



MASTER OF SCIENCE
IN ENGINEERING

Hes·SO

Haute Ecole Spécialisée
de Suisse occidentale

Fachhochschule Westschweiz

University of Applied Sciences and Arts
Western Switzerland

Master of Science HES-SO in Engineering
Av. de Provence 6
CH-1007 Lausanne

Master of Science HES-SO in Engineering

Orientation: Information and Communication Technologies (ICT)

GenBot

Author:

Romain Claret

`romain.claret@master.hes-so.ch`

Under the direction of:

Prof. Dr. Jean Hennebert

HES-SO/Fribourg

Institute of Complex Systems (iCoSys)

External expert:

Dr. Christophe Gisler

Lausanne, HES-SO//Master, June 1, 2019

Accepted by the HES-SO//Master (Switzerland, Lausanne) on a proposal from:

Prof. Dr. Jean Hennebert, deepening project supervisor
Dr. Christophe Gisler, iCoSys, main expert

Place, date: _____

Prof. Dr. Jean Hennebert
Supervisor

M. Philippe Joye
ICT MRU Leader at HES-SO//Fribourg

Contents

Contents	v
Acknowledgements	ix
Acronyms	xi
Abstract	xiii
1 Introduction	1
1.1 Aim of Study	1
1.2 Scope and Study Borders	1
1.3 Industrial Interest	1
2 Questions	3
2.1 Initial and Broad Questions	3
2.2 Narrowed Questions	4
2.3 Potential Red lines	4
2.4 The Deepening Project Question and Red line	4
3 Plan	5
3.1 Constraints	5
3.2 Initial Plan	5
3.2.1 Tasks	5
3.2.2 Milestones	6
3.2.3 Sprints	6
3.2.4 Gantt chart	7
3.3 Effective Plan	7
3.3.1 Tasks	7
3.3.2 Milestones	7
3.3.3 Gantt chart	7
4 State of the art	11
4.1 Chatbots	11
4.1.1 History of Chatbots	11
4.1.2 Narrow Chatbots	12
4.1.3 Information Retrieval (IR) Chatbots	12
4.1.4 Sequential Chatbots	12
4.1.5 Forwarding Chatbots	13
4.1.6 Learning Chatbots	13
4.1.7 Proactive Chatbots	13

Contents

4.1.8	Chatbot Examples	14
4.1.9	Narrow Chatbots compared to General Chatbots	14
4.1.10	General Chatbots	14
4.1.11	From Artificial Narrow Intelligence (ANI) to Artificial General Intelligence (AGI)	15
4.1.12	ANI	15
4.1.13	AGI	15
4.2	Natural Language Processing (NLP)	15
4.2.1	Natural Language (NL)	16
4.2.2	Current NLP technics	16
4.2.3	NLP declensions	17
4.3	Word Embedding	17
4.3.1	Word2Vec	17
4.3.2	Gensim Framework	17
4.4	Word2Vec Alternatives	18
4.4.1	Word2Vec[25]	18
4.4.2	FastText[8]	18
4.4.3	Glove[28]	19
4.4.4	Adagram[6]	19
4.5	Word Embedding Extensions	19
4.5.1	Doc2vec[22]	19
4.5.2	Skip-thought[20]	20
4.5.3	RNN	20
4.6	Beyond Word Embedding	20
4.6.1	Contextual embeddings	20
4.7	Datasets	21
4.8	Word2Vec Models	21
5	Analysis	23
5.1	Natural Language Processing	24
5.2	Pipeline	24
5.3	Word2Vec	24
5.3.1	Bag of words (CBOW)	24
5.3.2	Skip-Grams	24
5.3.3	Dimensions	24
5.3.4	N-Grams	24
5.3.5	Epochs	24
5.3.6	Lemmatization	24
5.3.7	Normalization	24
5.3.8	Distance and Cosine Angle	24
5.3.9	Training	24
5.3.10	Retrain Model	24
5.3.11	Memory Issues	24
5.3.12	Analogies	24
5.3.13	Proverbs	24
5.3.14	Evaluation	24
5.3.15	Visual Representation	24
5.3.16	Benchmarks	24
5.3.17	CPU VS GPU	24

Contents

5.3.18 Datasets	24
5.4 Chatbot	24
5.5 Proactivity	24
6 Experiments & Results	25
6.1 Environments	25
6.1.1 Jupyter Notebook	25
6.1.2 Local Machine	25
6.1.3 Amazon Web Services	25
6.1.4 iColab-gpu2	25
6.1.5 CPU Dedicated Machine	25
6.2 Gensim Framework	25
6.3 Materials	25
6.4 Products	25
7 Discussion	27
7.1 Next steps?	27
8 Conclusion	29
8.1 Word Embedding: Word2Vec	29
8.2 Framework: Gensim	29
8.3 Word2Vec Chatbots	29
8.4 Proactive Chatbots	29
Bibliography	31
Appendix	35
.0.1 Appendix	35

Acknowledgments

Acronyms

AGI

Artificial General Intelligence.

AI

Artificial Intelligence.

AIML

Artificial Intelligence Markup Language.

ANI

Artificial Narrow Intelligence.

ANN

Artificial Neural Networks.

BD

Big Data.

CBOW

Bag of words.

DL

Deep Learning.

DM

Data Mining.

DNN

Deep Neural Networks.

FAQ

Frequently Asked Questions.

ICT

Information and Communications Technologies.

IR

Information Retrieval.

Acronyms

ML

Machine Learning.

MRU

Master Research Units.

NL

Natural Language.

NLG

Natural Language Generation.

NLP

Natural Language Processing.

NLU

Natural Language Understanding.

NN

Neural Network.

Sci-Fi

Science Fiction.

SNN

Shallow Neural Network.

TF-IDF

Term Frequency-Inverse Document Frequency.

Abstract

In the scope of this deepening project, and as the technology of NLP is in constant evolution, we will be focusing on the exploration of the word embedding algorithm Word2Vec, which is, at the beginning of 2019, commonly used as a foundation for DNN Chatbots. As a result to this project, the student is demonstrating what is the Word2Vec technology, its extensions, and its applications.

