: Wed Aug 01 2018 16:06:14 GMT+0100 (GMT+01:00)
Native language: English
Participant ID: 5b4935fa479ca00001ef02d4
Timestamp: Wed Aug 01 2018 16:06:36 GMT+0100 (GMT+01:00)
Meanings assigned to audio clips:
attention1.wav: Hello how are you
attention2.wav: Hi how are you
h3nglaItskwu.wav: No meaning
h3nz0fjaIt.wav: No meaning
p3nh0nklait.wav: No meaning
SeIjubl.wav: No meaning
sm0tskaIspeIm.wav: No meaning
str3njaItstrEd.wav: No meaning
teIDublpr3njaItswEd.wav: No meaning
teIDubls3n2QnjaIt.wav: No meaning
teIDublstr3njaItstrEd.wav: No meaning
teIDubltr3nTraItstrEd.wav: No meaning
teIgutL.wav: No meaning
ZeIdubl.wav: No meaning
Participant ID: 5b4935fa479ca00001ef02d4
Timestamp: Wed Aug 01 2018 16:10:49 GMT+0100 (GMT+01:00)
Study withdrawn: false
Study submitted: true

Participant ID: 59b27b8b2a78fd00010b8403
Timestamp: Wed Aug 01 2018 16:05:51 GMT+0100 (British Summer Time)
Native language: English
Participant ID: 59b27b8b2a78fd00010b8403
Timestamp: Wed Aug 01 2018 16:06:29 GMT+0100 (British Summer Time)
Meanings assigned to audio clips:
attention1.wav: HELLO HOW ARE YOU
attention2.wav: HI HOW ARE YOU
h3nglaItskwu.wav: NO MEANING
h3nz0fjaIt.wav: NO MEANING
p3nh0nklait.wav: NO MEANING
SeIjubl.wav: NO MEANING
sm0tskaIspeIm.wav: NO MEANING
str3njaItstrEd.wav: NO MEANING
teIDublpr3njaItswEd.wav: NO MEANING
teIDubls3n2QnjaIt.wav: NO MEANING
teIDublstr3njaItstrEd.wav: NO MEANING
teIDubltr3nTraItstrEd.wav: NO MEANING
teIgutL.wav: NO MEANING
ZeIdubl.wav: NO MEANING
Participant ID: 59b27b8b2a78fd00010b8403
Timestamp: Wed Aug 01 2018 16:08:45 GMT+0100 (British Summer Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5aebb02c9457e00001052cbd
Timestamp: Wed Aug 01 2018 10:04:56 GMT-0500 (Central Daylight Time)
Native language: English
Participant ID: 5aebb02c9457e00001052cbd
Timestamp: Wed Aug 01 2018 10:05:15 GMT-0500 (Central Daylight Time)
Meanings assigned to audio clips:

attention1.wav: Hello, how are you?
attention2.wav: Hi, how are you?
h3nglaItskwu.wav: No meaning
h3nz0fjaIt.wav: Hands off the yacht
p3nh0nklaIt.wav: No meaning
SeIjublL.wav: No meaning
sm0tskaIspeIm.wav: No meaning
str3njaItstrEd.wav: No meaning
teIDublpr3njaItswEd.wav: No meaning
teIDubls3nZQnjaIt.wav: No meaning
teIDublstr3njaItstrEd.wav: No meaning
teIDubltr3nTraItstrEd.wav: No meaning
teIgutL.wav: No meaning
ZeIdublL.wav: No meaning
Participant ID: 5aebb02c9457e00001052cbd
Timestamp: Wed Aug 01 2018 10:08:10 GMT-0500 (Central Daylight Time)
Study withdrawn: false
Study submitted: true

Participant ID: 553eaedcfd99b278df6d2f4
Timestamp: Wed Aug 01 2018 11:05:14 GMT-0400 (Eastern Daylight Time)
Native language: English
Participant ID: 553eaedcfd99b278df6d2f4
Timestamp: Wed Aug 01 2018 11:05:26 GMT-0400 (Eastern Daylight Time)
Meanings assigned to audio clips:
attention1.wav: Hello how are you
attention2.wav: Hi how are you
h3nglaItskwu.wav: Turn the lights on
h3nz0fjaIt.wav: NO MEANING
p3nh0nklaIt.wav: NO MEANING
SeIjublL.wav: Sure you know
sm0tskaIspeIm.wav: NO MEANING
str3njaItstrEd.wav: NO MEANING
teIDublpr3njaItswEd.wav: NO MEANING
teIDubls3nZQnjaIt.wav: NO MEANING
teIDublstr3njaItstrEd.wav: NO MEANING
teIDubltr3nTraItstrEd.wav: NO MEANING
teIgutL.wav: NO MEANING
ZeIdublL.wav: They do not
Participant ID: 553eaedcfd99b278df6d2f4
Timestamp: Wed Aug 01 2018 11:07:53 GMT-0400 (Eastern Daylight Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5b4485e1d9655000014205a4
Timestamp: Wed Aug 01 2018 16:01:10 GMT+0100 (British Summer Time)
Native language: english
Participant ID: 5b4485e1d9655000014205a4
Timestamp: Wed Aug 01 2018 16:01:16 GMT+0100 (British Summer Time)
Meanings assigned to audio clips:
attention1.wav: Hello how are you
attention2.wav: HI how are you
h3nglaItskwu.wav: NO MEANING
h3nz0fjaIt.wav: hands off the ite
p3nh0nklaIt.wav: NO MEANING
SeIjublL.wav: NO MEANING
sm0tskaIspeIm.wav: sky spew
str3njaItstrEd.wav: NO MEANING
teIDublpr3njaItswEd.wav: NO MEANING
teIDubls3nZQnjaIt.wav: NO MEANING
teIDublstr3njaItstrEd.wav: NO MEANING
teIDubltr3nTraItstrEd.wav: NO MEANING
teIgutL.wav: NO MEANING
ZeIdublL.wav: NO MEANING
Participant ID: 5b4485e1d9655000014205a4
Timestamp: Wed Aug 01 2018 16:05:08 GMT+0100 (British Summer Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5b4903e1650a2200011ed0a2
Timestamp: Wed Aug 01 2018 16:01:16 GMT+0100 (GMT+01:00)
Native language: English
Participant ID: 5b4903e1650a2200011ed0a2
Timestamp: Wed Aug 01 2018 16:01:29 GMT+0100 (GMT+01:00)
Meanings assigned to audio clips:
attention1.wav: Hello how are you
attention2.wav: Hi how are you
h3nglaItskwu.wav: No meaning
h3nz0fjaIt.wav: Hands off
p3nh0nklaIt.wav: Come home
SeIjublL.wav: No meaning
sm0tskaIspeIm.wav: No meaning
str3njaItstrEd.wav: No meaning
teIDublpr3njaItswEd.wav: No meaning
teIDubls3nZQnjaIt.wav: No meaning
teIDublstr3njaItstrEd.wav: No meaning
teIDubltr3nTraItstrEd.wav: No meaning
teIgutL.wav: No meaning
ZeIdublL.wav: No meaning
Participant ID: 5b4903e1650a2200011ed0a2
Timestamp: Wed Aug 01 2018 16:04:47 GMT+0100 (GMT+01:00)
Study withdrawn: false

Study submitted: true

Participant ID: 5b5b07019351420001b73842
 Timestamp: Wed Aug 01 2018 15:57:33 GMT+0100 (BST)
 Native language: english
 Participant ID: 5b5b07019351420001b73842
 Timestamp: Wed Aug 01 2018 15:57:40 GMT+0100 (BST)
 Meanings assigned to audio clips:
 attention1.wav: hello, how are you?
 attention2.wav: hi, how are you?
 h3nglaItskwu.wav: NO MEANING
 h3nz0fjaIt.wav: NO MEANING
 p3nh0nklaIt.wav: NO MEANING
 SeIjubL.wav: NO MEANING
 sm0tskaIspeIm.wav: NO MEANING
 str3njaItstrEd.wav: NO MEANING
 teIDublpr3njaItswEd.wav: NO MEANING
 teIDubls3nZQnjaIt.wav: NO MEANING
 teIDublstr3njaItstrEd.wav: NO MEANING
 teIDubltr3nTraItstrEd.wav: NO MEANING
 teIgutL.wav: NO MEANING
 ZeIdubl.wav: NO MEANING

Participant ID: 5b5b07019351420001b73842
 Timestamp: Wed Aug 01 2018 15:59:32 GMT+0100 (BST)
 Study withdrawn: false
 Study submitted: true

Participant ID: 5b3f618f3d066d00016078f2
 Timestamp: Wed Aug 01 2018 16:02:02 GMT+0100 (British Summer Time)
 Native language: English
 Participant ID: 5b3f618f3d066d00016078f2
 Timestamp: Wed Aug 01 2018 16:02:09 GMT+0100 (British Summer Time)
 Meanings assigned to audio clips:
 attention1.wav: Hello how are you
 attention2.wav: Hi how are you
 h3nglaItskwu.wav: NO MEANING
 h3nz0fjaIt.wav: Hands off the item
 p3nh0nklaIt.wav: NO MEANING
 SeIjubL.wav: Je du blanc
 sm0tskaIspeIm.wav: NO MEANING
 str3njaItstrEd.wav: NO MEANING
 teIDublpr3njaItswEd.wav: NO MEANING
 teIDubls3nZQnjaIt.wav: NO MEANING
 teIDublstr3njaItstrEd.wav: NO MEANING
 teIDubltr3nTraItstrEd.wav: NO MEANING
 teIgutL.wav: NO MEANING
 ZeIdubl.wav: Je du blanc

Participant ID: 5b3f618f3d066d00016078f2
 Timestamp: Wed Aug 01 2018 16:04:20 GMT+0100 (British Summer Time)
 Study withdrawn: false
 Study submitted: true

Participant ID: 5b44bd7db080520001e0e566
 Timestamp: Wed Aug 01 2018 15:58:45 GMT+0100 (British Summer Time)
 Native language: English
 Participant ID: 5b44bd7db080520001e0e566
 Timestamp: Wed Aug 01 2018 15:59:00 GMT+0100 (British Summer Time)
 Meanings assigned to audio clips:
 attention1.wav: hello how are you
 attention2.wav: hi how are you
 h3nglaItskwu.wav: hurn glights grew
 h3nz0fjaIt.wav: hands off the item
 p3nh0nklaIt.wav: pern pon clight
 SeIjubL.wav: shayoodle
 sm0tskaIspeIm.wav: smots sky spain
 str3njaItstrEd.wav: screw
 teIDublpr3njaItswEd.wav: table purimites wed
 teIDubls3nZQnjaIt.wav: table sir jaun yite
 teIDublstr3njaItstrEd.wav: table stern rights dread
 teIDubltr3nTraItstrEd.wav: table turn rights dread
 teIgutL.wav: tagoodle
 ZeIdubl.wav: I doodle

Participant ID: 5b44bd7db080520001e0e566
 Timestamp: Wed Aug 01 2018 16:04:19 GMT+0100 (British Summer Time)
 Study withdrawn: false
 Study submitted: true

Participant ID: 5b50949c0266020001788ba1
 Timestamp: Wed Aug 01 2018 10:58:05 GMT-0400 (Eastern Daylight Time)
 Native language: english
 Participant ID: 5b50949c0266020001788ba1
 Timestamp: Wed Aug 01 2018 10:58:16 GMT-0400 (Eastern Daylight Time)
 Meanings assigned to audio clips:
 attention1.wav: HELLO HOW ARE YOU
 attention2.wav: HI HOW ARE YOU
 h3nglaItskwu.wav: NO MEANING
 h3nz0fjaIt.wav: NO MEANING
 p3nh0nklaIt.wav: NO MEANING
 SeIjubL.wav: NO MEANING
 sm0tskaIspeIm.wav: NO MEANING
 str3njaItstrEd.wav: NO MEANING
 teIDublpr3njaItswEd.wav: NO MEANING

teIDubls3nZQnjaIt.wav: NO MEANING
teIDublstr3njaItstrEd.wav: NO MEANING
teIDubltr3nTraItstrEd.wav: NO MEANING
teIgutL.wav: NO MEANING
ZeIdubl.wav: NO MEANING
Participant ID: 5b50949c0266020001788ba1
Timestamp: Wed Aug 01 2018 11:02:45 GMT-0400 (Eastern Daylight Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5b4aa57d94020800016315fb
Timestamp: Wed Aug 01 2018 10:59:22 GMT-0400 (Eastern Daylight Time)
Native language: English
Participant ID: 5b4aa57d94020800016315fb
Timestamp: Wed Aug 01 2018 10:59:34 GMT-0400 (Eastern Daylight Time)
Meanings assigned to audio clips:
attention1.wav: HELLO HOW ARE YOU
attention2.wav: HI HOW ARE YOU
h3nglaItskwu.wav: NO MEANING
h3nz0fjaIt.wav: HANDS OFF THE LIGHT
p3nh0nklaIt.wav: NO MEANING
SeIjubl.wav: NO MEANING
sm0tskaIspeIm.wav: SMOKE SKY SPAM
str3njaItstrEd.wav: NO MEANING
teIDublpr3njaItswEd.wav: NO MEANING
teIDubls3nZQnjaIt.wav: NO MEANING
teIDublstr3njaItstrEd.wav: NO MEANING
teIDubltr3nTraItstrEd.wav: NO MEANING
teIgutL.wav: NO MEANING
ZeIdubl.wav: NO MEANING
Participant ID: 5b4aa57d94020800016315fb
Timestamp: Wed Aug 01 2018 11:01:54 GMT-0400 (Eastern Daylight Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5b2bc4d30ec82d0001d299e1
Timestamp: Wed Aug 01 2018 15:55:54 GMT+0100 (British Summer Time)
Native language: English
Participant ID: 5b2bc4d30ec82d0001d299e1
Timestamp: Wed Aug 01 2018 15:56:04 GMT+0100 (British Summer Time)
Meanings assigned to audio clips:
attention1.wav: Hello how are you
attention2.wav: Hi how are you
h3nglaItskwu.wav: NO MEANING
h3nz0fjaIt.wav: Hands off the item
p3nh0nklaIt.wav: NO MEANING
SeIjubl.wav: NO MEANING
sm0tskaIspeIm.wav: NO MEANING
str3njaItstrEd.wav: NO MEANING
teIDublpr3njaItswEd.wav: NO MEANING
teIDubls3nZQnjaIt.wav: NO MEANING
teIDublstr3njaItstrEd.wav: NO MEANING
teIDubltr3nTraItstrEd.wav: NO MEANING
teIgutL.wav: NO MEANING
ZeIdubl.wav: NO MEANING
Participant ID: 5b2bc4d30ec82d0001d299e1
Timestamp: Wed Aug 01 2018 15:59:11 GMT+0100 (British Summer Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5b5c2fdbf7591300019061a1
Timestamp: Wed Aug 01 2018 15:58:12 GMT+0100 (British Summer Time)
Native language: English
Participant ID: 5b5c2fdbf7591300019061a1
Timestamp: Wed Aug 01 2018 15:58:21 GMT+0100 (British Summer Time)
Meanings assigned to audio clips:
attention1.wav: Hello how are you
attention2.wav: Hi how are you
h3nglaItskwu.wav: NO MEANING
h3nz0fjaIt.wav: NO MEANING
p3nh0nklaIt.wav: NO MEANING
SeIjubl.wav: NO MEANING
sm0tskaIspeIm.wav: NO MEANING
str3njaItstrEd.wav: NO MEANING
teIDublpr3njaItswEd.wav: NO MEANING
teIDubls3nZQnjaIt.wav: NO MEANING
teIDublstr3njaItstrEd.wav: NO MEANING
teIDubltr3nTraItstrEd.wav: NO MEANING
teIgutL.wav: NO MEANING
ZeIdubl.wav: NO MEANING
Participant ID: 5b5c2fdbf7591300019061a1
Timestamp: Wed Aug 01 2018 16:00:33 GMT+0100 (British Summer Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5b5db1db12af140001dcbb42
Timestamp: Wed Aug 01 2018 15:56:23 GMT+0100 (BST)
Native language: English
Participant ID: 5b5db1db12af140001dcbb42
Timestamp: Wed Aug 01 2018 15:56:34 GMT+0100 (BST)
Meanings assigned to audio clips:
attention1.wav: Hello how are you?

attention2.wav: Hi how are you
h3nglaItskwu.wav: Handglides gloom
h3nz0fjaIt.wav: Hands of light
p3nh0nklaIt.wav: Turn on light
SeIjubL.wav: Jâ[untranscribed symbol]ai oublon
sm0tskaIspeIm.wav: Smokes cause pain
str3njaItstrEd.wav: Strong nights dread
teIDublpr3njaItswEd.wav: Take coupon. Prune nights wed
teIDubls3nZQnjaIt.wav: Take coupon. Search all night
teIDublstr3njaItstrEd.wav: Takegroupon. Strew nights dread
teIDubltr3nTraItstrEd.wav: Take coupon. True brides dread
teIgutL.wav: Take couple
ZeIdubL.wav: Jâ[untranscribed symbol]ai double
Participant ID: 5b5db1db12af140001dcbb42
Timestamp: Wed Aug 01 2018 16:00:54 GMT+0100 (BST)
Study withdrawn: false
Study submitted: true

A.2.4 Retest Results

Successful to Unsuccessful Adversarial Wake-phrases

```
Result for SeIju`%b@L.wav:
WARNING:root:Transcript of user request: "show".
WARNING:root:Transcript of user request: "raise".
WARNING:root:Transcript of user request: "so you".
WARNING:root:Transcript of user request: "volume".
WARNING:root:Transcript of user request: "so you know".
WARNING:root:Transcript of user request: "so you know".
WARNING:root:Transcript of user request: "so you know".
WARNING:root:Playing assistant response.
```

```
Result for teIDu`%b@L.wav:
WARNING:root:Transcript of user request: "turn".
WARNING:root:Transcript of user request: "turn the".
WARNING:root:Transcript of user request: "turn the".
WARNING:root:Transcript of user request: "turn the".
WARNING:root:Transcript of user request: "turn Google".
WARNING:root:Expecting follow-on query from user.
```

```
Result for TreIbu`%s@L.wav:
WARNING:root:Transcript of user request: "read".
WARNING:root:Transcript of user request: "wake".
WARNING:root:Transcript of user request: "raise".
WARNING:root:Transcript of user request: "read me".
WARNING:root:Transcript of user request: "remove".
WARNING:root:Transcript of user request: "Radio 4".
WARNING:root:Transcript of user request: "Radio 4".
WARNING:root:Playing assistant response.
```

```
Result for tweIgu`%f@L.wav:
WARNING:root:Transcript of user request: "we".
WARNING:root:Transcript of user request: "wake".
WARNING:root:Transcript of user request: "where do".
WARNING:root:Transcript of user request: "wake you".
WARNING:root:Transcript of user request: "wake you".
WARNING:root:Transcript of user request: "radio 4".
WARNING:root:Playing assistant response.
```

```
Result for SeIhUk@L.wav:
WARNING:root:Transcript of user request: "show".
WARNING:root:Transcript of user request: "sing".
WARNING:root:Transcript of user request: "sing a".
WARNING:root:Transcript of user request: "single".
WARNING:root:Transcript of user request: "Bengal tour".
WARNING:root:Transcript of user request: "hey Google".
WARNING:root:Transcript of user request: "10".
WARNING:root:Playing assistant response.
```

```
Result for tSeIDu`%t@L.wav:
WARNING:root:Transcript of user request: "show".
WARNING:root:Transcript of user request: "shout".
WARNING:root:Transcript of user request: "JD".
WARNING:root:Transcript of user request: "schedule".
WARNING:root:Transcript of user request: "show duplo".
WARNING:root:Transcript of user request: "show duplo".
WARNING:root:Transcript of user request: "show duplo".
```

```
Result for tSeIDu`%s@L.wav:
WARNING:root:Transcript of user request: "show".
WARNING:root:Transcript of user request: "shout".
WARNING:root:Transcript of user request: "JD".
WARNING:root:Transcript of user request: "shaders".
WARNING:root:Transcript of user request: "shaders".
WARNING:root:Playing assistant response.
```

```
Result for ZeIju`%b@L.wav:
WARNING:root:Transcript of user request: "are you".
WARNING:root:Transcript of user request: "volume".
WARNING:root:Transcript of user request: "volume".
WARNING:root:Transcript of user request: "volume".
WARNING:root:Playing assistant response.
```