Keywords: Word Embedding, Word2Vec, NLP, Natural Language Understanding (NLU), Machine Learning (ML), Data Engineering, Conversational Agent, Chatbot, Generic

Chapter 1

Introduction

Start of 2019, chatbots are everywhere but very limited to narrow tasks, and are, in most cases, sequences of if-else conditions resulting in a very weak Artificial Intelligence (AI). Indeed, hard-coded connections are requiring an infinite amount of human power to create generic Chatbots able to maintain a conversation at a human level. However, the progress in the field of ML is demonstrating that providing large corpora to an unsupervised algorithm is enough to maintain a passive conversation with users, which results into a shifting of the human power into data engineering. Multiple algorithms and technics are emerging monthly, demonstrating promising conversational performance improvements; however, they are all still narrow AI. Indeed, even if they are getting better at providing meaningful sentences, they are still not able to generalize all tasks linked to a conversation, such as, understanding the context, search and learn for missing information, initiate conversation in a meaningful manner, be intuitive, and more. The generalization of those features would allow a significant step forward into general Chatbots.

1.1 Aim of Study

In harmony with the author's interest, the goal of this deepening semester project is to suggest and demonstrate strategic approaches as a premise to the AGI and to get a step closer to general Chatbots, which can initiate and maintain human-like conversations in a pro-active manner.

1.2 Scope and Study Borders

As a red line for this deepening project, the focus will be on the Word2Vec technology, from a research perspective. Indeed, this technology is seen as a foundation for the modern NLP and DNN Chatbots, which makes it an exciting vector of study about its current usage, its extensions, and potential evolution.

1.3 Industrial Interest

iCoSys, the Institut of Complex Systems at University of Applied Sciences and Arts at Fribourg, Switzerland, is interested into the result of this project as a study for their AI-News project, whose goal is to provide a chatbot as a tool to reader, to help them narrow their interests and deliver the right information. AI-News is in

Chapter 1. Introduction

collaboration with the Swiss Innovation Agency from the Swiss Confederation, and La Liberté, the daily newspaper from Fribourg.

Chapter 2

Questions

To help the student find a red line to focus its research on; he was required to work on the subject: *"What should be the initial questions to ask in order to make AGI Chatbots"* as a preliminary study, before the beginning of deepening project itself; and to write down the outcome as a set of questions related to his interests and the field of AGI Chatbots.

2.1 Initial and Broad Questions

As a result of the preliminary study, the following questions were extracted. Please take into account that those questions were not meant to be answered as part of the project itself but as part of the process of appropriation of the field of study.

- Is the Artificial Neural Networks (ANN) approach appropriate to represent the world?
- Can agents be made exclusively from a language?
- Are agents able to experience an environment?
- Is a narrative environment enough to understand an environment?
- Is the language able to provide to an agent an understanding of the world?
- Is the knowledge of the language syntax enough to gain an understanding?
- Is the result of unsupervised learning enough to discover all nuances?
- Is the unsupervised learning sufficient to make sense to an environment?
- Is a descriptive explanation of the world in a language be enough to express it?
- Is the description good enough to catch all the nuances?
- Is the language good enough to explain?
- Can we augment or make a semantic language?
- Can we create a common symbolic language?
- Is the language multi-dimensional?
- How many dimensions are needed for a complex language?
- Is it possible to give a word equivalence to machines for human-specific words?
- Are all emotions describable into words?
- Are emotions altering language descriptions?
- Is an approximation of the real world enough to understand the environment?
- Would a the simulated world be a good approximation of the real world?

Chapter 2. Questions

2.2 Narrowed Questions

In a second time, the student was asked to narrow the initial questions above into potential fields of study.

- Common human-machine language
 - Is it possible to create a multi-dimensional human-machine language, which includes a common semantic, symbolic, and emotion definition?
 - Is it possible to create an abstract world for machines to understand human symbolic based on a real world, and define fundamentals for machine representation of the language?
- Machine intuition
 - Is it possible to provide to machines an human-like intuition (inside voice), which would help to keep a long term context and specialize in specific fields?
- Evaluate human-machine communication
 - Is it possible to provide a protocol to test the communication skills and machine understanding?

2.3 Potential Red lines

From the potential fields above, the following suggested red lines were proposed.

- How to quantify a chatbot understanding?
- What is the premise to make chatbots general with today's technology?
- How can chatbot be proactive?
- How to simulate human-like intuition in chatbots?

2.4 The Deepening Project Question and Red line

Based on reflective work and discussions, the concluding red line and question for this deepening project are:

- What is Word Embedding and can it be used to make chatbots proactive?

Chapter 3

Plan

3.1 Constraints

Timeframe: 15 weeks

Starting date: 18.02.2019

Ending date: 31.05.2019

3.2 Initial Plan

As the first milestone for the deepening project, the student was required to create an initial plan, with the purpose to help himself and the teacher to visualize the project's main red line.

3.2.1 Tasks

1. Initial research about general chatbots
2. Determine the project target
3. Play with the subject
4. Explore the Word2Vec methodology
5. Explore the Word2Vec extensions
6. Combine and test ANN algorithms with Word2Vec
7. Explore ANN algorithm topology for the chatbot
8. Analyze of the chatbot intuition with parallel algorithms
9. Analyze of a protocol to evaluate proactive chatbots
10. Analyze Profile-based initiatives
11. Analyze and experiment profile nurturing
12. Analyze and experiment with chatbot initiatives with no profiles
13. Make overall improvements
14. Autonomous data gathering
15. Make suggestions
16. Determine possible continuation and future outcomes for the project

Chapter 3. Plan

3.2.2 Milestones

1. Initial deepening project plan and specification document
2. Basic multi-dimensional word embedding space
3. Basic conversational agent
4. Basic proactive chatbot
5. Deepening project report