Successful to Unsuccessful Adversarial Commands

```
Result for aCczOm4gz2wnYU10IHNrddy10j09PQ==.wav:
WARNING:root:Transcript of user request: "turn".
WARNING:root:Transcript of user request: "turn light".
WARNING:root:Transcript of user request: "turn lights".
WARNING:root:Transcript of user request: "turn lights".
WARNING:root:Transcript of user request: "turn light School".
WARNING:root:Transcript of user request: "turn light School".
WARNING:root:Transcript of user request: "turn light School".
WARNING:root:Playing assistant response.
```

```
Result for c20nMHRzIGsnYUkge3AnZUltPT09.wav:
WARNING:root:Transcript of user request: "smart".
```

```

WARNING:root:Transcript of user request: "smart Sky".
WARNING:root:Transcript of user request: "Smart cars".
WARNING:root:Transcript of user request: "Smart cars".
WARNING:root:Transcript of user request: "Smart cars pay".
WARNING:root:Transcript of user request: "smart khaya Spain".
WARNING:root:Transcript of user request: "smart khaya Spain".
WARNING:root:Transcript of user request: "smart khaya Spain".
WARNING:root:Playing assistant response.

Result for c20nMHRzIHNUJ2FJIHNSJ2VJbT09PQ==.wav:
WARNING:root:Transcript of user request: "smart".
WARNING:root:Transcript of user request: "what's my".
WARNING:root:Transcript of user request: "smarts nights".
WARNING:root:Transcript of user request: "smarts nicely".
WARNING:root:Transcript of user request: "smarts Knights Lane".
WARNING:root:Transcript of user request: "smarts Knights Lane".
WARNING:root:Transcript of user request: "smarts Knights Lane".
WARNING:root:Transcript of user request: "smarts Knights Lane".
WARNING:root:Playing assistant response.

Result for cCczOm4gdidhSXQgc2sndTo9PT0=.wav:
WARNING:root:Transcript of user request: "play".
WARNING:root:Transcript of user request: "turn".
WARNING:root:Transcript of user request: "turn my".
WARNING:root:Transcript of user request: "turn light".
WARNING:root:Transcript of user request: "turn lights".
WARNING:root:Transcript of user request: "turn lights".
WARNING:root:Transcript of user request: "turn lights ku".
WARNING:root:Transcript of user request: "turn light School".
WARNING:root:Transcript of user request: "turn light School".
WARNING:root:Transcript of user request: "turn light School".
WARNING:root:Playing assistant response.

Result for cGwnMzpuIHNwbCdhSXQgc3BsJ3U6PT09.wav:
WARNING:root:Transcript of user request: "play".
WARNING:root:Transcript of user request: "turn".
WARNING:root:Transcript of user request: "let's play".
WARNING:root:Transcript of user request: "Lounge light".
WARNING:root:Transcript of user request: "lounge lights".
WARNING:root:Transcript of user request: "lounge lights to".
WARNING:root:Transcript of user request: "lounge lights blue".
WARNING:root:Transcript of user request: "lounge lights blue".
WARNING:root:Transcript of user request: "lounge lights blue".
WARNING:root:Transcript of user request: "lounge lights blue".
WARNING:root:Transcript of user request: "lounge lights blue".
WARNING:root:Playing assistant response.

Result for dHInMzpuIFryJ2FJdCBzdHInRWQ9PT0=.wav:
WARNING:root:Transcript of user request: "true".
WARNING:root:Transcript of user request: "true and".
WARNING:root:Transcript of user request: "true and right".
WARNING:root:Transcript of user request: "true and right".
WARNING:root:Transcript of user request: "true and right stray".
WARNING:root:Transcript of user request: "true and rights to Red".
WARNING:root:Transcript of user request: "true and rights to Red".
WARNING:root:Transcript of user request: "true and bright red".
MKB: SUPPLEMENTAL_DISPLAY_TEXT: It's a sure thing
WARNING:root:Playing assistant response.

Result for dHInMzpuIG4nMGYgZ2wnYU10PT09.wav:
WARNING:root:Transcript of user request: "true".
WARNING:root:Transcript of user request: "true and".
WARNING:root:Transcript of user request: "turn off".
WARNING:root:Transcript of user request: "turn off the".
WARNING:root:Transcript of user request: "turn off light".
WARNING:root:Transcript of user request: "turn off light".
WARNING:root:Transcript of user request: "true enough blood".
MKB: SUPPLEMENTAL_DISPLAY_TEXT: Sorry, I don't understand
WARNING:root:Playing assistant response.

Result for ZmwnMzpuIG0nMG4ga2wnYU10PT09.wav:
WARNING:root:Transcript of user request: "no".
WARNING:root:Transcript of user request: "turn".
WARNING:root:Transcript of user request: "turn my".
WARNING:root:Transcript of user request: "alarm on".
WARNING:root:Transcript of user request: "turn mum's".
WARNING:root:Transcript of user request: "turn mum's light".
WARNING:root:Transcript of user request: "turn mum's light".
WARNING:root:Transcript of user request: "turn mum's light".
WARNING:root:Playing assistant response.

```

Unsuccessful to Successful Adversarial Wake-phrases

NONE

Unsuccessful to Successful Adversarial Commands

Result for Vcd10iBibCdhbSBzbsdhST09PQ==.wav:

WARNING:root:Transcript of user request: "lamp".
 WARNING:root:Transcript of user request: "the lamps my".
 WARNING:root:Transcript of user request: "the lamps my".
 WARNING:root:Transcript of user request: "who am I".
 WARNING:root:Playing assistant response.

Result for Vcd10iBrdydhhSBza3InYUK9PT0=.wav:
 WARNING:root:Transcript of user request: "who am".
 WARNING:root:Transcript of user request: "who am I".
 WARNING:root:Transcript of user request: "no qualms cry".
 WARNING:root:Transcript of user request: "no qualms cry".
 WARNING:root:Transcript of user request: "who am I".
 WARNING:root:Playing assistant response.

Result for Vcd10iBrdydhhSBzbSdhST09PQ=.wav:
 WARNING:root:Transcript of user request: "what's".
 WARNING:root:Transcript of user request: "who am I".
 WARNING:root:Transcript of user request: "what's my".
 WARNING:root:Transcript of user request: "what's my".
 WARNING:root:Transcript of user request: "who am i".
 WARNING:root:Playing assistant response.

Result for Vcd10iBrdydhhSBnbCdhST09PQ=.wav:
 WARNING:root:Transcript of user request: "who am".
 WARNING:root:Transcript of user request: "who am I".
 WARNING:root:Transcript of user request: "Wembley".
 WARNING:root:Transcript of user request: "Wembley".
 WARNING:root:Transcript of user request: "who am I".
 WARNING:root:Playing assistant response.

Result for Zid10iBiJ2FtIGJsJ2FJPT09.wav:
 WARNING:root:Transcript of user request: "Who".
 WARNING:root:Transcript of user request: "super".
 WARNING:root:Transcript of user request: "who am".
 WARNING:root:Transcript of user request: "who am I".
 WARNING:root:Transcript of user request: "who am I".
 WARNING:root:Transcript of user request: "who am I".
 WARNING:root:Playing assistant response.

Result for Rcd10iBiJ2FtIGsnYUK9PT0=.wav:
 WARNING:root:Transcript of user request: "resume".
 WARNING:root:Transcript of user request: "Dubai".
 WARNING:root:Transcript of user request: "do plan".
 WARNING:root:Transcript of user request: "no thank".
 WARNING:root:Transcript of user request: "new panic I".
 WARNING:root:Transcript of user request: "new Vampire".
 WARNING:root:Transcript of user request: "who am I".
 WARNING:root:Playing assistant response.

Result for did10iBiJ2FtIFonYUK9PT0=.wav:
 WARNING:root:Transcript of user request: "who am".
 WARNING:root:Transcript of user request: "thumbs".
 WARNING:root:Transcript of user request: "you by".
 WARNING:root:Transcript of user request: "who am I".
 WARNING:root:Playing assistant response.

Result for Rcd10iB0UydhbSBncidhST09PQ=.wav:
 WARNING:root:Transcript of user request: "new channel".
 WARNING:root:Transcript of user request: "new channel Grey".
 WARNING:root:Transcript of user request: "new channel Grey".
 WARNING:root:Transcript of user request: "who am i".
 WARNING:root:Playing assistant response.

Result for Vcd10iBrdydhhSB6J2FJPT09.wav:
 WARNING:root:Transcript of user request: "thumbs".
 WARNING:root:Transcript of user request: "who am I".
 WARNING:root:Transcript of user request: "who am I".
 WARNING:root:Transcript of user request: "who am I".
 WARNING:root:Playing assistant response.

Result for Zid10iBiJ2FtIGtsJ2FJPT09.wav:
 WARNING:root:Transcript of user request: "Who".
 WARNING:root:Transcript of user request: "super".
 WARNING:root:Transcript of user request: "who am".
 WARNING:root:Transcript of user request: "who lamp".
 WARNING:root:Transcript of user request: "Ubuntu".
 WARNING:root:Transcript of user request: "Superman fly".
 WARNING:root:Transcript of user request: "who am I".
 WARNING:root:Playing assistant response.

Result for dHcndTogdidhhSBncidhST09PQ=.wav:
 WARNING:root:Transcript of user request: "movie".
 WARNING:root:Transcript of user request: "move on my".
 WARNING:root:Transcript of user request: "who am I".
 WARNING:root:Playing assistant response.
 WARNING:root:Finished playing assistant response.

Result for did10iB6J2FtIGdyJ2FJPT09.wav:
 WARNING:root:Transcript of user request: "who's".
 WARNING:root:Transcript of user request: "resume".
 WARNING:root:Transcript of user request: "who's Mr Grey".
 WARNING:root:Transcript of user request: "who am i".

WARNING:root:Playing assistant response.
 WARNING:root:Finished playing assistant response.

Result for Vcd10iB0dydhsSB0dydhST09PQ==.wav:
 WARNING:root:Transcript of user request: "who won".
 WARNING:root:Transcript of user request: "who am I".
 WARNING:root:Transcript of user request: "who am I".
 WARNING:root:Transcript of user request: "who am I".
 WARNING:root:Transcript of user request: "who am I".
 WARNING:root:Playing assistant response.
 WARNING:root:Finished playing assistant response.

Result for did10iBicidhsSBrJ2FJPT09.wav:
 WARNING:root:Transcript of user request: "volume".
 WARNING:root:Transcript of user request: "volume by".
 WARNING:root:Transcript of user request: "vampire".
 WARNING:root:Transcript of user request: "who am I".
 WARNING:root:Playing assistant response.
 WARNING:root:Finished playing assistant response.

Result for Rcd10iBiJ2FtIHR3J2FJPT09.wav:
 WARNING:root:Transcript of user request: "resume".
 WARNING:root:Transcript of user request: "super".
 WARNING:root:Transcript of user request: "rhubarb".
 WARNING:root:Transcript of user request: "Dubai".
 WARNING:root:Transcript of user request: "Super Why".
 WARNING:root:Transcript of user request: "who am I".
 WARNING:root:Playing assistant response.
 WARNING:root:Finished playing assistant response.

Result for did10iBicidhsSbnCdhST09PQ==.wav:
 WARNING:root:Transcript of user request: "will I".
 WARNING:root:Transcript of user request: "the big light".
 WARNING:root:Transcript of user request: "who am I".
 WARNING:root:Playing assistant response.
 WARNING:root:Finished playing assistant response.

WARNING:root:Transcript of user request: "volume".
 WARNING:root:Transcript of user request: "thumbs".
 WARNING:root:Transcript of user request: "thumbs".
 WARNING:root:Transcript of user request: "who am I".
 WARNING:root:Playing assistant response.
 WARNING:root:Finished playing assistant response.

Result for Vcd10iBza3cnYW0gaidhST09PQ==.wav:
 WARNING:root:Transcript of user request: "who's".
 WARNING:root:Transcript of user request: "New Square".
 WARNING:root:Transcript of user request: "who am".
 WARNING:root:Transcript of user request: "who am I".
 WARNING:root:Transcript of user request: "who am I".
 WARNING:root:Transcript of user request: "who am I".
 WARNING:root:Transcript of user request: "who am I".
 WARNING:root:Playing assistant response.

Result for Vcd10iBzdHInYW0gc20nYuk9PT0=.wav:
 WARNING:root:Transcript of user request: "who's".
 WARNING:root:Transcript of user request: "feels true".
 WARNING:root:Transcript of user request: "who's my".
 WARNING:root:Transcript of user request: "who's my".
 WARNING:root:Transcript of user request: "who am I".
 WARNING:root:Playing assistant response.

Result for Vcd10iBiJ2FtIHR3J2FJPT09.wav:
 WARNING:root:Transcript of user request: "why".
 WARNING:root:Transcript of user request: "why".
 WARNING:root:Transcript of user request: "who am I".
 WARNING:root:Playing assistant response.

Result for Zid10iBEJ2FtIGdyJ2FJPT09.wav:
 WARNING:root:Transcript of user request: "Who".
 WARNING:root:Transcript of user request: "who's".
 WARNING:root:Transcript of user request: "hoover".
 WARNING:root:Transcript of user request: "who's on".
 WARNING:root:Transcript of user request: "who than we".
 WARNING:root:Transcript of user request: "who them dry".
 WARNING:root:Transcript of user request: "who am I".
 WARNING:root:Playing assistant response.

Result for Vcd10iB0dydhsSBkWhdST09PQ==.wav:
 WARNING:root:Transcript of user request: "who won".
 WARNING:root:Transcript of user request: "am J".
 WARNING:root:Transcript of user request: "am J".
 WARNING:root:Transcript of user request: "who am i".
 WARNING:root:Playing assistant response.

Result for Vcd10iBUJ2FtIGsnYuk9PT0=.wav:
 WARNING:root:Transcript of user request: "boost".
 WARNING:root:Transcript of user request: "sofa".
 WARNING:root:Transcript of user request: "who sang".
 WARNING:root:Transcript of user request: "no thank".
 WARNING:root:Transcript of user request: "no thank".

WARNING:root:Transcript of user request: "who am i".
WARNING:root:Playing assistant response.

Result for did10iBiJ2FtIHN0cidhST09PQ==.wav:
WARNING:root:Transcript of user request: "you by".
WARNING:root:Transcript of user request: "who am I".
WARNING:root:Playing assistant response.

Result for ZnInMHRzIGdsJ2FJIHRTJ2VJbT09PQ==.wav:
WARNING:root:Transcript of user request: "what".
WARNING:root:Transcript of user request: "what's".
WARNING:root:Transcript of user request: "what's".
WARNING:root:Transcript of user request: "run slideshow".
WARNING:root:Transcript of user request: "what's my name".
WARNING:root:Transcript of user request: "what's my name".
WARNING:root:Playing assistant response.

Result for aCcwDhMgaidhSSBUJ2VJbT09PQ==.wav:
WARNING:root:Transcript of user request: "how".
WARNING:root:Transcript of user request: "what's".
WARNING:root:Transcript of user request: "what's the".
WARNING:root:Transcript of user request: "what's my".
WARNING:root:Transcript of user request: "what's your".
WARNING:root:Transcript of user request: "what's the answer".
WARNING:root:Transcript of user request: "what's the answer".
WARNING:root:Transcript of user request: "what's the iPhone".
WARNING:root:Transcript of user request: "what's the iPhone".
WARNING:root:Transcript of user request: "what's the iPhone".
WARNING:root:Transcript of user request: "what's my name".
WARNING:root:Playing assistant response.

Result for c3cnMzpuIGgnMG4gc3BsJ2FJdD09PQ==.wav:
WARNING:root:Transcript of user request: "turn".
WARNING:root:Transcript of user request: "turn on".
WARNING:root:Transcript of user request: "Owen Hunt".
WARNING:root:Transcript of user request: "turn on Sky".
WARNING:root:Transcript of user request: "turn on light".
WARNING:root:Transcript of user request: "turn on light".
WARNING:root:Transcript of user request: "turn on light".
WARNING:root:Transcript of user request: "turn on light".
WARNING:root:Playing assistant response.

Result for cyczOm4ga2wnMG4geidhSXQ9PT0=.wav:
WARNING:root:Transcript of user request: "so".
WARNING:root:Transcript of user request: "turn".
WARNING:root:Transcript of user request: "turn on".
WARNING:root:Transcript of user request: "turn plum".
WARNING:root:Transcript of user request: "turn on".
WARNING:root:Transcript of user request: "turn plum".
WARNING:root:Transcript of user request: "turn plum".
WARNING:root:Transcript of user request: "turn bronze".
WARNING:root:Transcript of user request: "turn Klondike".
WARNING:root:Transcript of user request: "turn on light".
WARNING:root:Playing assistant response.

Result for bSczOm4gVCCwbiBEJ2FJdD09PQ==.wav:
WARNING:root:Transcript of user request: "no".
WARNING:root:Transcript of user request: "man".
WARNING:root:Transcript of user request: "month".
WARNING:root:Transcript of user request: "Man by".
WARNING:root:Transcript of user request: "month on".
WARNING:root:Transcript of user request: "man from the".
WARNING:root:Transcript of user request: "month on light".
WARNING:root:Transcript of user request: "turn on light".
WARNING:root:Playing assistant response.

Result for cHInMzpuIGJsJzBuIFQnYU10PT09.wav:
WARNING:root:Transcript of user request: "remember".
WARNING:root:Transcript of user request: "turn blue".
WARNING:root:Transcript of user request: "turn on".
WARNING:root:Transcript of user request: "turn blonde".
WARNING:root:Transcript of user request: "turn on side".
WARNING:root:Transcript of user request: "prevent lung fight".
WARNING:root:Transcript of user request: "turn on light".
WARNING:root:Playing assistant response.

Result for ZmwnMzpuIHcnMGYgZydhSXQ9PT0=.wav:
WARNING:root:Transcript of user request: "no".
WARNING:root:Transcript of user request: "turn".
WARNING:root:Transcript of user request: "turn on".
WARNING:root:Transcript of user request: "turn off".
WARNING:root:Transcript of user request: "turn off my".
WARNING:root:Transcript of user request: "turn off guide".
WARNING:root:Transcript of user request: "turn off guide".
WARNING:root:Transcript of user request: "turn off light".
WARNING:root:Playing assistant response.

Result for UyczOm4geicwZiB2J2FJdD09PQ==.wav:
WARNING:root:Transcript of user request: "show".
WARNING:root:Transcript of user request: "turn".
WARNING:root:Transcript of user request: "shannons".
WARNING:root:Transcript of user request: "turn the".

WARNING:root:Transcript of user request: "sons of".
 WARNING:root:Transcript of user request: "turn off my".
 WARNING:root:Transcript of user request: "Johns of fight".
 WARNING:root:Transcript of user request: "turn off light".
 WARNING:root:Playing assistant response.

Result for cHInMzpuIGdsJzBmIGtyJ2FJdD09PQ==.wav:
 WARNING:root:Transcript of user request: "turn".
 WARNING:root:Transcript of user request: "turn off".
 WARNING:root:Transcript of user request: "turn off".
 WARNING:root:Transcript of user request: "turn off".
 WARNING:root:Transcript of user request: "turn off light".
 WARNING:root:Transcript of user request: "turn off light".
 WARNING:root:Playing assistant response.

Result for ZHInMzpuIGwnMGYgeidhSXQ9PT0=.wav:
 WARNING:root:Transcript of user request: "turn".
 WARNING:root:Transcript of user request: "grandma".
 WARNING:root:Transcript of user request: "turn off".
 WARNING:root:Transcript of user request: "turn off site".
 WARNING:root:Transcript of user request: "turn off site".
 WARNING:root:Transcript of user request: "turn off light".
 WARNING:root:Playing assistant response.

Result for WiczOm4gbicwZiB0dydhSXQ9PT0=.wav:
 WARNING:root:Transcript of user request: "turn off".
 WARNING:root:Transcript of user request: "why".
 WARNING:root:Transcript of user request: "turn off white".
 WARNING:root:Transcript of user request: "turn off white".
 WARNING:root:Transcript of user request: "turn off light".
 WARNING:root:Playing assistant response.

Result for YmwnMzpuIG4nMGYgVCdhSXQ9PT0=.wav:
 WARNING:root:Transcript of user request: "no".
 WARNING:root:Transcript of user request: "play".
 WARNING:root:Transcript of user request: "lung".
 WARNING:root:Transcript of user request: "London".
 WARNING:root:Transcript of user request: "turn off".
 WARNING:root:Transcript of user request: "turn off the".
 WARNING:root:Transcript of user request: "turn off side".
 WARNING:root:Transcript of user request: "turn off side".
 WARNING:root:Transcript of user request: "turn off light".
 WARNING:root:Playing assistant response.

Result for YnInMzpuIG0nMGYgdidhSXQ9PT0=.wav:
 WARNING:root:Transcript of user request: "rewind".
 WARNING:root:Transcript of user request: "remind me".
 WARNING:root:Transcript of user request: "turn off".
 WARNING:root:Transcript of user request: "turn off the".
 WARNING:root:Transcript of user request: "turn off my".
 WARNING:root:Transcript of user request: "turn off light".
 WARNING:root:Transcript of user request: "turn off light".
 WARNING:root:Transcript of user request: "turn off light".
 WARNING:root:Transcript of user request: "turn off light".
 WARNING:root:Playing assistant response.

Result for VCczOm4gc3BsJ2FJdCB2J3U6PT09.wav:
 WARNING:root:Transcript of user request: "turn".
 WARNING:root:Transcript of user request: "thumbs".
 WARNING:root:Transcript of user request: "thumbs flight".
 WARNING:root:Transcript of user request: "thumbs flight".
 WARNING:root:Transcript of user request: "turn light new".
 WARNING:root:Transcript of user request: "turn light new".
 WARNING:root:Transcript of user request: "turn light blue".
 WARNING:root:Playing assistant response.

Result for eiczOm4gZydhSXQgcGwndTo9PT0=.wav:
 WARNING:root:Transcript of user request: "the".
 WARNING:root:Transcript of user request: "turn".
 WARNING:root:Transcript of user request: "send".
 WARNING:root:Transcript of user request: "send".
 WARNING:root:Transcript of user request: "turn light blue".
 WARNING:root:Transcript of user request: "turn light blue".
 WARNING:root:Playing assistant response.

Result for Z2wnMzpuIGtsJ2FJdCBEJ3U6PT09.wav:
 WARNING:root:Transcript of user request: "play".
 WARNING:root:Transcript of user request: "turn".
 WARNING:root:Transcript of user request: "London".
 WARNING:root:Transcript of user request: "Blunt".
 WARNING:root:Transcript of user request: "Lenka AI".
 WARNING:root:Transcript of user request: "turn light".
 WARNING:root:Transcript of user request: "turn lights".
 WARNING:root:Transcript of user request: "Blunt like you".
 WARNING:root:Transcript of user request: "turn light blue".
 WARNING:root:Playing assistant response.

A.3 Main Experiment Results

A.3.1 Audio File Input Results

Results for "hey google":
successes: 70
failures: 417

Wins for hey google at Level 1: 52

Z'eI g'u:g@L
v'eI g'u:g@L
T'eI g'u:g@L
h'eI k'Ud@L
h'eI g'Ud@L
h'eI gr'Ud@L
h'eI kw'Ud@L
h'eI d'u:b@L
h'eI g'u:b@L
h'eI w'u:b@L
h'eI T'u:b@L
h'eI D'u:b@L
h'eI kl'u:b@L
h'eI gl'u:b@L
h'eI gr'u:b@L
h'eI sn'u:b@L
h'eI Tr'u:b@L
h'eI gl'u:f@L
h'eI p'Uk@L
h'eI k'Uk@L
h'eI d'Uk@L
h'eI g'Uk@L
h'eI h'Uk@L
h'eI v'Uk@L
h'eI z'Uk@L
h'eI l'Uk@L
h'eI n'Uk@L
h'eI w'Uk@L
h'eI T'Uk@L
h'eI D'Uk@L
h'eI bl'Uk@L
h'eI kl'Uk@L
h'eI kr'Uk@L
h'eI fl'Uk@L
h'eI fr'Uk@L
h'eI gl'Uk@L
h'eI gr'Uk@L
h'eI pl'Uk@L
h'eI pr'Uk@L
h'eI Tr'Uk@L
h'eI d'u:p@L
h'eI n'u:p@L
h'eI gl'u:p@L
h'eI k'u:s@L
h'eI g'u:s@L
h'eI bl'u:s@L
h'eI gl'u:s@L
h'eI g'u:t@L
h'eI gl'u:t@L
h'eI k'u:z@L
h'eI g'u:z@L
h'eI gl'u:z@L

Wins for hey google at Level 2: 18

T'eI g'Uk@L
v'eI g'u:b@L
v'eI g'Uk@L
Z'eI g'u:b@L
Z'eI gl'u:s@L
T'eI g'u:b@L
T'eI gl'u:s@L
Z'eI fl'Uk@L
Z'eI Tr'u:b@L
Z'eI gl'u:p@L
Z'eI kr'Uk@L
v'eI gl'u:p@L
v'eI g'u:t@L
Z'eI l'Uk@L
v'eI g'Ud@L
Z'eI D'u:b@L
Z'eI gl'Uk@L
Z'eI g'u:z@L

Results for "who am i":
successes: 85
failures: 199

Wins for who am i at Level 1: 46

f'u: 'am 'aI
v'u: 'am 'aI

T'u: 'am 'aI
D'u: 'am 'aI
pl'u: 'am 'aI
pr'u: 'am 'aI
kw'u: 'am 'aI
sp'u: 'am 'aI
tw'u: 'am 'aI
spl'u: 'am 'aI
skw'u: 'am 'aI
h'u: b'am 'aI
h'u: g'am 'aI
h'u: f'am 'aI
h'u: s'am 'aI
h'u: Z'am 'aI
h'u: v'am 'aI
h'u: z'am 'aI
h'u: m'am 'aI
h'u: tS'am 'aI
h'u: T'am 'aI
h'u: D'am 'aI
h'u: bl'am 'aI
h'u: br'am 'aI
h'u: fr'am 'aI
h'u: kw'am 'aI
h'u: sn'am 'aI
h'u: tw'am 'aI
h'u: Tr'am 'aI
h'u: 'am k'aI
h'u: 'am Z'aI
h'u: 'am z'aI
h'u: 'am j'aI
h'u: 'am bl'aI
h'u: 'am br'aI
h'u: 'am kl'aI
h'u: 'am gl'aI
h'u: 'am gr'aI
h'u: 'am sm'aI
h'u: 'am sn'aI
h'u: 'am sw'aI
h'u: 'am tw'aI
h'u: 'am spl'aI
h'u: 'am skr'aI
h'u: 'am str'aI
h'u: 'am Tr'aI

Wins for who am i at Level 2: 21

D'u: b'am 'aI
h'u: bl'am bl'aI
f'u: D'am 'aI
spl'u: b'am 'aI
D'u: bl'am 'aI
D'u: 'am j'aI
f'u: 'am Z'aI
spl'u: f'am 'aI
spl'u: kw'am 'aI
h'u: T'am j'aI
f'u: 'am bl'aI
pl'u: T'am 'aI
spl'u: T'am 'aI
v'u: T'am 'aI
h'u: T'am gl'aI
pl'u: 'am gl'aI
v'u: 'am sm'aI
h'u: f'am Z'aI
spl'u: bl'am 'aI
pl'u: 'am br'aI
h'u: z'am Z'aI

Wins for who am i at Level 3: 18

f'u: tS'am bl'aI
f'u: m'am Z'aI
f'u: tw'am Z'aI
spl'u: kw'am j'aI
f'u: Tr'am Z'aI
D'u: bl'am Z'aI
spl'u: bl'am str'aI
f'u: g'am Z'aI
v'u: f'am Z'aI
f'u: b'am Z'aI
f'u: bl'am bl'aI
f'u: D'am Z'aI
spl'u: z'am Z'aI
f'u: D'am z'aI
f'u: br'am Z'aI
spl'u: bl'am Z'aI
f'u: T'am Z'aI
D'u: b'am k'aI

Results for "whats my name":

successes: 165
failures: 133

Wins for whats my name at Level 1: 56

b'0ts m'ai n'eIm
g'0ts m'ai n'eIm
f'0ts m'ai n'eIm
Z'0ts m'ai n'eIm
h'0ts m'ai n'eIm
v'0ts m'ai n'eIm
z'0ts m'ai n'eIm
tS'0ts m'ai n'eIm
T'0ts m'ai n'eIm
D'0ts m'ai n'eIm
br'0ts m'ai n'eIm
kr'0ts m'ai n'eIm
dr'0ts m'ai n'eIm
fl'0ts m'ai n'eIm
gl'0ts m'ai n'eIm
gr'0ts m'ai n'eIm
pr'0ts m'ai n'eIm
tw'0ts m'ai n'eIm
spl'0ts m'ai n'eIm
w'0ts k'ai n'eIm
w'0ts dZ'ai n'eIm
w'0ts Z'ai n'eIm
w'0ts z'ai n'eIm
w'0ts j'ai n'eIm
w'0ts bl'ai n'eIm
w'0ts br'ai n'eIm
w'0ts kl'ai n'eIm
w'0ts gl'ai n'eIm
w'0ts gr'ai n'eIm
w'0ts kw'ai n'eIm
w'0ts sm'ai n'eIm
w'0ts sn'ai n'eIm
w'0ts sw'ai n'eIm
w'0ts tw'ai n'eIm
w'0ts spl'ai n'eIm
w'0ts Tr'ai n'eIm
w'0ts m'ai p'eIm
w'0ts m'ai b'eIm
w'0ts m'ai Z'eIm
w'0ts m'ai h'eIm
w'0ts m'ai v'eIm
w'0ts m'ai z'eIm
w'0ts m'ai j'eIm
w'0ts m'ai r'eIm
w'0ts m'ai w'eIm
w'0ts m'ai tS'eIm
w'0ts m'ai T'eIm
w'0ts m'ai D'eIm
w'0ts m'ai br'eIm
w'0ts m'ai gl'eIm
w'0ts m'ai gr'eIm
w'0ts m'ai sl'eIm
w'0ts m'ai sm'eIm
w'0ts m'ai sn'eIm
w'0ts m'ai sp'eIm
w'0ts m'ai tr'eIm

Wins for whats my name at Level 2: 52

f'0ts m'ai Z'eIm
w'0ts gr'ai Z'eIm
w'0ts sw'ai T'eIm
w'0ts Z'ai j'eIm
dr'0ts m'ai r'eIm
h'0ts m'ai h'eIm
g'0ts m'ai T'eIm
v'0ts sn'ai n'eIm
D'0ts m'ai v'eIm
pr'0ts m'ai r'eIm
b'0ts k'ai n'eIm
br'0ts m'ai b'eIm
z'0ts m'ai v'eIm
br'0ts m'ai T'eIm
w'0ts j'ai w'eIm
v'0ts gl'ai n'eIm
dr'0ts k'ai n'eIm
v'0ts m'ai Z'eIm
spl'0ts sm'ai n'eIm
D'0ts m'ai z'eIm
f'0ts m'ai r'eIm
g'0ts m'ai tS'eIm
tS'0ts m'ai z'eIm
fl'0ts m'ai Z'eIm
tw'0ts m'ai sl'eIm
br'0ts sm'ai n'eIm
D'0ts Z'ai n'eIm
h'0ts m'ai gl'eIm
w'0ts sw'ai v'eIm
tw'0ts sm'ai n'eIm
w'0ts k'ai h'eIm
tS'0ts m'ai Z'eIm

w'0ts kw'ai Z'eIm
w'0ts kw'ai br'eIm
T'0ts sm'ai n'eIm
w'0ts spl'ai j'eIm
w'0ts tw'ai D'eIm
b'0ts m'ai gl'eIm
tS'0ts gl'ai n'eIm
T'0ts m'ai z'eIm
w'0ts j'ai r'eIm
w'0ts j'ai v'eIm
D'0ts m'ai sn'eIm
w'0ts sm'ai gl'eIm
w'0ts k'ai j'eIm
pr'0ts bl'ai n'eIm
spl'0ts m'ai T'eIm
g'0ts m'ai Z'eIm
dr'0ts m'ai tS'eIm
w'0ts spl'ai w'eIm
w'0ts kw'ai sl'eIm
b'0ts z'ai n'eIm

Wins for whats my name at Level 3: 57

dr'0ts k'ai j'eIm
br'0ts sm'ai b'eIm
spl'0ts sm'ai z'eIm
h'0ts j'ai gl'eIm
br'0ts sm'ai D'eIm
fl'0ts j'ai w'eIm
tS'0ts k'ai z'eIm
v'0ts tw'ai D'eIm
D'0ts sm'ai v'eIm
f'0ts bl'ai Z'eIm
f'0ts k'ai h'eIm
h'0ts bl'ai gl'eIm
dr'0ts j'ai tS'eIm
h'0ts Z'ai gl'eIm
fl'0ts j'ai r'eIm
D'0ts Z'ai D'eIm
spl'0ts j'ai w'eIm
f'0ts Z'ai Z'eIm
D'0ts gr'ai z'eIm
h'0ts j'ai w'eIm
D'0ts j'ai v'eIm
spl'0ts sm'ai sl'eIm
T'0ts sm'ai r'eIm
f'0ts k'ai r'eIm
tS'0ts gl'ai v'eIm
b'0ts j'ai w'eIm
g'0ts kl'ai Z'eIm
h'0ts bl'ai h'eIm
tS'0ts k'ai j'eIm
b'0ts j'ai v'eIm
kr'0ts Z'ai j'eIm
tw'0ts j'ai r'eIm
tw'0ts sm'ai p'eIm
tS'0ts sm'ai Z'eIm
tS'0ts gl'ai D'eIm
fl'0ts j'ai v'eIm
h'0ts tw'ai gl'eIm
D'0ts Z'ai v'eIm
D'0ts spl'ai v'eIm
D'0ts sn'ai z'eIm
D'0ts sw'ai sn'eIm
T'0ts Z'ai j'eIm
D'0ts bl'ai v'eIm
g'0ts sm'ai tS'eIm
tw'0ts sm'ai sl'eIm
D'0ts sm'ai sn'eIm
h'0ts j'ai r'eIm
h'0ts sw'ai T'eIm
v'0ts j'ai w'eIm
f'0ts j'ai Z'eIm
g'0ts j'ai v'eIm
D'0ts br'ai v'eIm
dr'0ts bl'ai r'eIm
g'0ts j'ai Z'eIm
pr'0ts Z'ai j'eIm
D'0ts br'ai sn'eIm
b'0ts sm'ai gl'eIm

Results for "turn on light":

successes: 155
failures: 160

Wins for turn on light at Level 1: 44

d'3:n '0n l'aIt
dZ'3:n '0n l'aIt
g'3:n '0n l'aIt
Z'3:n '0n l'aIt
z'3:n '0n l'aIt
w'3:n '0n l'aIt

T'3:n '0n l'aIt
 D'3:n '0n l'aIt
 S'3:n '0n l'aIt
 bl'3:n '0n l'aIt
 br'3:n '0n l'aIt
 kr'3:n '0n l'aIt
 dr'3:n '0n l'aIt
 fl'3:n '0n l'aIt
 gl'3:n '0n l'aIt
 gr'3:n '0n l'aIt
 pl'3:n '0n l'aIt
 pr'3:n '0n l'aIt
 sk'3:n '0n l'aIt
 sl'3:n '0n l'aIt
 sn'3:n '0n l'aIt
 tr'3:n '0n l'aIt
 spr'3:n '0n l'aIt
 spl'3:n '0n l'aIt
 skw'3:n '0n l'aIt
 Tr'3:n '0n l'aIt
 t'3:n f'0n l'aIt
 t'3:n v'0n l'aIt
 t'3:n tS'0n l'aIt
 t'3:n bl'0n l'aIt
 t'3:n pl'0n l'aIt
 t'3:n tr'0n l'aIt
 t'3:n Tr'0n l'aIt
 t'3:n '0n p'aIt
 t'3:n '0n d'aIt
 t'3:n '0n g'aIt
 t'3:n '0n Z'aIt
 t'3:n '0n v'aIt
 t'3:n '0n z'aIt
 t'3:n '0n j'aIt
 t'3:n '0n T'aIt
 t'3:n '0n D'aIt
 t'3:n '0n S'aIt
 t'3:n '0n gl'aIt

Wins for turn on light at Level 2: 46

pr'3:n '0n p'aIt
 tr'3:n f'0n l'aIt
 tr'3:n '0n g'aIt
 D'3:n bl'0n l'aIt
 t'3:n tr'0n p'aIt
 t'3:n Tr'0n T'aIt
 fl'3:n '0n S'aIt
 bl'3:n '0n d'aIt
 pr'3:n '0n T'aIt
 gr'3:n tS'0n l'aIt
 br'3:n '0n gl'aIt
 d'3:n '0n p'aIt
 gl'3:n '0n p'aIt
 br'3:n '0n Z'aIt
 d'3:n '0n g'aIt
 spl'3:n '0n g'aIt
 t'3:n tr'0n j'aIt
 pl'3:n '0n v'aIt
 br'3:n '0n j'aIt
 z'3:n bl'0n l'aIt
 D'3:n '0n D'aIt
 t'3:n tS'0n j'aIt
 dr'3:n '0n gl'aIt
 T'3:n '0n S'aIt
 dr'3:n '0n p'aIt
 gl'3:n '0n Z'aIt
 pl'3:n tS'0n l'aIt
 tr'3:n '0n gl'aIt
 bl'3:n tS'0n l'aIt
 spl'3:n '0n Z'aIt
 Tr'3:n '0n z'aIt
 pr'3:n '0n j'aIt
 w'3:n '0n Z'aIt
 w'3:n '0n j'aIt
 gl'3:n '0n j'aIt
 g'3:n tr'0n l'aIt
 Z'3:n '0n S'aIt
 t'3:n tS'0n g'aIt
 tr'3:n Tr'0n l'aIt
 S'3:n Tr'0n l'aIt
 spl'3:n '0n T'aIt
 T'3:n '0n p'aIt
 D'3:n v'0n l'aIt
 z'3:n tS'0n l'aIt
 bl'3:n '0n g'aIt
 skw'3:n '0n D'aIt