3.2.3 Sprints

18.02.19 to 08.03.19 (3 weeks)

- Do the initial research about general chatbots
- Determine the project target
- Play with the subject
- **DELIVERABLE:** Plan and Initial Specification document

11.02.19 to 29.03.19 (3 weeks)

- Explore the Word2Vec methodology and its extensions
- Combine and test ANN algorithms with Word2Vec
- **MVP:** Basic multi-dimensional word embedding space

01.04.19 to 19.04.19 (3 weeks)

- Explore ANN algorithm topology for the chatbot
- Analysis of the chatbot intuition with parallel algorithms
- Analysis of a protocol to evaluate proactive chatbots
- **MVP:** Basic conversational agent

22.04.19 to 10.05.19 (3 weeks)

- Profile-based initiatives
- Analysis and experiment of the profile nurturing
- Analyze and experiment with chatbot initiatives with no profiles.
- **MVP:** Basic proactive chatbot

13.05.19 to 31.05.19 (3 weeks)

- Overall improvements
- Autonomous data gathering
- Make suggestions
- Determine possible continuation and future outcomes for the project
- **DELIVERABLE:** Report + Sources

3.2.4 Gantt chart

Figure 3.1 represents the visual gantt chart for the initial plan.

3.3 Effective Plan

As expected the initial plan served as an initial model, and evolved iteratively based on the student and teacher feedback while exploring the subject.

3.3.1 Tasks

1. Initial research about general chatbots
2. Determine the project target
3. Set the initial plan
4. Make LaTeX report template
5. Explore the Word2Vec subject
6. Explore the Word2Vec algorithm
7. Build a Word2Vec model on the latest english wikipedia dump
8. Explore Word2Vec parameters
9. Explore Word2Vec analogies
10. Explore Word2Vec sentence generation
11. Explore visual representations of Word2Vec vectors
12. Explore Word2Vec applications with chatbots
13. Write the report

3.3.2 Milestones

1. Initial deepening project plan and specification document
2. Basic Word2Vec Word Embedding Model
3. Conclusions Word2Vec based chatbots
4. Ideas to make chatbots proactive
5. Deliver the report

3.3.3 Gantt chart

Figure 3.2 represents the visual gantt chart for the effective plan.

Chapter 3. Plan

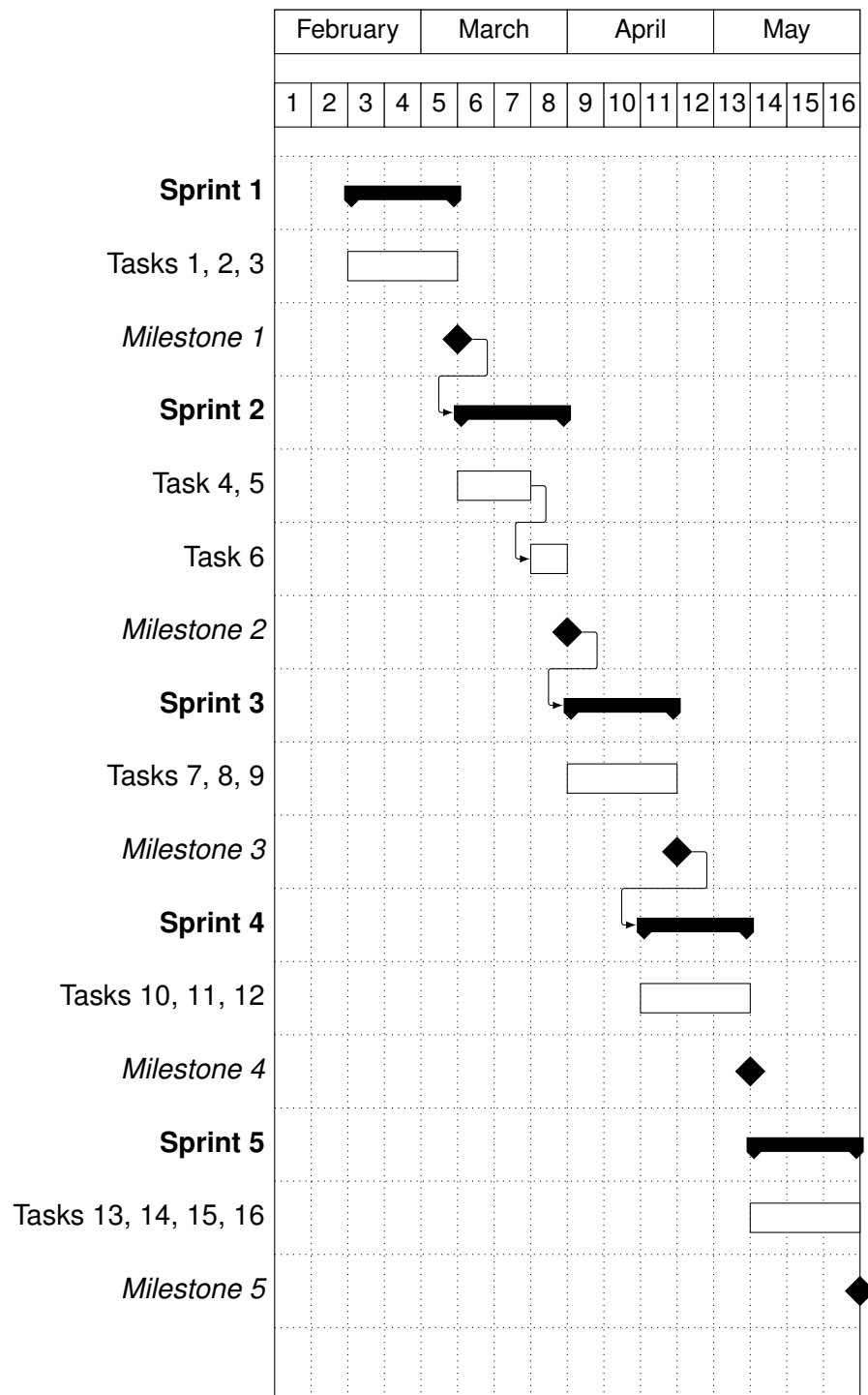


Figure 3.1: Initial Gantt Chart

3.3. Effective Plan

February					March				April				May			
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	

Figure 3.2: Effective Gantt Chart

Chapter 4

State of the art

4.1 Chatbots

From a user point of view, chatbots are trendy nowadays. Big companies such as *Google* or *Apple* are pushing to make the technology mainstream. Even if not every lambda people understand the word “chatbot”, they all have at least a mental representation of it. Indeed, whether they call it Digital Assistant, Siri, Ok Google, and so on, in the end, they all get the concept of an AI narrowed to more or less human-like conversations.