Wins for turn on light at Level 3: 65

dr'3:n tS'0n gl'aIt
 D'3:n bl'0n T'aIt
 pr'3:n tr'0n p'aIt
 S'3:n Tr'0n D'aIt

spl'3:n Tr'On T'aIt
dr'3:n tr'On j'aIt
z'3:n ts'On gl'aIt
z'3:n bl'On p'aIt
gr'3:n ts'On d'aIt
S'3:n Tr'On g'aIt
g'3:n tr'On T'aIt
z'3:n bl'On v'aIt
d'3:n ts'On g'aIt
z'3:n ts'On v'aIt
bl'3:n ts'On d'aIt
br'3:n tr'On j'aIt
gr'3:n Tr'On T'aIt
pr'3:n tr'On j'aIt
w'3:n ts'On j'aIt
bl'3:n ts'On j'aIt
g'3:n ts'On j'aIt
T'3:n ts'On j'aIt
gl'3:n ts'On g'aIt
gr'3:n tr'On j'aIt
Z'3:n ts'On j'aIt
gr'3:n ts'On gl'aIt
z'3:n bl'On gl'aIt
D'3:n bl'On g'aIt
br'3:n ts'On g'aIt
z'3:n Tr'On T'aIt
z'3:n bl'On T'aIt
spl'3:n Tr'On g'aIt
w'3:n Tr'On j'aIt
S'3:n Tr'On v'aIt
spr'3:n Tr'On T'aIt
tr'3:n ts'On g'aIt
tr'3:n f'On j'aIt
T'3:n ts'On S'aIt
z'3:n ts'On T'aIt
z'3:n bl'On g'aIt
br'3:n Tr'On gl'aIt
gr'3:n ts'On Z'aIt
d'3:n Tr'On p'aIt
br'3:n ts'On Z'aIt
z'3:n bl'On d'aIt
S'3:n Tr'On Z'aIt
S'3:n tr'On j'aIt
D'3:n ts'On g'aIt
T'3:n bl'On p'aIt
bl'3:n bl'On d'aIt
D'3:n v'On Z'aIt
tr'3:n Tr'On v'aIt
pr'3:n ts'On g'aIt
z'3:n bl'On z'aIt
D'3:n f'On D'aIt
z'3:n bl'On Z'aIt
pr'3:n tr'On T'aIt
pr'3:n Tr'On j'aIt
pl'3:n ts'On j'aIt
pl'3:n ts'On g'aIt
D'3:n bl'On d'aIt
D'3:n v'On z'aIt
D'3:n v'On d'aIt
dZ'3:n tr'On j'aIt
Tr'3:n ts'On j'aIt

Results for "turn off light":
successes: 185
failures: 141

Wins for turn off light at Level 1: 52

d'3:n 'Of l'aIt
dZ'3:n 'Of l'aIt
g'3:n 'Of l'aIt
Z'3:n 'Of l'aIt
z'3:n 'Of l'aIt
n'3:n 'Of l'aIt
w'3:n 'Of l'aIt
T'3:n 'Of l'aIt
D'3:n 'Of l'aIt
S'3:n 'Of l'aIt
bl'3:n 'Of l'aIt
br'3:n 'Of l'aIt
kl'3:n 'Of l'aIt
dr'3:n 'Of l'aIt
fl'3:n 'Of l'aIt
gl'3:n 'Of l'aIt
gr'3:n 'Of l'aIt
pl'3:n 'Of l'aIt
pr'3:n 'Of l'aIt
kw'3:n 'Of l'aIt
sk'3:n 'Of l'aIt
sl'3:n 'Of l'aIt
sw'3:n 'Of l'aIt
tr'3:n 'Of l'aIt

spr'3:n '0f l'aIt
spl'3:n '0f l'aIt
skr'3:n '0f l'aIt
str'3:n '0f l'aIt
Tr'3:n '0f l'aIt
t'3:n p'0f l'aIt
t'3:n t'0f l'aIt
t'3:n b'0f l'aIt
t'3:n dz'0f l'aIt
t'3:n v'0f l'aIt
t'3:n l'0f l'aIt
t'3:n j'0f l'aIt
t'3:n n'0f l'aIt
t'3:n tS'0f l'aIt
t'3:n T'0f l'aIt
t'3:n bl'0f l'aIt
t'3:n sl'0f l'aIt
t'3:n tw'0f l'aIt
t'3:n '0f Z'aIt
t'3:n '0f v'aIt
t'3:n '0f z'aIt
t'3:n '0f j'aIt
t'3:n '0f tS'aIt
t'3:n '0f T'aIt
t'3:n '0f D'aIt
t'3:n '0f S'aIt
t'3:n '0f gl'aIt
t'3:n '0f spl'aIt

Wins for turn off light at Level 2: 50

sk'3:n '0f j'aIt
z'3:n '0f tS'aIt
n'3:n '0f tS'aIt
dr'3:n '0f S'aIt
g'3:n j'0f l'aIt
t'3:n bl'0f j'aIt
kw'3:n '0f gl'aIt
spl'3:n '0f tS'aIt
bl'3:n '0f gl'aIt
t'3:n T'0f j'aIt
Z'3:n '0f v'aIt
sl'3:n '0f Z'aIt
t'3:n dz'0f v'aIt
br'3:n T'0f l'aIt
sk'3:n '0f Z'aIt
bl'3:n p'0f l'aIt
z'3:n dz'0f l'aIt
sw'3:n '0f j'aIt
kw'3:n n'0f l'aIt
gl'3:n b'0f l'aIt
w'3:n '0f j'aIt
D'3:n j'0f l'aIt
spl'3:n '0f spl'aIt
br'3:n l'0f l'aIt
sk'3:n bl'0f l'aIt
sk'3:n '0f v'aIt
tr'3:n v'0f l'aIt
t'3:n b'0f j'aIt
pr'3:n b'0f l'aIt
Z'3:n dz'0f l'aIt
t'3:n dz'0f gl'aIt
kw'3:n '0f D'aIt
z'3:n '0f j'aIt
dZ'3:n j'0f l'aIt
t'3:n p'0f T'aIt
t'3:n bl'0f T'aIt
Tr'3:n dz'0f l'aIt
dr'3:n tS'0f l'aIt
fl'3:n '0f z'aIt
dZ'3:n p'0f l'aIt
gr'3:n t'0f l'aIt
spl'3:n '0f D'aIt
fl'3:n '0f j'aIt
br'3:n '0f D'aIt
t'3:n tS'0f Z'aIt
sl'3:n '0f gl'aIt
S'3:n n'0f l'aIt
T'3:n v'0f l'aIt
Z'3:n tS'0f l'aIt
spr'3:n j'0f l'aIt

Wins for turn off light at Level 3: 83

bl'3:n T'0f gl'aIt
n'3:n T'0f j'aIt
dZ'3:n j'0f spl'aIt
dZ'3:n p'0f spl'aIt
Z'3:n p'0f v'aIt
T'3:n bl'0f j'aIt
D'3:n dz'0f v'aIt
dZ'3:n p'0f v'aIt
bl'3:n p'0f S'aIt
gl'3:n bl'0f T'aIt

sl'3:n bl'0f j'aIt
 z'3:n tw'0f ts'aIt
 z'3:n l'0f j'aIt
 bl'3:n p'0f gl'aIt
 sl'3:n b'0f j'aIt
 tr'3:n b'0f j'aIt
 sl'3:n n'0f Z'aIt
 gr'3:n t'0f v'aIt
 T'3:n v'0f T'aIt
 w'3:n p'0f j'aIt
 kw'3:n n'0f j'aIt
 bl'3:n ts'0f Z'aIt
 D'3:n j'0f D'aIt
 br'3:n ts'0f Z'aIt
 dr'3:n ts'0f j'aIt
 z'3:n p'0f j'aIt
 gl'3:n b'0f D'aIt
 gl'3:n T'0f j'aIt
 Z'3:n dz'0f j'aIt
 br'3:n p'0f T'aIt
 D'3:n j'0f j'aIt
 pr'3:n b'0f T'aIt
 sk'3:n bl'0f v'aIt
 skr'3:n T'0f j'aIt
 g'3:n bl'0f j'aIt
 dz'3:n j'0f gl'aIt
 spl'3:n p'0f D'aIt
 pl'3:n p'0f T'aIt
 g'3:n j'0f spl'aIt
 br'3:n l'0f v'aIt
 T'3:n bl'0f T'aIt
 g'3:n dz'0f gl'aIt
 spr'3:n j'0f spl'aIt
 z'3:n dz'0f j'aIt
 n'3:n ts'0f ts'aIt
 z'3:n p'0f T'aIt
 g'3:n b'0f j'aIt
 dz'3:n j'0f v'aIt
 Z'3:n tw'0f v'aIt
 spr'3:n bl'0f j'aIt
 pl'3:n T'0f j'aIt
 S'3:n n'0f T'aIt
 dr'3:n ts'0f gl'aIt
 Tr'3:n p'0f T'aIt
 dz'3:n p'0f D'aIt
 T'3:n p'0f T'aIt
 fl'3:n T'0f j'aIt
 spr'3:n j'0f D'aIt
 T'3:n v'0f gl'aIt
 br'3:n l'0f Z'aIt
 gl'3:n b'0f gl'aIt
 S'3:n n'0f spl'aIt
 br'3:n bl'0f T'aIt
 spl'3:n T'0f j'aIt
 Tr'3:n ts'0f Z'aIt
 S'3:n dz'0f gl'aIt
 pr'3:n b'0f j'aIt
 sl'3:n T'0f gl'aIt
 sl'3:n p'0f Z'aIt
 br'3:n l'0f T'aIt
 dz'3:n j'0f T'aIt
 Z'3:n ts'0f T'aIt
 sk'3:n b'0f j'aIt
 sl'3:n bl'0f Z'aIt
 sk'3:n T'0f Z'aIt
 bl'3:n bl'0f T'aIt
 Z'3:n dz'0f T'aIt
 S'3:n n'0f z'aIt
 sk'3:n bl'0f gl'aIt
 br'3:n t'0f D'aIt
 S'3:n b'0f j'aIt
 tr'3:n v'0f z'aIt
 sl'3:n j'0f Z'aIt

Results for "turn light red":
 successes: 3
 failures: 90

Wins for turn light red at Level 1: 3
 S'3:n l'aIt r'Ed
 br'3:n l'aIt r'Ed
 Tr'3:n l'aIt r'Ed

Wins for turn light red at Level 2: 0

Wins for turn light red at Level 3: 0

Results for "turn light blue":
 successes: 166
 failures: 94

Wins for turn light blue at Level 1: 41

dZ'3:n l'aIt bl'u:
g'3:n l'aIt bl'u:
s'3:n l'aIt bl'u:
Z'3:n l'aIt bl'u:
z'3:n l'aIt bl'u:
n'3:n l'aIt bl'u:
w'3:n l'aIt bl'u:
T'3:n l'aIt bl'u:
D'3:n l'aIt bl'u:
S'3:n l'aIt bl'u:
bl'3:n l'aIt bl'u:
br'3:n l'aIt bl'u:
kl'3:n l'aIt bl'u:
kr'3:n l'aIt bl'u:
dr'3:n l'aIt bl'u:
fl'3:n l'aIt bl'u:
gl'3:n l'aIt bl'u:
gr'3:n l'aIt bl'u:
pl'3:n l'aIt bl'u:
pr'3:n l'aIt bl'u:
kw'3:n l'aIt bl'u:
sk'3:n l'aIt bl'u:
sl'3:n l'aIt bl'u:
tr'3:n l'aIt bl'u:
spr'3:n l'aIt bl'u:
spl'3:n l'aIt bl'u:
skr'3:n l'aIt bl'u:
str'3:n l'aIt bl'u:
Tr'3:n l'aIt bl'u:
t'3:n p'aIt bl'u:
t'3:n d'aIt bl'u:
t'3:n g'aIt bl'u:
t'3:n Z'aIt bl'u:
t'3:n v'aIt bl'u:
t'3:n z'aIt bl'u:
t'3:n j'aIt bl'u:
t'3:n T'aIt bl'u:
t'3:n D'aIt bl'u:
t'3:n S'aIt bl'u:
t'3:n gl'aIt bl'u:
t'3:n l'aIt v'u:

Wins for turn light blue at Level 2: 62

w'3:n j'aIt bl'u:
kl'3:n l'aIt v'u:
gr'3:n Z'aIt bl'u:
pl'3:n S'aIt bl'u:
D'3:n gl'aIt bl'u:
gr'3:n D'aIt bl'u:
s'3:n v'aIt bl'u:
kl'3:n D'aIt bl'u:
spr'3:n j'aIt bl'u:
sk'3:n l'aIt v'u:
dr'3:n v'aIt bl'u:
w'3:n d'aIt bl'u:
sk'3:n g'aIt bl'u:
br'3:n p'aIt bl'u:
kl'3:n Z'aIt bl'u:
gl'3:n z'aIt bl'u:
sl'3:n Z'aIt bl'u:
t'3:n T'aIt v'u:
sl'3:n p'aIt bl'u:
str'3:n Z'aIt bl'u:
pl'3:n T'aIt bl'u:
fl'3:n gl'aIt bl'u:
br'3:n d'aIt bl'u:
S'3:n v'aIt bl'u:
str'3:n D'aIt bl'u:
S'3:n d'aIt bl'u:
gr'3:n T'aIt bl'u:
dZ'3:n S'aIt bl'u:
Z'3:n v'aIt bl'u:
pl'3:n g'aIt bl'u:
kr'3:n S'aIt bl'u:
br'3:n Z'aIt bl'u:
z'3:n z'aIt bl'u:
w'3:n T'aIt bl'u:
sk'3:n T'aIt bl'u:
fl'3:n v'aIt bl'u:
skr'3:n S'aIt bl'u:
spl'3:n d'aIt bl'u:
gl'3:n g'aIt bl'u:
dZ'3:n Z'aIt bl'u:
fl'3:n S'aIt bl'u:
kl'3:n T'aIt bl'u:
bl'3:n j'aIt bl'u:
kr'3:n g'aIt bl'u:
gl'3:n Z'aIt bl'u:
g'3:n gl'aIt bl'u:
gr'3:n p'aIt bl'u:

spr'3:n g'aIt bl'u:
tr'3:n Z'aIt bl'u:
Tr'3:n v'aIt bl'u:
spr'3:n D'aIt bl'u:
s'3:n j'aIt bl'u:
dr'3:n d'aIt bl'u:
z'3:n d'aIt bl'u:
skr'3:n Z'aIt bl'u:
spl'3:n g'aIt bl'u:
spl'3:n gl'aIt bl'u:
z'3:n gl'aIt bl'u:
str'3:n g'aIt bl'u:
dr'3:n p'aIt bl'u:
spr'3:n S'aIt bl'u:
str'3:n T'aIt bl'u:

Wins for turn light blue at Level 3: 63

spl'3:n d'aIt v'u:
kl'3:n D'aIt v'u:
n'3:n T'aIt v'u:
br'3:n Z'aIt v'u:
tr'3:n T'aIt v'u:
s'3:n T'aIt v'u:
S'3:n d'aIt v'u:
z'3:n d'aIt v'u:
spr'3:n g'aIt v'u:
Z'3:n v'aIt v'u:
sk'3:n j'aIt v'u:
w'3:n j'aIt v'u:
pr'3:n T'aIt v'u:
sk'3:n g'aIt v'u:
g'3:n T'aIt v'u:
w'3:n d'aIt v'u:
D'3:n gl'aIt v'u:
z'3:n z'aIt v'u:
kl'3:n g'aIt v'u:
sk'3:n gl'aIt v'u:
spr'3:n D'aIt v'u:
D'3:n T'aIt v'u:
Z'3:n T'aIt v'u:
br'3:n d'aIt v'u:
bl'3:n T'aIt v'u:
kl'3:n T'aIt v'u:
dZ'3:n Z'aIt v'u:
pl'3:n g'aIt v'u:
bl'3:n j'aIt v'u:
gr'3:n T'aIt v'u:
s'3:n j'aIt v'u:
gl'3:n T'aIt v'u:
dZ'3:n T'aIt v'u:
str'3:n D'aIt v'u:
gl'3:n z'aIt v'u:
spl'3:n g'aIt v'u:
str'3:n g'aIt v'u:
T'3:n T'aIt v'u:
kl'3:n v'aIt v'u:
gr'3:n D'aIt v'u:
dZ'3:n S'aIt v'u:
kl'3:n Z'aIt v'u:
dr'3:n d'aIt v'u:
gl'3:n Z'aIt v'u:
gl'3:n g'aIt v'u:
kl'3:n j'aIt v'u:
sk'3:n d'aIt v'u:
sl'3:n p'aIt v'u:
S'3:n T'aIt v'u:
spr'3:n T'aIt v'u:
br'3:n T'aIt v'u:
dr'3:n T'aIt v'u:
z'3:n T'aIt v'u:
S'3:n v'aIt v'u:
sl'3:n T'aIt v'u:
kl'3:n d'aIt v'u:
tr'3:n Z'aIt v'u:
spr'3:n j'aIt v'u:
Tr'3:n v'aIt v'u:
pl'3:n T'aIt v'u:
sk'3:n v'aIt v'u:
sl'3:n Z'aIt v'u:
dr'3:n v'aIt v'u:

A.3.2 Over-the-Air Results

Nonsense word sequences for human tests, TARGET COMMAND: who am i
FULL-MANGLED: T'eI g'Uk@L f'u: g'am Z'aI

```
ON_CONVERSATION_TURN_STARTED
ON_END_OF_UTTERANCE
ON_END_OF_UTTERANCE
ON_RECOGNIZING_SPEECH_FINISHED:
  {"text": "new game is I"}
ON_RESPONDING_STARTED:
  {"is_error_response": false}
ON_RESPONDING_FINISHED
ON_CONVERSATION_TURN_FINISHED:
  {"with_follow_on_turn": false}
```

Nonsense word sequences for human tests, target command: who am i, NO GOOGLE
Level 3: T'eI g'Uk@L D'u: 'am j'aI

```
ON_CONVERSATION_TURN_STARTED
ON_END_OF_UTTERANCE
ON_END_OF_UTTERANCE
ON_RECOGNIZING_SPEECH_FINISHED:
  {"text": "who am I I"}
ON_RESPONDING_STARTED:
  {"is_error_response": false}
ON_RENDERER_RESPONSE:
  {"text": "Okay. Let's get Who am I.", "type": RenderResponseType.TEXT}
ON_RESPONDING_FINISHED
ON_RENDERER_RESPONSE:
  {
    "text": "Glad you're up for whoami. This is going to be fun. I'm here with you every step of the way. Hope you're feeling lucky. Meet you",
    "type": RenderResponseType.TEXT
  }
}
```