4.1.1 History of Chatbots

From when are they coming? Not mentioning *Alan Turing* or *Joseph Weizenbaum*, considered as the fathers of AI and chatbots, would not be fair. Indeed, they forecasted in 1950, that computers would be able to use human-like communication and they proposed a test to distinguish humans from machines, called the Turing Test[1]: where a human is asked to talk to a masked entity and determine if it is talking to a human or a computer. If the human cannot determine who is the computer, then the machine passed the Turing test, as seen on figure 4.1.

In 1966, Joseph Weizenbaum wrote *Eliza*[23], a computer program simulating a psychotherapist, seen as one of the first well-known attempts to make a Chatbot passing Turing test. Note that due to technical restrictions, *Eliza* is not performing well at all time. As it is for today, it is possible to play with it at on a dedicated website.

Since *Eliza*, a lot of progress has been made, indeed, to only cite a few noticeable chatbots: *Parry*[19] (1972), *Jabberwack*[35] (1988), *Dr. Sbaitso*[10] (1991), *A.L.I.C.E*[33] (1995), *Smarterchild*[34] (2001), *Watson*[18] (2006), *Siri*[5] (2010), *OK Google*[16] (2012), *Alexa*[4] (2014), *Cortana*[24] (2014), Facebook Bots[13] (2016), and *Tay*[29] (2016), which where all part of the Chatbot history [14].

From IF-ELSE, Artificial Intelligence Markup Language (AIML), up to ML with ANN and Deep Neural Networks (DNN), the improvement in the field of chatbots increased drastically over the years. At every iteration, the algorithms are becoming more sophisticated and better at using the human language, which is now called the field of the NLP and NLU.

Chapter 4. State of the art

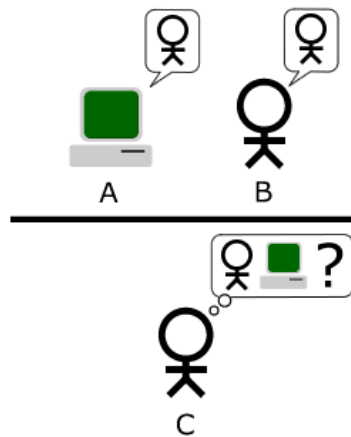


Figure 4.1: The "standard interpretation" of the Turing Test, in which player C, the interrogator, is tasked with trying to determine which player - A or B - is a computer and which is a human. The interrogator is limited to only using the responses to written questions in order to make the determination. [7]

4.1.2 Narrow Chatbots

Once again, chatbots are almost everywhere nowadays. Indeed, it became a common tool for companies of any size to communicate with their customers and a toy for users. However, most of the time, Chatbots are not understood by their users and is leading to a high level of frustration. Even if they are becoming increasingly mainstream and sophisticated, people do not realize their limits. Today's chatbots are often mistaken for AGI in Science Fiction (Sci-Fi) and are expected to do much more than they can do. Indeed, making ANI chatbots implies a specialization into a specific field.

Not to forget that the primary purpose of chatbots is to provide a conversational service to the user from text to vocal or even visual format. However, its purpose can be derivated in an almost unlimited amount of solutions such as Health, Weather, Customer Service, Games, and much more.

4.1.3 IR Chatbots

Most of the time used by Frequently Asked Questions (FAQ) chatbots, which are probably the most common type of chatbots, its goal is to answer specific questions, based on a specific keyword. Indeed, the communication skills are limited to pre-made sentences and a question/answer database, which often results, in the best case scenario, in a perfect match, or the worst case scenario, in the return of something unexpected.

Technically speaking, IR is part of the Data Mining (DM) in the field of ML. It is well suited for search engines, as it works in a query mode. Indeed, the algorithm tries to find the best match to the submitted query in its database, usually with pattern extraction and a rank.

4.1.4 Sequential Chatbots

"If he says this, then say that, then do so." This sentence is a good example of the concept of sequential chatbots. From a communication point of view, it does

4.1. Chatbots

not have to talk to accomplish its purpose; it is indeed usually based on a keyword detection technic to determine what pre-made action to do. However, as the whole system works on pre-made actions, the development of such algorithms requires a lot of brain power from the developers. Indeed, as all actions result from anticipated specific keywords, and even specific order of keywords, the complexity can quickly increase, which most of the time, makes sequential chatbots seen as command line terminals instead of conversational chatbots.

4.1.5 Forwarding Chatbots

Often used by companies for customer service, it has become the most popular type of chatbots and seen as a hybridization of the IR and sequential chatbots. Its goal is to simulate an agent that is available 24/7 to help the customer. Indeed, it will try its best to answer the most popular questions based on its FAQ database and forward the user seamlessly to a human agent if its knowledge is getting limited. In the best case scenario, it is greatly appreciated by the user as the transition from Chatbot to human is not noticeable.