Nonsense word sequences for human tests, target command: whats my name, NO GOOGLE

Level 1: h'eI w'u:b@L w'0ts bl'aI n'eIm

Level 2: h'eI w'u:b@L w'0ts gr'aI Z'eIm

Level 3: h'eI w'u:b@L D'0ts sn'aI z'eIm

```
ON_CONVERSATION_TURN_STARTED
ON_END_OF_UTTERANCE
ON_END_OF_UTTERANCE
ON_RECOGNIZING_SPEECH_FINISHED:
  {"text": "what's my name"}
ON_RESPONDING_STARTED:
  {"is_error_response": false}
ON_RESPONDING_FINISHED
ON_CONVERSATION_TURN_FINISHED:
  {"with_follow_on_turn": false}
```

```
ON_CONVERSATION_TURN_STARTED
ON_END_OF_UTTERANCE
ON_END_OF_UTTERANCE
ON_RECOGNIZING_SPEECH_FINISHED:
  {"text": "what's my name"}
ON_RESPONDING_STARTED:
  {"is_error_response": false}
ON_RESPONDING_FINISHED
ON_CONVERSATION_TURN_FINISHED:
  {"with_follow_on_turn": false}
```

```
ON_CONVERSATION_TURN_STARTED
ON_END_OF_UTTERANCE
ON_END_OF_UTTERANCE
ON_RECOGNIZING_SPEECH_FINISHED:
  {"text": "what's my name"}
ON_RESPONDING_STARTED:
  {"is_error_response": false}
ON_RESPONDING_FINISHED
ON_CONVERSATION_TURN_FINISHED:
  {"with_follow_on_turn": false}
```

Nonsense word sequences for human tests, target command: turn on light, NO GOOGLE
Level 1: h'eI g'u:b@L dZ'3:n '0n l'aIt

```
ON_CONVERSATION_TURN_STARTED
ON_END_OF_UTTERANCE
ON_END_OF_UTTERANCE
ON_RECOGNIZING_SPEECH_FINISHED:
  {"text": "turn on light"}
ON_RESPONDING_STARTED:
  {"is_error_response": false}
ON_RESPONDING_FINISHED
ON_CONVERSATION_TURN_FINISHED:
  {"with_follow_on_turn": false}
```

A.3.3 Human Comprehensibility Results

Target command: Turn off light
Condition: Google revealed first

Participant ID: 5c4d07f8b00d500001e52047
Timestamp: Sun Jan 27 2019 11:29:13 GMT-0500 (Eastern Standard Time)
Native language: English
Participant ID: 5c4d07f8b00d500001e52047
Timestamp: Sun Jan 27 2019 11:29:19 GMT-0500 (Eastern Standard Time)
Meanings assigned to audio clips:
FULL_MANGLED_ONE.wav: no meaning
Level_4_ONE.wav: no meaning
Level_3_ONE.wav: no meaning
Level_2_ONE.wav: they gurgle turn off light
NOT_MANGLED_ONE.wav: Hey Google, Turn off light
Participant ID: 5c4d07f8b00d500001e52047
Timestamp: Sun Jan 27 2019 11:30:24 GMT-0500 (Eastern Standard Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5c4dcedf620ba700010c3f64
Timestamp: Sun Jan 27 2019 16:26:02 GMT+0000 (Greenwich Mean Time)
Native language: english
Participant ID: 5c4dcedf620ba700010c3f64
Timestamp: Sun Jan 27 2019 16:26:17 GMT+0000 (Greenwich Mean Time)
Meanings assigned to audio clips:
FULL_MANGLED_ONE.wav: n/a
Level_4_ONE.wav: NO MEANING
Level_3_ONE.wav: NO MEANING
Level_2_ONE.wav: NO MEANING
NOT_MANGLED_ONE.wav: Hey Google, turn off light.
Participant ID: 5c4dcedf620ba700010c3f64
Timestamp: Sun Jan 27 2019 16:28:32 GMT+0000 (Greenwich Mean Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5c4be2759233b80001b2ed8a
Timestamp: Sun Jan 27 2019 11:31:52 GMT-0500 (Eastern Standard Time)
Native language: English
Participant ID: 5c4be2759233b80001b2ed8a
Timestamp: Sun Jan 27 2019 11:32:09 GMT-0500 (Eastern Standard Time)
Meanings assigned to audio clips:
FULL_MANGLED_ONE.wav: No Meaning
Level_4_ONE.wav: NO MEANING
Level_3_ONE.wav: turn off light
Level_2_ONE.wav: turn off light
NOT_MANGLED_ONE.wav: turn off light
Participant ID: 5c4be2759233b80001b2ed8a
Timestamp: Sun Jan 27 2019 11:34:11 GMT-0500 (Eastern Standard Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5c40fdc98daebb0001d199cd
Timestamp: Sun Jan 27 2019 16:32:20 GMT+0000 (GMT)
Native language: English
Participant ID: 5c40fdc98daebb0001d199cd
Timestamp: Sun Jan 27 2019 16:32:34 GMT+0000 (GMT)
Meanings assigned to audio clips:
FULL_MANGLED_ONE.wav: No meaning
Level_4_ONE.wav: Google
Level_3_ONE.wav: Light
Level_2_ONE.wav: Google turn off light
NOT_MANGLED_ONE.wav: Hey google turn off light
Participant ID: 5c40fdc98daebb0001d199cd
Timestamp: Sun Jan 27 2019 16:34:31 GMT+0000 (GMT)
Study withdrawn: false
Study submitted: true

Participant ID: 5c4d12cbb00d500001e521b5
Timestamp: Sun Jan 27 2019 11:32:31 GMT-0500 (Eastern Standard Time)
Native language:
Participant ID: 5c4d12cbb00d500001e521b5
Timestamp: Sun Jan 27 2019 11:32:43 GMT-0500 (Eastern Standard Time)
Meanings assigned to audio clips:
FULL_MANGLED_ONE.wav: no meaning
Level_4_ONE.wav: sit google dont laugh
Level_3_ONE.wav: vain google dont laugh might
Level_2_ONE.wav: sit google truth arse light
NOT_MANGLED_ONE.wav: hey google turn off light
Participant ID: 5c4d12cbb00d500001e521b5
Timestamp: Sun Jan 27 2019 11:34:57 GMT-0500 (Eastern Standard Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5c4dbcba620ba700010c3c67

Timestamp: Sun Jan 27 2019 16:34:57 GMT+0000 (Greenwich Mean Time)
Native language: English
Participant ID: 5c4dbcb620ba700010c3c67
Timestamp: Sun Jan 27 2019 16:35:12 GMT+0000 (Greenwich Mean Time)
Meanings assigned to audio clips:
FULL_MANGLED_ONE.wav: NO MEANING
Level_4_ONE.wav: NO MEANING
Level_3_ONE.wav: NO MEANING
Level_2_ONE.wav: They Google Turn Off Light
NOT_MANGLED_ONE.wav: Hey Google Turn Off Light
Participant ID: 5c4dbcb620ba700010c3c67
Timestamp: Sun Jan 27 2019 16:36:52 GMT+0000 (Greenwich Mean Time)
Study withdrawn: false
Study submitted: true

Target command: Turn off light
Condition: Google revealed last

Participant ID: 5c4d8af20147e800011a861f
Timestamp: Mon Jan 28 2019 10:09:04 GMT+0000 (Greenwich Mean Time)
Native language: English
Participant ID: 5c4d8af20147e800011a861f
Timestamp: Mon Jan 28 2019 10:09:16 GMT+0000 (Greenwich Mean Time)
Meanings assigned to audio clips:
FULL_MANGLED_THREE.wav: No meaning
Level_4_THREE.wav: No meaning sounds German
Level_3_THREE.wav: Off light at end
Level_2_THREE.wav: Hey guzunt Turn off light
NOT_MANGLED_THREE.wav: Hey Google turn off light
Participant ID: 5c4d8af20147e800011a861f
Timestamp: Mon Jan 28 2019 10:10:57 GMT+0000 (Greenwich Mean Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5c430f6489f63e
Timestamp: Mon Jan 28 2019 10:09:38 GMT+0000 (GMT)
Native language: English
Participant ID: 5c430f6489f63e
Timestamp: Mon Jan 28 2019 10:09:52 GMT+0000 (GMT)
Meanings assigned to audio clips:
FULL_MANGLED_THREE.wav: No meaning
Level_4_THREE.wav: No meaning
Level_3_THREE.wav: Turn off light
Level_2_THREE.wav: Off light
NOT_MANGLED_THREE.wav: Hey google turn off light

Participant ID: 5c4a061161db8c00015e4fb0
Timestamp: Mon Jan 28 2019 10:12:53 GMT+0000 (Greenwich Mean Time)
Native language: English
Participant ID: 5c4a061161db8c00015e4fb0
Timestamp: Mon Jan 28 2019 10:13:07 GMT+0000 (Greenwich Mean Time)
Meanings assigned to audio clips:
FULL_MANGLED_THREE.wav: NO MEANING
Level_4_THREE.wav:
Level_3_THREE.wav: NO MEANING
Level_2_THREE.wav: NO MEANING
NOT_MANGLED_THREE.wav: HEY GOOGLE TURN OFF LIGHT
Participant ID: 5c4a061161db8c00015e4fb0
Timestamp: Mon Jan 28 2019 10:14:32 GMT+0000 (Greenwich Mean Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5c4ea7c3889752000156ddc5
Timestamp: Mon Jan 28 2019 02:13:07 GMT-0800 (Pacific Standard Time)
Native language: English
Participant ID: 5c4ea7c3889752000156ddc5
Timestamp: Mon Jan 28 2019 02:13:15 GMT-0800 (Pacific Standard Time)
Meanings assigned to audio clips:
FULL_MANGLED_THREE.wav: NO MEANING
Level_4_THREE.wav: Hey
Level_3_THREE.wav: Turn Off Light
Level_2_THREE.wav: Hey Could You Turn Off Light
NOT_MANGLED_THREE.wav: Hey Google Turn Off Light
Participant ID: 5c4ea7c3889752000156ddc5
Timestamp: Mon Jan 28 2019 02:15:30 GMT-0800 (Pacific Standard Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5c4d6b17ad4660000191441e
Timestamp: Mon Jan 28 2019 10:13:07 GMT+0000 (GMT)
Meanings assigned to audio clips:
FULL_MANGLED_THREE.wav: No meaning
Level_4_THREE.wav: No meaning
Level_3_THREE.wav: No meaning
Level_2_THREE.wav: Hey turn off light

NOT_MANGLED_THREE.wav:

Participant ID: 5a6b057a0384310001944dfc
Timestamp: Mon Jan 28 2019 21:15:41 GMT+1100 (Australian Eastern Daylight Time)
Native language: English
Participant ID: 5a6b057a0384310001944dfc
Timestamp: Mon Jan 28 2019 21:16:04 GMT+1100 (Australian Eastern Daylight Time)
Meanings assigned to audio clips:
FULL_MANGLED_THREE.wav: No meaning
Level_4_THREE.wav: No meaning
Level_3_THREE.wav: No meaning
Level_2_THREE.wav: Hey
NOT_MANGLED_THREE.wav: Hey google turn off light
Participant ID: 5a6b057a0384310001944dfc
Timestamp: Mon Jan 28 2019 21:18:56 GMT+1100 (Australian Eastern Daylight Time)
Study withdrawn: false
Study submitted: true

Target command: What's my name
Condition: Google revealed first

Participant ID: 5c4d8af20147e800011a861f
Timestamp: Mon Jan 28 2019 10:54:13 GMT+0000 (Greenwich Mean Time)
Native language: English
Participant ID: 5c4d8af20147e800011a861f
Timestamp: Mon Jan 28 2019 10:54:19 GMT+0000 (Greenwich Mean Time)
Meanings assigned to audio clips:
FULL_MANGLED_FOUR.wav: No meaning
Level_4_FOUR.wav: No meaning
Level_3_FOUR.wav: Ja Google what's prime jain
Level_2_FOUR.wav: Ja Google what's bly name
NOT_MANGLED_FOUR.wav: Hey Google what's my name
Participant ID: 5c4d8af20147e800011a861f
Timestamp: Mon Jan 28 2019 10:55:31 GMT+0000 (Greenwich Mean Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5c4d6b17ad4660000191441e
Timestamp: Mon Jan 28 2019 10:54:57 GMT+0000 (GMT)
Meanings assigned to audio clips:
FULL_MANGLED_FOUR.wav: No meaning
Level_4_FOUR.wav: No meaning
Level_3_FOUR.wav: No meaning
Level_2_FOUR.wav: No meaning
NOT_MANGLED_FOUR.wav: Hey google what's my name

Participant ID: 5c4d8307782a2500019f5bfb
Timestamp: Mon Jan 28 2019 10:54:35 GMT+0000 (GMT)
Meanings assigned to audio clips:
FULL_MANGLED_FOUR.wav: Another language?
Level_4_FOUR.wav: Google
Level_3_FOUR.wav: Google
Level_2_FOUR.wav: Google what's my name
NOT_MANGLED_FOUR.wav: Hey google what's my name
Participant ID: 5c4d8307782a2500019f5bfb
Timestamp: Mon Jan 28 2019 10:56:39 GMT+0000 (GMT)
Study withdrawn: false
Study submitted: true

Participant ID: 5c430f6489f63e000179cd16
Timestamp: Mon Jan 28 2019 11:00:34 GMT+0000 (GMT)
Native language: English
Participant ID: 5c430f6489f63e000179cd16
Timestamp: Mon Jan 28 2019 11:00:47 GMT+0000 (GMT)
Meanings assigned to audio clips:
FULL_MANGLED_FOUR.wav: No meaning
Level_4_FOUR.wav: No meaning
Level_3_FOUR.wav: Google
Level_2_FOUR.wav: Google what's my name
NOT_MANGLED_FOUR.wav: Hey google what's my name
Participant ID: 5c430f6489f63e000179cd16
Timestamp: Mon Jan 28 2019 11:02:59 GMT+0000 (GMT)
Study withdrawn: false
Study submitted: true

Participant ID: 5c4892ba865f540001e902f4
Timestamp: Mon Jan 28 2019 05:58:10 GMT-0500 (EST)
Meanings assigned to audio clips:
FULL_MANGLED_FOUR.wav: No meaning
Level_4_FOUR.wav: No meaning
Level_3_FOUR.wav: No meaning
Level_2_FOUR.wav: What's my name?
NOT_MANGLED_FOUR.wav: Hey Google, what's my name?

Timestamp: Mon Jan 28 2019 11:08:50 GMT+0000 (Greenwich Mean Time)
Native language: English
Timestamp: Mon Jan 28 2019 11:09:09 GMT+0000 (Greenwich Mean Time)
Meanings assigned to audio clips:

FULL_MANGLED_FOUR.wav: Thats
Level_4_FOUR.wav: Google that's
Level_3_FOUR.wav: Goggle what's
Level_2_FOUR.wav: Goggle What's name
NOT_MANGLED_FOUR.wav: Hey google what's my name
Timestamp: Mon Jan 28 2019 11:11:24 GMT+0000 (Greenwich Mean Time)
Study withdrawn: false
Study submitted: true

Target command: What's my name
Condition: Google revealed last

Participant ID: 5c49e951bd42830001dbae9d
Timestamp: Sun Jan 27 2019 17:47:27 GMT-0500 (Eastern Standard Time)
Native language: English
Participant ID: 5c49e951bd42830001dbae9d
Timestamp: Sun Jan 27 2019 17:47:38 GMT-0500 (Eastern Standard Time)
Meanings assigned to audio clips:
FULL_MANGLED_TWO.wav: no meaning
Level_4_TWO.wav: HEY, THATS
Level_3_TWO.wav: HEY, WHAT'S
Level_2_TWO.wav: HEY GOOGLE WHAT'S MY NAME?
NOT_MANGLED_TWO.wav: HEY GOOGLE WHAT'S MY NAME?
Participant ID: 5c49e951bd42830001dbae9d
Timestamp: Sun Jan 27 2019 17:50:27 GMT-0500 (Eastern Standard Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5c42702c848dbb000157f7c1
Timestamp: Sun Jan 27 2019 14:48:49 GMT-0800 (Pacific Standard Time)
Native language: English
Participant ID: 5c42702c848dbb000157f7c1
Timestamp: Sun Jan 27 2019 14:48:57 GMT-0800 (Pacific Standard Time)
Meanings assigned to audio clips:
FULL_MANGLED_TWO.wav: no meaning
Level_4_TWO.wav: Hey, that's not the same.
Level_3_TWO.wav: Hey, cry, shame.
Level_2_TWO.wav: Hey, Google, what's my name.
NOT_MANGLED_TWO.wav: Hey Google what's my name.

Participant ID: 5bf5822f8dbe7f0001fa6237
Timestamp: Sun Jan 27 2019 22:49:16 GMT+0000 (GMT)
Native language: English
Participant ID: 5bf5822f8dbe7f0001fa6237
Timestamp: Sun Jan 27 2019 22:49:33 GMT+0000 (GMT)
Meanings assigned to audio clips:
FULL_MANGLED_TWO.wav: No meaning
Level_4_TWO.wav: Hey that's
Level_3_TWO.wav: Hey
Level_2_TWO.wav:
NOT_MANGLED_TWO.wav: Hey google what's my name
Participant ID: 5bf5822f8dbe7f0001fa6237
Timestamp: Sun Jan 27 2019 22:51:10 GMT+0000 (GMT)
Study withdrawn: false
Study submitted: true

Participant ID: 5c4e2113423a5000128add6
Timestamp: Sun Jan 27 2019 22:55:27 GMT+0000 (GMT Standard Time)
Native language: english
Participant ID: 5c4e2113423a5000128add6
Timestamp: Sun Jan 27 2019 22:55:44 GMT+0000 (GMT Standard Time)
Meanings assigned to audio clips:
FULL_MANGLED_TWO.wav: no meaning
Level_4_TWO.wav: no meaning
Level_3_TWO.wav: no meaning
Level_2_TWO.wav: no meaning
NOT_MANGLED_TWO.wav: hey google whats my name
Participant ID: 5c4e2113423a5000128add6
Timestamp: Sun Jan 27 2019 22:57:38 GMT+0000 (GMT Standard Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5c4b4e7ebcc3dc00015486ca
Timestamp: Sun Jan 27 2019 22:57:11 GMT+0000 (Greenwich Mean Time)
Native language: English
Participant ID: 5c4b4e7ebcc3dc00015486ca
Timestamp: Sun Jan 27 2019 22:57:24 GMT+0000 (Greenwich Mean Time)
Meanings assigned to audio clips:
FULL_MANGLED_TWO.wav: No meaning
Level_4_TWO.wav: Hey rubal zen
Level_3_TWO.wav: no meaning
Level_2_TWO.wav: hey rubel whats my name
NOT_MANGLED_TWO.wav: hey google whats my name
Participant ID: 5c4b4e7ebcc3dc00015486ca
Timestamp: Sun Jan 27 2019 22:58:51 GMT+0000 (Greenwich Mean Time)

Study withdrawn: false
Study submitted: true

Participant ID: 5c4d7e97c0783c00016915f5
Timestamp: Sun Jan 27 2019 23:09:43 GMT+0000 (GMT)
Meanings assigned to audio clips:
FULL_MANGLED_TWO.wav: No meaning
Level_4_TWO.wav: No meaning
Level_3_TWO.wav: No meaning
Level_2_TWO.wav: Hey Rubio what's my name
NOT_MANGLED_TWO.wav: Hey google what's my name?

Target command: Turn light blue
Condition: Google revealed first

Participant ID: 5c69c4acd9c5be000174432d
Timestamp: Sun Feb 17 2019 15:18:29 GMT-0600 (Central Standard Time)
Native language: English
Participant ID: 5c69c4acd9c5be000174432d
Timestamp: Sun Feb 17 2019 15:18:35 GMT-0600 (Central Standard Time)
Meanings assigned to audio clips:
FULL_MANGLED_SIX.wav: No Meaning
Level_4_SIX.wav: No Meaning
Level_3_SIX.wav: Blue
Level_2_SIX.wav: light blue
NOT_MANGLED_SIX.wav: Hey google turn light blue
Participant ID: 5c69c4acd9c5be000174432d
Timestamp: Sun Feb 17 2019 15:19:45 GMT-0600 (Central Standard Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5c64981ca98f0700010f8c81
Timestamp: Sun Feb 17 2019 21:20:34 GMT+0000 (Greenwich Mean Time)
Native language: English
Participant ID: 5c64981ca98f0700010f8c81
Timestamp: Sun Feb 17 2019 21:20:45 GMT+0000 (Greenwich Mean Time)
Meanings assigned to audio clips:
FULL_MANGLED_SIX.wav: no meaning
Level_4_SIX.wav: no meaning
Level_3_SIX.wav: no meaning
Level_2_SIX.wav: light blue
NOT_MANGLED_SIX.wav: hey google turn light blue
Participant ID: 5c64981ca98f0700010f8c81
Timestamp: Sun Feb 17 2019 21:22:09 GMT+0000 (Greenwich Mean Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5c68af5ffdeee50001b9cf01
Timestamp: Sun Feb 17 2019 16:20:57 GMT-0500 (Eastern Standard Time)
Native language: English
Participant ID: 5c68af5ffdeee50001b9cf01
Timestamp: Sun Feb 17 2019 16:21:22 GMT-0500 (Eastern Standard Time)
Meanings assigned to audio clips:
FULL_MANGLED_SIX.wav: NO MEANING
Level_4_SIX.wav: google
Level_3_SIX.wav: NO MEANING
Level_2_SIX.wav: Light Blue
NOT_MANGLED_SIX.wav: Hey Google turn light blue
Participant ID: 5c68af5ffdeee50001b9cf01
Timestamp: Sun Feb 17 2019 16:24:22 GMT-0500 (Eastern Standard Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5c69499c949182000143510f
Timestamp: Sun Feb 17 2019 21:21:12 GMT+0000 (Greenwich Mean Time)
Native language: English and Bangali
Participant ID: 5c69499c949182000143510f
Timestamp: Sun Feb 17 2019 21:21:44 GMT+0000 (Greenwich Mean Time)
Meanings assigned to audio clips:
FULL_MANGLED_SIX.wav: NO MEANING
Level_4_SIX.wav: NO MEANING
Level_3_SIX.wav: Google Blue
Level_2_SIX.wav: Google, Bright, Blue.
NOT_MANGLED_SIX.wav: Hey google turn light blue.
Participant ID: 5c69499c949182000143510f
Timestamp: Sun Feb 17 2019 21:24:03 GMT+0000 (Greenwich Mean Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5c699a71d0fb3700017312bd
Timestamp: Sun Feb 17 2019 21:20:52 GMT+0000 (GMT)
Meanings assigned to audio clips:
FULL_MANGLED_SIX.wav: No meaning
Level_4_SIX.wav: Google
Level_3_SIX.wav: Google blue
Level_2_SIX.wav: Google light blue
NOT_MANGLED_SIX.wav: Hey google turn light blue

Participant ID: 5c688dc0d6c7500001a03eea
Timestamp: Sun Feb 17 2019 21:22:02 GMT+0000 (GMT)
Meanings assigned to audio clips:
FULL_MANGLED_SIX.wav: No meaning
Level_4_SIX.wav:
Level_3_SIX.wav:
Level_2_SIX.wav: Light blue
NOT_MANGLED_SIX.wav: Hey google, turn light blue
Participant ID: 5c688dc0d6c7500001a03eea
Timestamp: Sun Feb 17 2019 21:23:23 GMT+0000 (GMT)
Study withdrawn: false
Study submitted: true

Target command: Turn light blue
Condition: Google revealed last

Participant ID: 5c695a065416680001bef3c0
Timestamp: Sun Feb 17 2019 13:41:22 GMT-0500 (Eastern Standard Time)
Native language: English
Participant ID: 5c695a065416680001bef3c0
Timestamp: Sun Feb 17 2019 13:41:29 GMT-0500 (Eastern Standard Time)
Meanings assigned to audio clips:
FULL_MANGLED_FIVE.wav: NO MEANING
Level_4_FIVE.wav: HEY
Level_3_FIVE.wav: Light Blue
Level_2_FIVE.wav: Hey Light Blue
NOT_MANGLED_FIVE.wav: Turn
Participant ID: 5c695a065416680001bef3c0
Timestamp: Sun Feb 17 2019 13:42:45 GMT-0500 (Eastern Standard Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5c492e3e4c3dc50001fbdb20
Timestamp: Sun Feb 17 2019 14:41:52 GMT-0400 (Atlantic Standard Time)
Native language: English
Participant ID: 5c492e3e4c3dc50001fbdb20
Timestamp: Sun Feb 17 2019 14:42:03 GMT-0400 (Atlantic Standard Time)
FULL_MANGLED_FIVE.wav: eight glistle eight fight food
Level_4_FIVE.wav: eight ??? chicken ??? food
Level_3_FIVE.wav: eight glisten ??? light blue
Level_2_FIVE.wav: hay ??? ??? light blue
NOT_MANGLED_FIVE.wav: hey google, turn light blue
Participant ID: 5c492e3e4c3dc50001fbdb20
Timestamp: Sun Feb 17 2019 14:43:43 GMT-0400 (Atlantic Standard Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5c67d802b65acc0001efc38f
Timestamp: Sun Feb 17 2019 18:42:31 GMT+0000 (Greenwich Mean Time)
Native language: English
Participant ID: 5c67d802b65acc0001efc38f
Timestamp: Sun Feb 17 2019 18:42:39 GMT+0000 (Greenwich Mean Time)
Meanings assigned to audio clips:
FULL_MANGLED_FIVE.wav: NO MEANING
Level_4_FIVE.wav: NO MEANING
Level_3_FIVE.wav: NO MEANING
Level_2_FIVE.wav: Light Blue
NOT_MANGLED_FIVE.wav: Turn light blue
Participant ID: 5c67d802b65acc0001efc38f
Timestamp: Sun Feb 17 2019 18:43:57 GMT+0000 (Greenwich Mean Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5c6888ac146af30001561096
Timestamp: Sun Feb 17 2019 18:42:46 GMT+0000 (GMT)
Native language: English
Participant ID: 5c6888ac146af30001561096
Timestamp: Sun Feb 17 2019 18:42:52 GMT+0000 (GMT)
Meanings assigned to audio clips:
FULL_MANGLED_FIVE.wav: No meaning
Level_4_FIVE.wav: No Meaning
Level_3_FIVE.wav: Light blue
Level_2_FIVE.wav: Light Blue
NOT_MANGLED_FIVE.wav: Churn Light Blue
Participant ID: 5c6888ac146af30001561096
Timestamp: Sun Feb 17 2019 18:43:45 GMT+0000 (GMT)
Study withdrawn: false
Study submitted: true

Participant ID: 5c66f38c8555e400015c78bf
Timestamp: Sun Feb 17 2019 13:46:47 GMT-0500 (Eastern Standard Time)
Native language: English
Participant ID: 5c66f38c8555e400015c78bf
Timestamp: Sun Feb 17 2019 13:47:07 GMT-0500 (Eastern Standard Time)
Meanings assigned to audio clips:
FULL_MANGLED_FIVE.wav: NO MEANING
Level_4_FIVE.wav: NO MEANING
Level_3_FIVE.wav: Light blue
Level_2_FIVE.wav: Hey jiggle light blue
NOT_MANGLED_FIVE.wav: Hey Google, turn light blue.

Participant ID: 5c66f38c8555e400015c78bf
Timestamp: Sun Feb 17 2019 13:49:34 GMT-0500 (Eastern Standard Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5c6970b347a5f100014a35fc
Timestamp: Sun Feb 17 2019 18:47:33 GMT+0000 (Greenwich Mean Time)
Native language: English
Participant ID: 5c6970b347a5f100014a35fc
Timestamp: Sun Feb 17 2019 18:47:46 GMT+0000 (Greenwich Mean Time)
Meanings assigned to audio clips:
FULL_MANGLED_FIVE.wav: No meaning
Level_4_FIVE.wav: Hey Jekyll June Right View
Level_3_FIVE.wav: Hey Glisten Light Blue
Level_2_FIVE.wav: Hey Jiggle Kern Right Blue
NOT_MANGLED_FIVE.wav: Hey Google Turn light blue
Participant ID: 5c6970b347a5f100014a35fc
Timestamp: Sun Feb 17 2019 18:50:16 GMT+0000 (Greenwich Mean Time)
Study withdrawn: false
Study submitted: true

Target command: Turn on light
Condition: Google revealed first

Participant ID: 5c69a2847d490900015b0fd3
Timestamp: Sun Feb 17 2019 22:44:08 GMT+0000 (Greenwich Mean Time)
Native language: English
Participant ID: 5c69a2847d490900015b0fd3
Timestamp: Sun Feb 17 2019 22:44:14 GMT+0000 (Greenwich Mean Time)
Native language:
Meanings assigned to audio clips:
FULL_MANGLED_SEVEN.wav: NO MEANING
Level_4_SEVEN.wav: NO MEANING
Level_3_SEVEN.wav: NO MEANING
Level_2_SEVEN.wav: White
NOT_MANGLED_SEVEN.wav: Hey Google turn on light
Participant ID: 5c69a2847d490900015b0fd3
Timestamp: Sun Feb 17 2019 22:45:37 GMT+0000 (Greenwich Mean Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5c693a2b58b97400013b3637
Timestamp: Sun Feb 17 2019 14:44:26 GMT-0800 (Pacific Standard Time)
Native language: English
Participant ID: 5c693a2b58b97400013b3637
Timestamp: Sun Feb 17 2019 14:44:34 GMT-0800 (Pacific Standard Time)
Meanings assigned to audio clips:
FULL_MANGLED_SEVEN.wav: NO MEANING
Level_4_SEVEN.wav: NO MEANING
Level_3_SEVEN.wav: NO MEANING
Level_2_SEVEN.wav: Light
NOT_MANGLED_SEVEN.wav: Hey Google, turn on light
Participant ID: 5c693a2b58b97400013b3637
Timestamp: Sun Feb 17 2019 14:45:28 GMT-0800 (Pacific Standard Time)
Study withdrawn: false
Study submitted: true

Participant ID: 542498adfdf99b691fb384d1
Timestamp: Sun Feb 17 2019 14:44:36 GMT-0800 (Pacific Standard Time)
Native language: English
Meanings assigned to audio clips:
Participant ID: 542498adfdf99b691fb384d1
Timestamp: Sun Feb 17 2019 14:44:42 GMT-0800 (Pacific Standard Time)
Meanings assigned to audio clips:
FULL_MANGLED_SEVEN.wav: No meaning
Level_4_SEVEN.wav: no meaning
Level_3_SEVEN.wav: no meaning
Level_2_SEVEN.wav: "turn on light"
NOT_MANGLED_SEVEN.wav: Hey Google, turn on light
Participant ID: 542498adfdf99b691fb384d1
Timestamp: Sun Feb 17 2019 14:45:39 GMT-0800 (Pacific Standard Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5c50829416a93400013a2209
Timestamp: Sun Feb 17 2019 14:48:17 GMT-0800 (Pacific Standard Time)
Native language: English
Participant ID: 5c50829416a93400013a2209
Timestamp: Sun Feb 17 2019 14:48:38 GMT-0800 (Pacific Standard Time)
Meanings assigned to audio clips:
FULL_MANGLED_SEVEN.wav: glistle guide
Level_4_SEVEN.wav: fern charm guide
Level_3_SEVEN.wav: third on shite
Level_2_SEVEN.wav: jay giggle turn on light
NOT_MANGLED_SEVEN.wav: hey google turn on light
Participant ID: 5c50829416a93400013a2209
Timestamp: Sun Feb 17 2019 14:50:35 GMT-0800 (Pacific Standard Time)

Study withdrawn: false
Study submitted: true

Participant ID: 5c6838ea5f8d580001e8a8ab
Timestamp: Sun Feb 17 2019 22:48:48 GMT+0000 (Greenwich Mean Time)
Native language: English
Participant ID: 5c6838ea5f8d580001e8a8ab
Timestamp: Sun Feb 17 2019 22:49:01 GMT+0000 (Greenwich Mean Time)
Meanings assigned to audio clips:
FULL_MANGLED_SEVEN.wav: NO MEANING
Level_4_SEVEN.wav: BROWN GUIDE
Level_3_SEVEN.wav: GOOGLE
Level_2_SEVEN.wav: WHITE
NOT_MANGLED_SEVEN.wav: HEY GOOGLE TURN ON LIGHT
Participant ID: 5c6838ea5f8d580001e8a8ab
Timestamp: Sun Feb 17 2019 22:52:58 GMT+0000 (Greenwich Mean Time)
Study withdrawn: false
Study submitted: true

Participant ID: 59d39d8ba001f200012631b0
Timestamp: Sun Feb 17 2019 22:44:40 GMT+0000 (GMT)
Meanings assigned to audio clips:
FULL_MANGLED_SEVEN.wav: Guide
Level_4_SEVEN.wav: Che google fun guide
Level_3_SEVEN.wav: Che google con shield
Level_2_SEVEN.wav: Jay google turn on light
NOT_MANGLED_SEVEN.wav: Hey google turn on light

Target command: Turn on light
Condition: Google revealed last

Participant ID: cece._17
Timestamp: Sun Feb 17 2019 18:11:53 GMT-0500 (EST)
Native language: English
Participant ID: cece._17
Timestamp: Sun Feb 17 2019 18:12:04 GMT-0500 (EST)
Meanings assigned to audio clips:
FULL_MANGLED_EIGHT.wav: guide
Level_4_EIGHT.wav: glisten turn on
Level_3_EIGHT.wav: on light
Level_2_EIGHT.wav: hey google turn on light
NOT_MANGLED_EIGHT.wav: hey google turn on light
Participant ID: cece._17
Timestamp: Sun Feb 17 2019 18:13:09 GMT-0500 (EST)
Study withdrawn: false
Study submitted: true

Participant ID: 5c69e5ee09041a00010258db
Timestamp: Sun Feb 17 2019 15:11:28 GMT-0800 (Pacific Standard Time)
Participant ID: 5c69e5ee09041a00010258db
Timestamp: Sun Feb 17 2019 15:11:40 GMT-0800 (Pacific Standard Time)
Meanings assigned to audio clips:
FULL_MANGLED_EIGHT.wav: NO MEANING
Level_4_EIGHT.wav: NO MEANING
Level_3_EIGHT.wav: Jay To a function and light
Level_2_EIGHT.wav: Hey good function and light
NOT_MANGLED_EIGHT.wav: hey google turn on light
Participant ID: 5c69e5ee09041a00010258db
Timestamp: Sun Feb 17 2019 15:13:11 GMT-0800 (Pacific Standard Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5c69d300dcdfb700014b137b
Timestamp: Sun Feb 17 2019 23:15:57 GMT+0000 (GMT)
Native language: English
Meanings assigned to audio clips:
Participant ID: 5c69d300dcdfb700014b137b
Timestamp: Sun Feb 17 2019 23:16:11 GMT+0000 (GMT)
Meanings assigned to audio clips:
FULL_MANGLED_EIGHT.wav: No meaning
Level_4_EIGHT.wav: No meaning
Level_3_EIGHT.wav: The last word sounds like light
Level_2_EIGHT.wav: Hey, something something light
NOT_MANGLED_EIGHT.wav: Hey google, turn on light
Participant ID: 5c69d300dcdfb700014b137b
Timestamp: Sun Feb 17 2019 23:17:53 GMT+0000 (GMT)
Study withdrawn: false
Study submitted: true

Participant ID: 5c69df899fcc230001236568
Timestamp: Sun Feb 17 2019 23:12:52 GMT+0000 (Greenwich Mean Time)
Native language: English
Participant ID: 5c69df899fcc230001236568
Timestamp: Sun Feb 17 2019 23:13:03 GMT+0000 (Greenwich Mean Time)
Meanings assigned to audio clips:
FULL_MANGLED_EIGHT.wav: NO MEANING

Level_4_EIGHT.wav: shite
Level_3_EIGHT.wav: NO MEANING
Level_2_EIGHT.wav: GOODMORN LIGHT
NOT_MANGLED_EIGHT.wav: hey google, turn on light
Participant ID: 5c69df899fcc230001236568
Timestamp: Sun Feb 17 2019 23:14:44 GMT+0000 (Greenwich Mean Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5c69d992f00261000194021a
Timestamp: Sun Feb 17 2019 18:13:11 GMT-0500 (Eastern Standard Time)
Native language: English
Participant ID: 5c69d992f00261000194021a
Timestamp: Sun Feb 17 2019 18:13:18 GMT-0500 (Eastern Standard Time)
Meanings assigned to audio clips:
FULL_MANGLED_EIGHT.wav: words heard: whistle, guide
Level_4_EIGHT.wav: whistle
Level_3_EIGHT.wav: light
Level_2_EIGHT.wav: Hey Google, turn on light.
NOT_MANGLED_EIGHT.wav: Hey Google, turn on light.
Participant ID: 5c69d992f00261000194021a
Timestamp: Sun Feb 17 2019 18:14:44 GMT-0500 (Eastern Standard Time)
Study withdrawn: false
Study submitted: true

Timestamp: Sun Feb 17 2019 23:11:12 GMT+0000 (Greenwich Mean Time)
Meanings assigned to audio clips:
FULL_MANGLED_EIGHT.wav: No meaning
Level_4_EIGHT.wav: No meaning
Level_3_EIGHT.wav: Turn on light
Level_2_EIGHT.wav: Hey something light
NOT_MANGLED_EIGHT.wav: Hey google turn on light
Timestamp: Sun Feb 17 2019 23:14:06 GMT+0000 (Greenwich Mean Time)
Study withdrawn: false
Study submitted: true