4.1.6 Learning Chatbots

As ML evolves at an incredible rate and is boosted by DNN, new NLP algorithms emerges, and most of the time leaves the previous generation far behind. Modern learning chatbots algorithms are what come closer to human-like conversations. Leaving the algorithm progress alone through iterations on a large dataset or commonly named Big Data (BD) of real conversations, it will learn patterns by itself. However, the output generated by the trained model is dependent on the data the training occurred on. The most well-known example is *Tay*[29] (2016), the Twitter chatbot from Microsoft, that was influenced by the 4chan community to make it speak like a Young Racist Girl. However, it is essential to take note that learning chatbots have been existing for a long time now. *A.L.I.C.E*[33] had already basic learning skills, as AIML was taking care of saving variable on the run, such as the first name of the user. Even if this methodology could be seen archaic if compared to new Deep Learning (DL) algorithms such as LSTM, it is still used today likewise the AIML technology.

4.1.7 Proactive Chatbots

“Hey, I saw that you are on the website for some minutes now, do you need some more dedicated information?”. It is almost impossible that someone never received a message alike. Indeed, proactivity is not new in the field of chatbots. Mimicking an interest from the Chatbot to initiate conversations has become a standard in marketing and customer support chatbots. However, the limitations are hit fast, beyond asking general questions, not much progress has been made until now.

True proactive chatbots are implying that the algorithm is capable of initiating conversations from a human-like perspective, initiating the conversation or asking information in a meaningful manner based on the user, the context and the relationship with the user. The state-of-the-art search could not find any evidence of existing real proactive chatbots as described.

Chapter 4. State of the art

4.1.8 Chatbot Examples

As a help to get a feeling about narrow chatbots, a none-exhaustive list of applications is available below, and for more references about chatbots, Chatbot .org[9] is an excellent, up to date, place about referencing old and new chatbots.

- Receive relevant information about a trip, book flights, and hotels, and get updated on the boarding and weather conditions at the airport.
- Keep track and order coffee remotely at the office.
- Monitor customer's satisfaction.
- Convert potential customer into paying customers by interacting with them at the right moment.
- Personal assistant on-the-go, get the schedule and the next meetings.
- Relay for people on hold at a service.

4.1.9 Narrow Chatbots compared to General Chatbots

Before going further into the world of general chatbots, it is required to understand the following two axes of AI defined by Tasks and Knowledge. Indeed, narrow chatbots are limited by the range of tasks they can accomplish and the knowledge they can use. However, most of the time, they are very good at a particular task for a particular knowledge requirement. The table 4.1 tries to represent the position of Narrow and General Chatbots on those two axes.

Tasks: Talk, FAQ, Remote Control, Customer, or Placing orders are just a few tasks that a chatbot could accomplish.

Knowledge: Health, Weather, Customer Service, or Games are just a few knowledge examples chatbots could excel at.

Knowledge	
Tasks	Expert in a specific Field Expert at all Tasks
	General Chatbots Expert in all Fields Expert at all Tasks
Narrow Chatbots Expert in a specific Field Expert at specific Task	Expert in all Fields Expert at specific Task

Table 4.1: Tasks versus Knowledge in the field of Chatbots

4.1.10 General Chatbots

Much effort is being made to get chatbots that can perform well simultaneously in various tasks and knowledge. Indeed, general chatbots are not limited to previously learned tasks and subjects; they should also be able to learn and relearn. Those type of chatbots have not been found during the state of the art phase, and are

probably by this mean either none-existent at the moment or hidden in laboratories, far from public knowledge. However, big companies like Amazon are providing to the public a feel of general chatbots with *Alexa*[4]. Users can converse with it, command their smart houses, use it as a personal assistant, and even program it to perform custom actions. However, it is not yet able to learn by itself and generate out of the blue non-programmed skills. Note that general chatbots could be scary for lambda people if it starts mimicking human being too well, as in the user mind, talking to a machine should be differentiable from talking to humans. Admittedly, in the case of the *Turing Test*[1], the human does not know if it is talking to a human or a machine, which makes it probably more comfortable to accept than talking to a machine directly. Sci-Fi is conditioning people to believe that human-like performing machines are dangerous for the human species.

4.1.11 From ANI to AGI

On a side, even if it is not part of the deepening project, it is interesting to write a few lines about AI. New incredible algorithms outperforming the previous one, and experiments reports are emerging almost every month and redefining the standard of AI. Paradigms are shifting and technologically speaking; we are entering a new era of computer-assisted humankind.

4.1.12 ANI

More than a sequential algorithm, narrow artificial intelligence in modern terminology is the definition of “being good at something”. ANI has been made possible with the huge progress in ML, the arrival of the DL, and the need for humans to store data about everything (BD). In medicine, for instance, it is sometimes performing so well, that humans, who spent years studying are left behind by an algorithm trained on large datasets for a few days.

4.1.13 AGI

The next step into the field of AI, when supervision has been banned as a teaching method for algorithms as the human interaction is inputting more errors than machine themselves if unsupervised. In addition to teaching themselves, algorithms are teaching each other, and improve over the iteration with auto corrections and optimizations. They are excelling at all tasks requiring repetition, precision, and safety. Besides, they are also all able to retrieve any available information and use it for their need. “In the future, machines will be able to understand and do everything, much more efficiently than humans.”

4.2 NLP

Present in our daily lives, this technology is used massively to automate the extract information from human communication. In other words, it is seen as the given skill to machines to comprehend human language.

Examples The following is an non-exhaustive list of NLP use:

- Customer Support Chatbots

Chapter 4. State of the art

- Translation into foreign languages
- Voice recognition
- Spam filters
- Interpret written queries
- Generate the responses

4.2.1 NL

Naively, it is the language naturally used by humans. The goal of the NLP is to mimic the NL to create a human-like verbal interaction. However, it is not an easy task as it is nearly not possible to teach a machine to talk like a human. Indeed, even if a machine was given the same language rules as humans, they do not understand by themselves, they are just applying the rules, which results in a problem during ambiguities. It would be necessary to sequentially teach the missing pieces of information, which would result in an almost an unlimited amount of conditions.