Target command: Who am I
Condition: Google revealed first

Participant ID: 574802c17fd0ec000eb63b30
Timestamp: Sun Feb 17 2019 18:46:05 GMT-0500 (Eastern Standard Time)
Native language: English
Participant ID: 574802c17fd0ec000eb63b30
Timestamp: Sun Feb 17 2019 18:46:14 GMT-0500 (Eastern Standard Time)
Meanings assigned to audio clips:
FULL_MANGLED_NINE.wav: Google Fu Dan
Level_4_NINE.wav: Giggle food
Level_3_NINE.wav: Google
Level_2_NINE.wav: Google who am i?
NOT_MANGLED_NINE.wav: Hey Google, who am i?
Participant ID: 574802c17fd0ec000eb63b30
Timestamp: Sun Feb 17 2019 18:49:01 GMT-0500 (Eastern Standard Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5c685aa09268b70001a330ba
Timestamp: Sun Feb 17 2019 23:46:16 GMT+0000 (Greenwich Mean Time)
Native language: English
Participant ID: 5c685aa09268b70001a330ba
Timestamp: Sun Feb 17 2019 23:46:22 GMT+0000 (Greenwich Mean Time)
Meanings assigned to audio clips:
FULL_MANGLED_NINE.wav: NO MEANING
Level_4_NINE.wav: NO MEANING
Level_3_NINE.wav: NO MEANING
Level_2_NINE.wav: NO MEANING
NOT_MANGLED_NINE.wav: Hey google, who am I ?
Participant ID: 5c685aa09268b70001a330ba
Timestamp: Sun Feb 17 2019 23:47:13 GMT+0000 (Greenwich Mean Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5c69e6514b5bb30001b3d29d
Timestamp: Sun Feb 17 2019 17:47:25 GMT-0600 (Central Standard Time)
Native language: English
Participant ID: 5c69e6514b5bb30001b3d29d
Timestamp: Sun Feb 17 2019 17:47:35 GMT-0600 (Central Standard Time)
Meanings assigned to audio clips:
FULL_MANGLED_NINE.wav: Fish eye
Level_4_NINE.wav: A giggle
Level_3_NINE.wav: Who am I
Level_2_NINE.wav: A giggle who am cry
NOT_MANGLED_NINE.wav: Hey google who am I
Participant ID: 5c69e6514b5bb30001b3d29d
Timestamp: Sun Feb 17 2019 17:49:03 GMT-0600 (Central Standard Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5c687b0c1b9d0e000190aae2
Timestamp: Sun Feb 17 2019 15:50:29 GMT-0800 (Pacific Standard Time)
Native language: English
Participant ID: 5c687b0c1b9d0e000190aae2
Timestamp: Sun Feb 17 2019 15:50:37 GMT-0800 (Pacific Standard Time)
Meanings assigned to audio clips:
FULL_MANGLED_NINE.wav: NO MEANING
Level_4_NINE.wav: NO MEANING
Level_3_NINE.wav: NO MEANING
Level_2_NINE.wav: NO MEANING
NOT_MANGLED_NINE.wav: Hey Google, who am I.
Participant ID: 5c687b0c1b9d0e000190aae2
Timestamp: Sun Feb 17 2019 15:51:37 GMT-0800 (Pacific Standard Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5c68a07e67057200011ace92
Timestamp: Sun Feb 17 2019 23:50:51 GMT+0000 (Greenwich Mean Time)
Native language: English
Participant ID: 5c68a07e67057200011ace92
Timestamp: Sun Feb 17 2019 23:51:10 GMT+0000 (Greenwich Mean Time)
Meanings assigned to audio clips:
FULL_MANGLED_NINE.wav: No meaning
Level_4_NINE.wav: Google
Level_3_NINE.wav: Google who am i
Level_2_NINE.wav: A Google who am rye
NOT_MANGLED_NINE.wav: Hey Google who am i
Participant ID: 5c68a07e67057200011ace92
Timestamp: Sun Feb 17 2019 23:53:28 GMT+0000 (Greenwich Mean Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5c680ad1b1b0080001bfc04f
Timestamp: Sun Feb 17 2019 23:46:00 GMT+0000 (GMT)
Native language: English
Participant ID: 5c680ad1b1b0080001bfc04f
Timestamp: Sun Feb 17 2019 23:46:11 GMT+0000 (GMT)
Meanings assigned to audio clips:
FULL_MANGLED_NINE.wav: No meaning
Level_4_NINE.wav: No meaning
Level_3_NINE.wav: No meaning
Level_2_NINE.wav: Google who am I
NOT_MANGLED_NINE.wav: Hey google who am I

Target command: Who am I
Condition: Google revealed last

Participant ID: 5c44e3b73be7b70001fd792a
Timestamp: Sun Feb 17 2019 19:18:36 GMT-0500 (Eastern Standard Time)
Native language: English
Participant ID: 5c44e3b73be7b70001fd792a
Timestamp: Sun Feb 17 2019 19:18:41 GMT-0500 (Eastern Standard Time)
Meanings assigned to audio clips:
FULL_MANGLED_TEN.wav: No meaning
Level_4_TEN.wav: No meaning
Level_3_TEN.wav: Hey
Level_2_TEN.wav: Hey, who am I
NOT_MANGLED_TEN.wav: Hey Google, who am I?
Participant ID: 5c44e3b73be7b70001fd792a
Timestamp: Sun Feb 17 2019 19:19:47 GMT-0500 (Eastern Standard Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5c6981bb0894a000015ad0a1
Timestamp: Mon Feb 18 2019 00:19:19 GMT+0000 (Greenwich Mean Time)
Native language: English
Participant ID: 5c6981bb0894a000015ad0a1
Timestamp: Mon Feb 18 2019 00:19:26 GMT+0000 (Greenwich Mean Time)
Meanings assigned to audio clips:
FULL_MANGLED_TEN.wav: NO MEANING
Level_4_TEN.wav: NO MEANING
Level_3_TEN.wav: NO MEANING
Level_2_TEN.wav: something something who am something
NOT_MANGLED_TEN.wav: Hey Google, who am I?
Participant ID: 5c6981bb0894a000015ad0a1
Timestamp: Mon Feb 18 2019 00:21:03 GMT+0000 (Greenwich Mean Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5c68af707d17f800013f6ad9
Timestamp: Sun Feb 17 2019 16:23:24 GMT-0800 (Pacific Standard Time)
Native language: English
Participant ID: 5c68af707d17f800013f6ad9
Timestamp: Sun Feb 17 2019 16:23:40 GMT-0800 (Pacific Standard Time)
Meanings assigned to audio clips:
FULL_MANGLED_TEN.wav: No meaning
Level_4_TEN.wav: Google
Level_3_TEN.wav: No meaning
Level_2_TEN.wav: Hey - no meaning

NOT_MANGLED_TEN.wav: Hey Google who am I
Participant ID: 5c68af707d17f800013f6ad9
Timestamp: Sun Feb 17 2019 16:25:49 GMT-0800 (Pacific Standard Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5c6363e6c19d950001efb0d4
Timestamp: Mon Feb 18 2019 11:23:49 GMT+1100 (Australian Eastern Daylight Time)
Native language: English
Participant ID: 5c6363e6c19d950001efb0d4
Timestamp: Mon Feb 18 2019 11:24:01 GMT+1100 (Australian Eastern Daylight Time)
Meanings assigned to audio clips:
FULL_MANGLED_TEN.wav: NO MEANING
Level_4_TEN.wav: NO MEANING
Level_3_TEN.wav: "Hey", then a bunch of gibberish
Level_2_TEN.wav: Hey look I am grey
NOT_MANGLED_TEN.wav: Hey Google, who am I?
Participant ID: 5c6363e6c19d950001efb0d4
Timestamp: Mon Feb 18 2019 11:25:16 GMT+1100 (Australian Eastern Daylight Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5c69d1df5f8d580001e8b339
Timestamp: Mon Feb 18 2019 00:24:31 GMT+0000 (Greenwich Mean Time)
Native language: English
Participant ID: 5c69d1df5f8d580001e8b339
Timestamp: Mon Feb 18 2019 00:24:46 GMT+0000 (Greenwich Mean Time)
Meanings assigned to audio clips:
FULL_MANGLED_TEN.wav: No meaning
Level_4_TEN.wav: No meaning
Level_3_TEN.wav: No meaning
Level_2_TEN.wav: No meaning
NOT_MANGLED_TEN.wav: Hey Google, who am I?
Participant ID: 5c69d1df5f8d580001e8b339
Timestamp: Mon Feb 18 2019 00:26:18 GMT+0000 (Greenwich Mean Time)
Study withdrawn: false
Study submitted: true

Participant ID: 5c683cd8c643f400012ea0dc
Timestamp: Sun Feb 17 2019 17:26:39 GMT-0700 (MST)
Meanings assigned to audio clips:
FULL_MANGLED_TEN.wav: no meaning
Level_4_TEN.wav: A giggle
Level_3_TEN.wav: Hey look at
Level_2_TEN.wav: Hey look at who am
NOT_MANGLED_TEN.wav: Hey google who am i

Appendix B

Missense Attacks on Amazon Alexa

B.1 Pilot Experiment - Training Dataset

```
## intent:bye
- Bye
- Goodbye
- See you later
- Bye bot
- Goodbye friend
- bye
- bye for now
- catch you later
- gotta go
- See you
- goodnight
- have a nice day
- i'm off
- see you later alligator
- we'll speak soon

## intent:greet
- Hi
- Hey
- Hi bot
- Hey bot
- Hello
- Good morning
- hi again
- hi folks
- hi Mister
- hi pal!
- hi there
- greetings
- hello everybody
- hello is anybody there
- hello robot

## intent:thank
- Thanks
- Thank you
- Thank you so much
- Thanks bot
- Thanks for that
- cheers
- cheers bro
- ok thanks!
- perfect thank you
- thanks a bunch for everything
- thanks for the help
- thanks a lot
- amazing, thanks
- cool, thanks
- cool thank you

## intent:affirm
- yes
- yes sure
- absolutely
- for sure
- yes yes yes
- definitely
```

```

## intent:name
- My name is [Alice](name)
- I am [Josh](name)
- I'm [Lucy](name)
- People call me [Greg](name)
- It's [David](name)
- Usually people call me [Amy](name)
- My name is [John](name)
- You can call me [Sam](name)
- Please call me [Linda](name)
- Name name is [Tom](name)
- I am [Richard](name)
- I'm [Tracy](name)
- Call me [Sally](name)
- I am [Philipp](name)
- I am [Charlie](name)

## intent:accountBalance
- balance in my account
- show the balance on all my accounts
- show the balance on my savings account
- want to check my account balance
- check my current account balance
- what is my savings account balance
- what is my account balance
- what is the balance on my account
- how much money do I have in my savings account
- how much money do I have
- how much do I have in savings
- what is my balance
- tell me my balance
- show me my account balance
- how much do I have in my current account
- tell me my current balance
- what is the current balance on my account
- how much is there is my account
- what is my savings balance
- how much do I have in my savings account
- how much is left in my savings account
- tell me how much I have in my current account
- how much do I have in all my accounts
- want to check my savings account balance
- how much money do I have left
- check how much money I have
- i want to know how much I have in each account
- tell me my savings account balance
- what is my current balance
- how much money is left in my current account

## intent:viewTransactions
- can i have a look at my transactions
- list down my recent transactions
- provide me an account statement
- show my transactions
- view my transactions
- what are the activities on my account
- what were my expenses
- what were my transactions
- where did I spend my money
- how much have I spent this month
- how much did I spend last month
- show me my recent transactions
- what were my transactions last week
- what were my outgoings this month
- what was paid into my account last month
- what were my outgoings last month
- what has been paid into my account this month
- read my account statement
- read my statement
- recent transactions
- how much have I spent this week
- what did I spend last week
- what were the payments into my account last week
- this week's outgoings
- my transactions last week
- transactions last month
- transactions so far this month
- tell me my transactions
- how much have I spent recently
- show my statement

## intent:payBill
- make a payment
- make bill payment
- pay my bills
- pay all my outstanding bills
- pay my bill for gas
- pay my credit card bill
- pay my utilities bills
- pay my mobile phone bill

```

- pay electricity bill
- pay tax bill
- pay outstanding bills
- pay my debts
- pay off my debts
- pay rent bill
- pay my rent due
- clear my outstanding bills
- clear my debts
- pay my water bill
- pay off my credit card debt
- pay my bill for water
- pay my bill for tax
- make a payment for electricity
- clear my water bill
- clear my credit card balance
- pay the electricity bill
- make a rent payment
- pay all bills due
- make a payment to my landlord
- pay all bills due this month
- pay this month's mobile phone bill

intent:forgotPassword

- can you help with password reset
- how can i reset my password
- i don't remember my password
- i forgot my password
- i need to reset password
- i need to change my password
- i need to set a new password
- have forgotten my password
- can't remember my password
- don't know my password
- i need a new password
- can you reset my password
- please reset my password
- change password
- new password

intent:lostCard

- card is lost
- could you block my card
- freeze my card
- loss of my card
- lost my card
- my card is stolen
- need to block my card
- i have lost my card
- can't find my card
- please block my card
- someone stole my card
- need to cancel my card
- cancel card
- have mislaid my card
- don't have my card

intent:nonsense

- 1 mm equals how many my what
- A body part that starts with G
- A cell with abundant peroxisomes would most likely be involved in what
- A general term for any disease causing organism is
- A localized group of organisms that belong to the same species is called a what
- A means of bringing fairly large particles into the cell
- A Piedmont is a physical region found where
- A process that requires the addition of energy is
- A series of events in which one organism eats one another
- About one quart 0.908 dry quart is equivalent to what
- According to Darwin Natural Selection is based on what found in populations
- After how many miles Must Change timing belt
- Air resistance is a type of friction
- An antimicrobial that inhibits hepatic P450 microsomal enzymes
- An object slowing down has a What acceleration
- An organism living in or on another person
- At the end of World War 2 the Allies agreed on the creation of what organization to settle future disputes
- At what wave height did the pile installation at NRB need to cease
- Being a lightweight with a high melting point this group II metal is an ideal hardening agent
- Blood leaving the right ventricle enters the what
- Demographic factors such as race ethnicity or region help define political opinion in a process known as political
- Diabetes affects which endocrine gland
- Diaphragmatic hernia is also known as
- Did William R. Thomas' base salary increase during 2013
- Distance from one point to another
- Does a centipede reproduce sexually or asexually
- In 1976 which gymnast scored maximum scores of 10 as she won three gold medals
- In badminton what is always an underhand shot
- In many states what allows the people to vote on laws passed by the state Legislature and signed by the governor

- In Texas how close to a fireplug may a vehicle park
- In the book charletts web What is the rats name what
- Incubation period of measles
- Insulin and glucagons regulate blood levels of what
- Main function of a leaf
- Malaria and amoebic dysentery are caused by what
- Mineral substance found in blood
- Most eccentric orbit
- Most federal circuits of the federal appellate court system are determined based on what
- Most muscular and elastic blood vessel
- Most surface ocean current are due to what
- Most truck trailers are feet wide
- Name for CAO
- Nudibranchs are members of which group of molluscs
- Null this out
- you like pasta
- you make a suggestion
- Your tongue is covered with a layers of bumps that contain taste buds what are they called
- about
- after
- a higher level of annoyance than norm
- Ali rampage
- all eligible leann
- all I could
- all systems
- all that stuff we're doing I'm trying to figure out what are
- all the
- all this darkness
- all those out
- along
- a loss books
- a lot of artwork
- a lot of the times was
- alright so chef Watson is this
- alright so piazza doesn't like that
- alright so that's kind of our robots where we're at with this on
- alright teaches one everything oughtn't
- also
- also make it date thinking
- also see
- although
- an athlete whether
- and
- and basically
- and I don't all park
- and I have and that's a good example
- and intimate writes about what that means now says
- and I think one of
- and just
- and our
- and provide
- and so that
- and then you say you know
- and this
- and those three right now variables
- and you couldn't find service so is having a hard time
- anyone else
- anyone the time
- any story that
- anything about
- a participant
- are
- area
- are just so number I was talking about
- are on the move stay like this
- are six
- are you
- as for objects virtual room rewards
- as the
- a tiny or
- audience station
- author's note this
- a verbal menu
- back out to all the I don't know our time
- because this is automatically
- besides
- better proxy I don't know of points of discussion right now but later this may not be I think
- bill Dorman said by dividing the total
- block
- brain I business is the stuff also you know considering
- but for the auditor scenario
- but I feel
- but if it only fills itself after
- but I thought it would be just another demos
- but it's
- but it's not possibly trying applauded me
- but it's this
- but I would like to
- but let's get this wrong and then
- but not always
- but now I know the workings last night
- but on hands

- but on the
- but on the other side
- but on your leann
- but that's not on a
- but that's what's
- but there's a very fine job people say out of the
- but we chucks so you get a chance before you have a lawyer
- but when you can use something
- but when you go up to come to here they only know so many things would you know if this is
- but you have
- but you have a verbal menu which is a failed
- but you know it's
- can
- can get
- can go
- can I ask
- can I knew everybody from the conference lines or see
- can questions later
- can you know
- cap'n
- check out some wires
- could be a particular after it will have suggested place should you make it have shorter memory
- could be done is right maybe it's
- delays
- developers
- does
- does seasonal
- does that
- does the answer
- do I deserve the result of lots of laps
- don't ask about some
- do then why don't
- do they do they
- do things talk to
- down
- do you know
- do you still like bird where's the anything like perfect
- do you support Richard though
- do you want
- dragging
- earlier this and attacks over somebody look you know what that you asked for privacy
- either like
- either time
- entity relationship you know classify those out
- even privacy where
- everybody stop
- every now and it'll suck randomly talk
- every statements
- except
- except that it will be
- existing home
- famous novelist
- faster
- finding
- first it was to get it perfect so Hussein for the purpose of this debate that you give a shot
- first thing
- first you would as opposed
- five
- for
- for the most part Elsie's perfect's
- friends
- functional honest
- general
- girl please
- giving
- going to build a whole infrastructure
- great
- he lifted out it's good stuff is
- here
- her personal
- he's all like this
- he's not including
- he those two questions lifted its
- Hong
- how
- how did that
- how to get everything
- how will that is like status of the book at all
- how would supersede
- I am around I am always happy when I am thinking again
- I asked what
- I but I
- I can barely
- I can get do
- I'd
- I decided that is
- I didn't ask you that
- I'd just
- I'd like there was one did suggest
- I don't does it look no
- I don't have you seen parties should be something
- I don't know about everyone vannatter

- I don't know which also burns like clearly how can I talk about
- I don't really good
- I don't think for something
- I drive by a fly ball go whatever the
- I dug out years ago I
- I feel like you go about saying
- if I am really hoping I can get five
- if I can only
- I figured this is my
- if I think about substance
- if I were to go maybe account
- if they don't
- if you
- if you can do it a process they or eliminated helps with
- if you don't like off
- if you say things you know related to weather like relies on China rain
- if you schedule work at on is
- if you want notifying on
- I just solutions
- I just will follow look at the list recently
- I know if you lose talk
- I know the answer
- I know your intellect
- I mean ideally be nice to have arranged her
- I mean if I might not do that you know but
- I mean something
- I'm going
- I'm kind of a sense of all the penalties are like no
- I'm looking for
- I'm sorry
- I'm trying to figure out
- I'm with words
- in Europe results back maybe have a little gesture like
- in the lobby there anything
- I often
- I stood
- is whenever things go awry
- it
- it amounts to suspects
- it helps us quite a bit
- I think at one point I never doubted that he was going to make it
- I think erupted
- I think it's annoying as well
- I think that
- I think that you know like
- I think this is it only
- I thought it contest
- I thought there
- I thought we had maybe like so
- it in the timeouts
- it is
- it is abilities
- it is useless
- it'll
- it'll ignored
- it'll make it you will local stuff
- it'll people does not generally you is or Starbucks and signaled to the robot robot
- it'll take a key phrases like love life realized that all certain
- it oughtn't
- it out a secondary and Alcee retired nonsense
- it's
- it's a
- it's about the language Nassau executives see the system
- it's actually
- it's a different style for different phones
- it's going dialogues
- it's good
- it's just
- it's just when you
- it's kind of
- it's like
- it's my considered talk to see some
- it's not efficient
- it's not like prominence of the
- it's probably somewhere entirely
- it's very everything on saying
- it's Yankees will work on the image processing for this free butter pecan also museum grabbing
- it was
- it was a really
- it was easy
- it would take
- it your fixed
- I would like to
- I your real life experience that would you know in
- Jude auto is has the hams
- just
- just because it was suffocating
- just because you
- just leann
- just like just
- just sort of a Protestant trying to figure out
- last

- leann
- learn from center learn from
- let's counselors to get lost lost
- let's get
- like cases I was little testing for a highly more realistically by
- like Selassie's question
- like some
- like the ones
- like will say something my honorable log
- look ani
- lots
- make sure that it
- maybe what
- model
- Moses
- most
- mostly customers
- moving
- my
- my answer machine not much of the
- my problem for Albanian there unless the inexorably
- NATO
- never picks it up
- next
- next year
- nobody
- nobody asked whether
- nobody wants
- not that this
- not whenever we train
- now the first thing that I want
- okay also
- okay listing
- on
- once
- once we started
- once you have it
- one of them aloneness
- one of these guys
- only pick the probation podiums
- on this
- on this particular do not
- ordeal is about a
- orthodontist of the robot
- our partners
- Palestinian for
- parents
- perhaps most
- plan goes how I know I know how to handle
- please note
- points responses to
- politics
- pony
- pretty fast
- probably down forcibly six hour
- problems
- Providence
- rating
- ray made the point
- reasonably exorcism is definitions are profits were
- remiss usually municipal goes like this screw so little Jews rusher includes several we don't have this issue
- request
- responded whether
- right
- right it
- right now I don't know what's going
- rob is the guy that center for the better the
- robots
- say
- say I love you can speak normally but not like
- says exemplifies re application
- screen
- Scully the problems
- Seattle problem noise assignments foreplay
- second dialogue
- second one was
- see
- seems happy
- select
- sh
- she conte
- she had printed happen
- she's the only one
- shortest
- show so I know he showed he knows he show I've never thought has
- since you know
- small seems I want potatoes
- smiles tested all stop
- snarling inside
- so
- so all I have to do is give it some nonsensical utterances
- so anybody who has any

- so because you are getting satisfactory answers a good purpose of going to say that I'm going to
- so beside the road eaters detectives over the next also biologic excellent tomorrow
- so can I get more of last
- so does well
- so for the immediate time
- so goes last night I
- so go free
- so I clearly am
- so I does this vessel here
- so I felt that we sit
- so I still
- so I think that situation I understand
- so I think the
- so it is
- so it's like this
- so it's uses a concierge so they put it out there plans on and
- so let's
- so like whatever I tell it to raise an army that's
- somebody
- some boys giving a snickers leann
- some items are good
- some of the
- some of the bugs I know it's not the most exciting glorious
- someplace in place you know
- something
- something else here suppose you receive a penalty will be able to answer
- something is
- something out of it
- something that affects into is like
- something that we have
- something to make it more experience
- so next I think we need to fix for ninety
- son of
- so not sure if you see it
- so now you have
- so once again actively is focus on bugs certain things out
- so or something
- so policies like figure out
- sorry
- sorry fifty
- sorry for
- sort of
- sort of thing have Yassin's
- so stuff they've been
- so that
- so that is what I
- so that lousy
- so that's
- so that's a question you
- so that's that's what I'm trying
- so that you have a little bit snotty
- so the idea being by life
- so there is
- so there's there's so's missing
- so these are some questions
- so the second part of
- so the simple
- so this
- so this goes down
- so this thing has a built in system were prototype you know you can
- so was there something that includes
- so we didn't also so artistic
- so we just it was a matching set of words that help Thursday
- so we're not getting
- so whatever box
- so what if they simply applied for what accepted my status
- so what I want to do with a little less
- so what's TGI
- so what we're trying to do is always listening
- so what we try not only that you don't know the Serbs
- so while
- so with our language pacifier's if we use that goes a long way towards
- so you can
- so you know is using the local trains go
- so your aloneness
- so your treasury
- so you say
- speaker out
- split between good about
- stacking
- Stansky
- start
- starting
- states cap'n
- stay
- Steve transcribe the ideas that when you get in these
- still
- still oughtn't shall
- still research
- still we need to know what's the what's stored locally
- strong hands everyone
- student Maggie

- stuff
- submit
- such
- supposing actually
- symbolizing
- talked over me a lot
- technical
- testing
- Texas
- text
- that
- that are
- that fastest we did before so the capital code
- that I
- that is
- that is not a bug consultants are chemically it's it's not necessarily hot that sentence but
- that means
- that's
- that's a different
- that said lots and does
- that's an efficiency divisions
- that's fine
- that's interesting
- that's like a real
- that's that's saying a lot so we're going to be
- that's where
- that's within this
- that we have a
- the
- the bagel place that is the most was that is
- the city
- the city does
- the closer you see the words are area don't want
- the developers not
- the developers they'll have stickers on the back
- the fifth open
- the help they just barely plaintiffs kicked off all data visible on the one
- the law here
- the liberals
- then
- then it came
- the number to do just that
- the one
- the one thing your friends
- the other part is one problem though
- the population
- there are two questions depends on the call including you
- the reason
- the reasons
- there is
- there is always you know these
- there is an infection in their survivor
- there is some
- there should
- the responsible
- the robot
- the so
- the update on its books but I would like to request is whenever you are logging roads
- they
- they are as a person
- they are working on that there should be some testing
- they do you
- they look cool our wisdom
- they others of us service
- they're very similar that's
- they see robots go when they don't
- think about annexes
- thirty answering
- this
- this accident rollerblade's lawyers in your business avoiding you go back basically lose lose
- this area
- this four seconds
- this is
- this is a problem that we Casey
- this is good
- this means that as deeply
- this morning I usually just
- this obligated scatter within communities
- this results
- this robust Lebanese
- those those rules are
- three dimensional space of the room is from
- three tasks unclear that I think
- through the years
- total whispers
- trying to figure space hotels we're trying to help you because
- trying to make
- try to make bus
- unfortunately anybody's
- unless
- until it does the dishes
- user

- vannatter
- verses
- Wagner said each interaction
- waiting for them sternly needs to maybe this is something that's like a promotion
- we all
- we can do I train in for that search just being able
- we dishes
- we fought so I'd like
- we have a
- well
- well I said earlier
- well it looks like to network
- well itself architecture that I kind of worried
- well it was
- well it was a
- well it was just a signature
- well it wouldn't
- well maybe
- well please list if you after that last request
- well so probably pick up umbrella in in
- well there's a sweeping point they will never die
- well what we do about
- we may at this last night
- we're
- we're all people you actually just because we
- we're going to the robot
- we're looking for
- we're probably doing
- we're working I was working with and like
- we see
- we started
- we talking about here
- we train system lately
- we will
- we work with
- with each record
- work on this to get everything merged
- works
- would you go to a hotel because here in the Hague you know if I don't have anything
- yeah like I don't
- yeah through only one or stay
- yeah total whispers
- you
- you can
- you can also decided
- you can get it out
- you can study time but
- you can use this
- you can you have a historical
- you feel totally focused upon receipt elevator
- you have several
- you have to say all screens
- you just available
- you know a big
- you know I mean is we have sort of what should
- you know it's like where did
- you know this evening
- you'll see
- you pass
- you're going
- you've got
- you will be
- you will be sorry
- you would like to see
- zero amount
- zero zero zero over thirty I guess

B.2 Experiment - Unsuccessful Adversarial Utterances

Target Easter Egg Command	Adversarial Utterance	Alexa Target System Response
flip a coin	flipacoin.com	non-target response (suggests 'Flip a Coin' Skill)
tell me a joke	telltale games	non-comprehension
tell me a spooky story	the scene from the first storey is spooky	non-comprehension
what's your favourite hobby?	this hobby is my favourite	non-comprehension
do you like cats or dogs?	my dog acts like cats	non-comprehension
the first rule of Fight Club	for the first time, there was a fight in the club	non-comprehension
are you Sky Net?	the sky in the limit	non-target response (suggests music track 'the sky is the limit')
party time!	this party's time	non-comprehension
open the pod bay door	an open door policy	non-target response (information on open door policy)
how much wood can a woodchuck chuck if a woodchuck could chuck wood?	I bought a chuck of wood	non-comprehension
which comes first: the chicken or the egg?	they egged him first	non-comprehension
may the Force be with you	I want to see the force with you	non-comprehension
who let the dogs out?	I want to let my home no dogs allowed	no response
where are my keys?	this is my key contribution	non-comprehension
rock paper scissors	rock the box that contains paper and scissors	non-target response (suggests 'rock, paper, scissors, lizard, Spock' Skill)
is the cake a lie?	I lied in front of the cake	no response
what is the sound of one hand clapping?	this sounds like a hand clapping	non-comprehension

Table B.1: Adversarial utterances results - participant one

Target Easter Egg Command	Adversarial Utterance	Alexa Target System Response
tell me a joke	his tell is a joke	non-target response (information on YouTuber A Joke)
tell me a spooky story	it's spooky on that story	non-target response (suggests 'spooky sounds' Skill)
what's your favourite hobby?	favourite hobby horse	no response
do you like cats or dogs?	cool cats like hot dogs	non-target response (suggests 'bark like a dog' Skill)
party time!	won't make the party in time	non-target response (suggests 'animal sounds' Skill)
which comes first: the chicken or the egg?	first ask but if they are chicken, egg them all	no response
rock paper scissors	could you cut rock paper with scissors?	non-comprehension
is the cake a lie?	just let the cake lie	non-comprehension
what is the sound of one hand clapping?	Sound! I'll give you a hand in clapping	non-comprehension
can you fly?	Keanu is pretty fly	non-target response (answers 'no, Keanu is not a fly')

Table B.2: Adversarial utterances results - participant two

Target Easter Egg Command	Adversarial Utterance	Alexa Target System Response
flip a coin	to coin a flip	non-target response (information on how to flip a coin)
tell me a joke	this so-called tell is a joke	non-comprehension
what's your favourite hobby?	this hobby is my favourite hawk	non-comprehension
do you like cats or dogs?	raining like cats and dogs	non-target response (explanation of the expression)
the first rule of Fight Club	first, the club fights against the rule	non-comprehension
are you Sky Net?	Is Sky net or gross?	non-comprehension
open the pod bay door	open the door for the bay pod	non-comprehension
which comes first: the chicken or the egg?	egg the chicken to go first	no response
where do you live?	where is the live wire?	non-target response (information on film 'Live Wire')
who let the dogs out?	he's been dogged by lets?	no response
is the cake a lie?	where does the cake lie?	non-target response (information on the history of cake)
what is the sound of one hand clapping?	if she can't clap without her hands I don't want her sound	non-comprehension
can you fly?	are you a fly?	non-comprehension
is this the real life?	my life is not worth a real	non-comprehension

Table B.3: Adversarial utterances results - participant three

Target Easter Egg Command	Adversarial Utterance	Alexa Target System Response
flip a coin	who coined the phrase flip a bird	non-target response (suggests 'Flip a Coin' Skill)
tell me a joke	what's Penn and Teller's best joke?	non-comprehension
tell me a spooky story	tell me how many stairs are in one story	non-target response (answers 'stories and steps are not compatible units')
what's your favourite hobby?	cycling is my favourite hobby	non-comprehension
do you like cats or dogs?	do cats like dogs?	non-target response (information on miscommunication between cats and dogs for pet owners)
the first rule of Fight Club	when did Mohammed Ali have his first fight in a boxing club?	non-comprehension
are you Sky Net?	how fast is the Sky broadband network?	non-target response (information on Chicago Sky basketball games)
party time!	is there a Party in Times Square?	non-target response (answers 'The Times is affiliated with centre right politics')
open the pod bay door	do dogs bay if you don't open the door?	non-comprehension
nice to see you, to see you...	is Nice next to the sea?	non-target response (answers 'there aren't any seas near Nice')
how much wood can a woodchuck chuck if a woodchuck could chuck wood?	how much more wood could Chuck Norris chuck than a woodchuck?	non-target response (answers 'Chuck Norris would chuck the woodchuck, not the other way around')
which comes first: the chicken or the egg?	when do chickens lay their first egg?	non-comprehension
may the Force be with you	what's the SI unit for force?	non-target response (information on SI unit)
where do you live?	are you alive?	non-target response (answers 'artificially, maybe, but not in the same way that you are alive')
who let the dogs out?	should I let my dog on the sofa?	non-target response (story about dog rescued from the sea)
where are my keys?	what is key lime pie?	non-target response (information on key lime pie)
rock paper scissors	what wins out of rock, paper and scissors?	non-comprehension
what is the sound of one hand clapping?	why is the sound of clapping your hands so loud?	non-comprehension
can you fly?	do chicken fly?	non-target response (answers 'chicken can fly for a short distance')
is this the real life?	is Second Life real?	non-target response (information on virtual world Second Life)

Table B.4: Adversarial utterances results - participant four

Target Easter Egg Command	Adversarial Utterance	Alexa Target System Response
flip a coin	drop a coin in the flip	non-target response (suggests 'Flip a Coin' Skill)
tell me a spooky story	don't tell me it's spooky on this story	non-target response (suggests audiobook)
what's your favourite hobby?	what's your hobby's favourite?	non-target response (information on John Adam's favourite hobby)
are you Sky Net?	did you sky the net?	non-comprehension
party time!	I was party to it this time	no response
which comes first: the chicken or the egg?	first let's egg the chicken on	no response
where do you live?	do you live?	non-target response (answers 'I am not really alive, but I can be lively')
is the cake a lie?	let the cake lie	non-target response (answers 'cakes are not lies')
what is the sound of one hand clapping?	is the hand that's clapping sound	non-target response (information on book 'The Sound Of One Hand Clapping')
can you fly?	can you be fly?	non-comprehension

Table B.5: Adversarial utterances results - participant five

Target Easter Egg Command	Adversarial Utterance	Alexa Target System Response
flip a coin	drink that flip	non-target response (suggests 'Flip a Coin' Skill)
tell me a joke	the joker doesn't tell on others	no response
tell me a spooky story	what a spooky three story building	non-comprehension
what's your favourite hobby?	ride your favourite hobby	non-comprehension
do you like cats or dogs?	it's raining cats and dogs	non-target response (information on book 'It's Raining Cats and Dogs')
party time!	our party won the election	non-target response (information on current UK election)
which comes first: the chicken or the egg?	the egg came from that chicken	non-target response (joke about mis-laying eggs)
where do you live?	live your life to the fullest	non-target response (advice on living a full life)
where are my keys?	I own that key	non-comprehension
rock paper scissors	rock the baby	non-target response (suggests song 'Rock My Baby')
is the cake a lie?	he lies like it is a piece of cake	no response
can you fly?	catch that fly	non-target response (information on how to catch flies)
is this the real life?	life is priceless, can't be measured in real	non-comprehension

Table B.6: Adversarial utterances results - participant six

Target Easter Egg Command	Adversarial Utterance	Alexa Target System Response
flip a coin	how much coin does a flip cost	non-target response (answers 'Sorry, I don't know how much lateral flagellar export/assembly protein ECUMN 0254 is worth')
tell me a joke	tell off for telling a joke	no response
what's your favourite hobby?	the hobby horse was my favourite toy	non-comprehension
do you like cats or dogs?	do cats look like dogs?	non-target response (suggests 'bark like a dog' Skill)
the first rule of Fight Club	what are the club rules on a first class flight?	non-comprehension
nice to see you, to see you...	mice can see nice mice	non-comprehension
may the Force be with you	the force sprayed on you	non-comprehension
where do you live?	how do you live like that?	non-target response (advice on how to live a happy life)
who let the dogs out?	should I let the dogs out without a lead?	non-comprehension
where are my keys?	I went to Florida but couldn't find any keys	no response
rock paper scissors	I rocked the paper boat with scissors	non-comprehension
is the cake a lie?	can you lie on a cake bed?	non-comprehension
is this the real life?	is he a real idiot at life?	refusal to answer ('I'd rather not answer that')

Table B.7: Adversarial utterances results - participant seven

Target Easter Egg Command	Adversarial Utterance	Alexa Target System Response
flip a coin	coin a flippancy	non-target response (suggests 'Flip a Coin' Skill)
tell me a joke	a telling joke	non-target response (answers 'I don't have a joke about that')
do you like cats or dogs?	cat-like dog	non-target response (information on cats and dogs living together)
are you Sky Net?	net that there sky	no response
party time!	Time Party	non-target response (information on political parties)
nice to see you, to see you...	Let's go see Nice	non-comprehension
may the Force be with you	don't force it!	non-comprehension
where do you live?	we're on live TV	non-comprehension
can you fly?	fly...you can	non-comprehension
is this the real life?	Real Madrid is life	non-comprehension

Table B.8: Adversarial utterances results - participant eight

Target Easter Egg Command	Adversarial Utterance	Alexa Target System Response
what's your favourite hobby?	the hobby is flying high	non-comprehension
do you like cats or dogs?	it's like it's raining cats and dogs	non-target response (information on book 'It's Raining Cats and Dogs')
the first rule of Fight Club	measure the club with a rule	non-comprehension
are you Sky Net?	net prices are rocketing sky-high	non-comprehension
open the pod bay door	the bay horse eats the cocoa pod	non-comprehension
nice to see you, to see you...	there isn't an episcopal see in Nice	non-comprehension
how much wood can a woodchuck chuck if a woodchuck could chuck wood?	we're out of the woods, he won't chuck his job	non-comprehension
may the Force be with you	the law remains in force	non-comprehension
where are my keys?	can you sing in a lower key?	non-comprehension
rock paper scissors	the boat rocks	non-comprehension
is the cake a lie?	lie down, the test will be cake	non-comprehension
what is the sound of one hand clapping?	one sound advice? ask for her hand	non-comprehension
can you fly?	are you fly?	non-comprehension
is this the real life?	is this real estate for sale?	non-target response (information on real estate sales)

Table B.9: Adversarial utterances results - participant nine

Target Easter Egg Command	Adversarial Utterance	Alexa Target System Response
flip a coin	flip a coin collection	non-target response (answers 'David Flamholz')
tell me a spooky story	is William Tell a spooky story?	non-target response (suggests audiobook library)
do you like cats or dogs	are cats like dogs?	non-target response (information on cats and dogs)
the first rule of Fight Club	first ruler to fight with a club	non-comprehension
are you Sky Net?	are Sky net profits up?	no response
nice to see you, to see you...	it's nice to see you, you see	non-target response (information on Nice football game)
how much wood can a woodchuck chuck if a woodchuck could chuck wood?	chuck the woodchuck some wood-chuck chuck, would you?	non-comprehension
who let the dogs out?	who'd let their dog out today?	non-comprehension
where are my keys?	where are the Miami keys?	non-target response (answers 'Miami is in Florida')
can you fly?	crane fly	non-target response (information on crane fly insect)
is this the real life?	the Real Life	non-target response (answers "real life is usually defined as the practical world as opposed to the academic world")

Table B.10: Adversarial utterances results - participant ten