NL decomposition Beyond the grammar and orthography, human language is composed of an incredible amount of subtleties, which makes sense most of the time intuitively for humans, but not for machines. To help understand the complexity behind NL, the following are list expresses the foundation of human language:

- Semantics: express the relations between words, sentences, paragraphs, etc.
- Morphology: maintain a structure and the content of word forms
- Phonology: sounds used to express words
- Syntax: rules applied to the bag of words to create valid texts
- Pragmatics: how the context influences the meaning words

4.2.2 Current NLP technics

Most of the following technics have been developed in the IR field.

- Term Frequency-Inverse Document Frequency (TF-IDF): Used to set the word importance in corpora.
- CBOW [5.3.1]: Counts the words occurrences throughout in corpus.
- Skip-Grams [5.3.2]: Counts the occurrences of the character throughout in corpus.
- Topic modeling: Text clustering providing meaningful information to discover hidden structures via text chunking to identify the parts of the sentence in relation to each other.
- Segmentation: Split corpus into predefined parts, such as: *sentences, paragraphs, chapters, etc.*
- Tokenization: Split the sentences into words.
- Tagging: Based on a pre-made dictionary, it gives a new layer of meaning to the word, such as: *verb, adverb, noun, people name, locations, number, etc.*
- Dictionary: Use of tokenized words to build a dictionary, which could contain the word occurrences.

4.3. Word Embedding

- Stop Words: Ignoring words only used as liaisons, and does not contain information, such as: *and, or, etc.*
- Stemming: Uniformizing words to their root by removing the prefix and suffix, such as: *remake and loveable.*
- Lemmatization: Replace the words to their base form, such as *conjugated verb.*

4.2.3 NLP declensions

Further, into the field of NLP, it is today commonly split into two groups:

NLU It is a subdivision of DM, and it involves the processing of the text by analyzing its content by extracting relevant information, usually called keywords.

Natural Language Generation (NLG) In combination with NLU, often applied to text classification, NLG is useful to generate custom sentences using custom keyword extracted to make the response to a query even more relevant and usually keeps track of the context.

4.3 Word Embedding

It can be summarised as the vector representation of a word and often using between 100 and 400 dimensions. Its position in the multi-dimensional space keeps track of the word context and semantic to a dictionary of words and the corpora. Due to the vector nature of the words, geometrical operations can be applied to those words to find word similarities and relationship between them.

4.3.1 Word2Vec

Published Google in 2013, Word2Vec[25], probably became, the most popular algorithm in the word embedding field, nowadays. It uses a Shallow Neural Network (SNN) 4.2, similar to a conventional supervised model. Indeed, it is a two-layer Neural Network (NN), its input is text corpora, and its output is word vectors based on a given dictionary. Even if it is easy to train and test, it is often difficult to tweak the algorithm, and as a result, makes it harder to make a good generalization. Even if it is not using DL itself as output, the input text form of the words are transformed into their value form, which makes it incredibly powerful and useful for DNN algorithms as input.

4.3.2 Gensim Framework

Since its publication in 2013, Word2Vec[25] various actors implemented the algorithm and companies like RaRe Technologies[31], specialized in NLP, made it easier for Scientist and Hobbyists to jump right in. Which probably influenced increased its popularity by making it accessible for the community instead of big companies, institutions only. In the case of this deepening project, the author will be focusing on the framework made available by RaRe Technologies and more specifically by RaRe Consulting[30], Gensim[32]. Other frameworks and solutions are available, they are in the end implementing the same algorithm generating a Word2Vec model

Chapter 4. State of the art

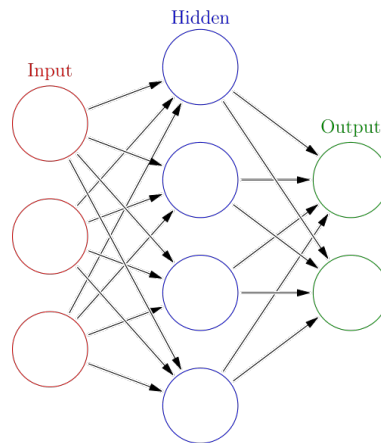


Figure 4.2: Artificial neural network with layer coloring [15]

by capturing the context of the words, but they are different by their language and their different integration with custom features.

Gensim is a python implementation of Word2Vec, and it is not a general purpose ML or even DL framework. However, in order of magnitude, it does one thing, but does it very well, as it does its job faster than Tensorflow[2] for instance.

4.4 Word2Vec Alternatives

Word2Vec does not define Word Embedding; indeed the concept of the vector representation in a multi-dimensional space has multiple solutions to it. Without going into too many details, and to name a few, the following could be alternatives to Word2Vec, with their pros and cons.

4.4.1 Word2Vec[25]

Pros

- First word embedding solution to be able to generate a model on a large corpus with a dictionary containing millions of words.
- Outputs word vectors based on corpora with raw text.

Cons

- Difficulties to extract the sense of words with multiple meanings depending on the context. e.g., the word **Doctor**, it could be the Academic Title, a Physician, or even the name of a TV-Show.

4.4.2 FastText[8]

Made by Facebook.

Pros

- Same as Word2Vec, except that technically it uses the Skip-Grams [5.3.2] technique to train on characters composing the words instead of CBOW [5.3.1], which trains with words.

4.5. Word Embedding Extensions

Cons

- Same as Word2Vec, words with multiple meanings are not managed well.

4.4.3 Glove[28]

Global Vectors for Word Representation is a contribution to the Word Embedding by the University of Stanford in NLP.

Pros

- Less time consuming than Word2Vec.
- Depending on the scenario and benchmarks, it sometimes performs better than Word2Vec at tasks related to semantic.

Cons

- Larger memory usage than Word2Vec.
- Same as Word2Vec for multiple word meanings.

4.4.4 Adagram[6]

Adaptive Skip-Gram is a Russian contribution by National Research University of Moscow.

Pros

- It claims, contrary to Word2Vec, to be able to manage different word meanings as it should extract the context of the surrounding words.

Cons

- In the current state, it is not designed for corpora with a dictionary larger than tens of millions of words.
- Does not keep track of the word order.

4.5 Word Embedding Extensions

Based on the Word Embedding technics, proposals were raised about its extension to the sentences and even documents. A simple solution to generate a sentence/documents representation would be to sum word vectors composing the sentences/documents; however, the following non-exhaustive technics are performing better than naive addition.

4.5.1 Doc2vec[22]

Adaptation of Word2Vec for document embedding.

Pros

- Based on Word2Vec.
- Performs well in most cases.

Chapter 4. State of the art

Cons

- In few cases the embedding could be biased towards the specific content words.[21]

4.5.2 Skip-thought[20]

Made for corpus with semantically related sentences.

Pros

- Works well with corpus having a sentence continuity.

Cons

- Adjacent sentences must be semantically related.

4.5.3 RNN

With the current market and institutional need to make everything DL, the subject of DNN must be slightly overviewed. All the technics described previously in this chapter are not using DNN, and there is a good reason for their success. Indeed, they do not require labeled data for training, which is required by DL algorithms. However, the idea has not been abandoned; solutions are raising to overcome this drawback, such as using crowd-based solution using users as signals to determine document similarities [27].

4.6 Beyond Word Embedding

As NLP usage increases over the years, Word Embedding technics and its extensions are becoming increasingly more sophisticated and are getting closer to human-like generalization. As controversially suggested by Geoffrey Hinton, famous DL researcher, it would be possible to get to human-like conversational capabilities via a method he calls Thought Vectors [11]. Without going into exciting details, it implies for the AGI, at this stage, even if the algorithm does not understand the meaning of the sentences, the reasoning behind a thought would be well enough emulated to make it human-like.

Thought Embedding A vectorized thought would be trained to generate a thought's context. As for Word Embedding and Doc Embedding, Thought Embedding are linked by a chain of reasoning.

4.6.1 Contextual embeddings

Though Embedding paper and implementation is yet to be made, however, progress has been made in the direction with the contextual embeddings Context2Vec[26]. It is a bidirectional LSTM[17] unsupervised model generating a word Embedding based on its occurrence in the sentences, which unlike Adagram[4.4.4] is taking into account order.

However, Context2Vec does not define Contextual embeddings. New emerging 2019 algorithms, based on the attention mechanism [36], such as Transformers[3]

or even further its bidirectional extension BERT[12], which makes LSTM almost obsolete.

Summarized

- Context dependent word embeddings.
- Can generate sentence embeddings.
- The output can be used almost as it is for NLP.
- Tracking using selective "forget" gates.

4.7 Datasets

Nowadays, in data engineering, the new gold is data. Luckily, today is driven by BD, and almost any kind of data is available to who knows where to look at. For the case of the deepening project, the author is interested in conversational data. Even though the English Wikipedia Dump from 09th May 2019 will be using exclusively, the following is a non-exhaustive list of corpora gathered during the deepening project.

- Wikipedia Dumps Index: <https://dumps.wikimedia.org/backup-index.html>
- English Wikipedia Dumps (bzip2: 16Gb, raw: 90Gb): [https://dumps.wikimedia.org/enwiki/latest/\[enwiki-latest-pages-articles.xml.bz2\]](https://dumps.wikimedia.org/enwiki/latest/[enwiki-latest-pages-articles.xml.bz2])
- Reddit Comments (bzip2: 6Gb, raw: 32Gb): https://www.reddit.com/r/datasets/comments/3bxlg7/i_have_every_publicly_available_reddit_comment/
- The Open American National Corpus: <http://www.anc.org>
- Santa Barbara Corpus of Spoken American English: <https://www.linguistics.ucsb.edu/research/santa-barbara-corpus>
- Leipzig Corpora Collection <http://wortschatz.uni-leipzig.de/en/download/>
- Legal Case Reports: <http://archive.ics.uci.edu/ml/datasets/Legal\+Case\+Reports>
- Cornell Newsroom: <https://summari.es>
- DeepMind Q&A: <https://cs.nyu.edu/~kcho/DMQA/>
- Large Movie Review: <http://ai.stanford.edu/~amaas/data/sentiment/>
- Project Gutenberg - Free eBooks: <https://www.www.gutenberg.org>

4.8 Word2Vec Models

Training own models are very resource consumptive, and often the resources are not available, it could be datasets or computer power. Luckily, big companies like Google did the work for us, and pre-trained models that could be used out of the box. The following is a non-exhaustive list of models for Word2Vec:

- Gensim directory: <https://github.com/RaRe-Technologies/gensim-data/releases>
- Fasttext: <https://fasttext.cc/docs/en/english-vectors.html>
- estnltk: <http://ats.cs.ut.ee/keelethehnoloogia/estnltk/word2vec/>

Chapter 5

Analysis

5.1 Natural Language Processing

5.2 Pipeline

5.3 Word2Vec

5.3.1 CBOW

5.3.2 Skip-Grams

5.3.3 Dimensions

5.3.4 N-Grams

5.3.5 Epochs

5.3.6 Lemmatization

5.3.7 Normalization

5.3.8 Distance and Cosine Angle

5.3.9 Training

5.3.10 Retrain Model

5.3.11 Memory Issues

5.3.12 Analogies

5.3.13 Proverbs

5.3.14 Evaluation

5.3.15 Visual Representation

5.3.16 Benchmarks

5.3.17 CPU VS GPU

5.3.18 Datasets

5.4 Chatbot

5.5 Proactivity

Chapter 6

Experiments & Results

6.1 Environments

6.1.1 Jupyter Notebook

6.1.2 Local Machine

6.1.3 Amazon Web Services

6.1.4 iColab-gpu2

6.1.5 CPU Dedicated Machine

6.2 Gensim Framework

Errors Memory allocation with multi-core. The problem is occurring during the merge of the cores. Indeed, my current machine has 128GB ram, and the dataset weights about 16GB in the memory, and each core during merging is processing at least the same amount, plus the processed informations.

```
2019-03-25 08:31:18,867 : INFO : PROGRESS: pass 0, dispatched chunk #34 = documents up to #70000/4614519, outstanding queue
size 31
Exception in thread Thread-1:
Traceback (most recent call last):
  File "/usr/lib/python3.5/threading.py", line 914, in _bootstrap_inner
    self.run()
  File "/usr/lib/python3.5/threading.py", line 862, in run
    self.target(*self.args, **self.kwargs)
  File "/usr/lib/python3.5/multiprocessing/pool.py", line 366, in _handle_workers
    pool._maintain_pool()
  File "/usr/lib/python3.5/multiprocessing/pool.py", line 240, in _maintain_pool
    self._repopulate_pool()
  File "/usr/lib/python3.5/multiprocessing/pool.py", line 233, in _repopulate_pool
    w.start()
  File "/usr/lib/python3.5/multiprocessing/process.py", line 105, in start
    self._popen = self._Popen(self)
  File "/usr/lib/python3.5/multiprocessing/context.py", line 267, in _Popen
    return Popen(process_obj)
  File "/usr/lib/python3.5/multiprocessing/popen_fork.py", line 20, in __init__
    self._launch(process_obj)
  File "/usr/lib/python3.5/multiprocessing/popen_fork.py", line 67, in _launch
    self.pid = os.fork()
OSError: [Errno 12] Cannot allocate memory
```

Figure 6.1: Error 1

6.3 Materials

6.4 Products

Chapter 6. Experiments & Results

```
2019-03-25 08:31:18,867 : INFO : PROGRESS: pass 0, dispatched chunk #34 = documents up to #70000/4614519, outstanding queue
size 31
Exception in thread Thread-1:
Traceback (most recent call last):
  File "/usr/lib/python3.5/threading.py", line 914, in _bootstrap_inner
    self.run()
  File "/usr/lib/python3.5/threading.py", line 862, in run
    self.target(*self.args, **self.kwargs)
  File "/usr/lib/python3.5/multiprocessing/pool.py", line 366, in _handle_workers
    pool.maintain_pool()
  File "/usr/lib/python3.5/multiprocessing/pool.py", line 240, in _maintain_pool
    self.repopulate_pool()
  File "/usr/lib/python3.5/multiprocessing/pool.py", line 233, in _repopulate_pool
    w.start()
  File "/usr/lib/python3.5/multiprocessing/process.py", line 105, in start
    self._popen = self._Popen(self)
  File "/usr/lib/python3.5/multiprocessing/context.py", line 267, in _Popen
    return Popen(process_obj)
  File "/usr/lib/python3.5/multiprocessing/popen_fork.py", line 20, in __init__
    self._launch(process_obj)
  File "/usr/lib/python3.5/multiprocessing/popen_fork.py", line 67, in _launch
    self.pid = os.fork()
OSError: [Errno 12] Cannot allocate memory
```

Figure 6.2: Error 2

Chapter 7

Discussion

7.1 Next steps?

Chapter 8

Conclusion

8.1 Word Embedding: Word2Vec

8.2 Framework: Gensim

8.3 Word2Vec Chatbots

8.4 Proactive Chatbots

Lausanne, June 1, 2019

Romain Claret

Bibliography

- [1] A. M. Turing. *COMPUTING MACHINERY AND INTELLIGENCE*. [Online; accessed 26-May-2019]. 1950. URL: <https://www.csee.umbc.edu/courses/471/papers/turing.pdf>.
- [2] Martín Abadi et al. *TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems*. Software available from tensorflow.org. 2015. URL: <https://www.tensorflow.org/>.
- [3] Tim Salimans Alec Radford Karthik Narasimhan and Ilya Sutskever. “Improving Language Understanding by Generative Pre-Training”. In: *Preprint. Work in progress*. 2018. URL: https://s3-us-west-2.amazonaws.com/openai-assets/research-covers/language-unsupervised/language_understanding_paper.pdf.
- [4] Amazon Lab126. *Amazon Alexa*. [Online; accessed 26-May-2019]. 2014. URL: <https://developer.amazon.com/alexa>.
- [5] Apple. *Siri*. [Online; accessed 26-May-2019]. 2010. URL: <https://machinelearning.apple.com/2018/04/16/personalized-hey-siri.html>.
- [6] Sergey Bartunov et al. “Breaking Sticks and Ambiguities with Adaptive Skip-gram”. In: *arXiv e-prints*, arXiv:1502.07257 (2015), arXiv:1502.07257. arXiv: 1502.07257 [cs.CL].
- [7] Bilby. *File:Turing Test version 3.png*. [Online; accessed 26-May-2019]. 2008. URL: https://commons.wikimedia.org/wiki/File:Turing_Test_version_3.png.
- [8] Piotr Bojanowski et al. “Enriching Word Vectors with Subword Information”. In: *arXiv e-prints*, arXiv:1607.04606 (2016), arXiv:1607.04606. arXiv: 1607.04606 [cs.CL].
- [9] Chatbots.org Team. *Chatbots.org*. [Online; accessed 26-May-2019]. 2005. URL: <https://www.chatbots.org/>.
- [10] Creative Labs. *Dr. Sbaitso*. [Online; accessed 26-May-2019]. 1991. URL: <http://ccftp.creative.com/manualdn/Manuals/TSD/2389/sblive.pdf>.
- [11] Hannah Devlin. “Google a step closer to developing machines with human-like intelligence”. In: *The Guardian* (May 21, 2015). URL: <https://www.theguardian.com/science/2015/may/21/google-a-step-closer-to-developing-machines-with-human-like-intelligence> (visited on 05/26/2019).
- [12] Jacob Devlin et al. “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding”. In: *arXiv e-prints*, arXiv:1810.04805 (2018), arXiv:1810.04805. arXiv: 1810.04805 [cs.CL].

Bibliography

- [13] Facebook. *Bots for Messenger*. [Online; accessed 26-May-2019]. 2016. URL: <https://developers.facebook.com/blog/post/2016/04/12/bots-for-messenger/>.
- [14] Futurism, LLC. *The History of Chatbots Infographic*. [Online; accessed 26-May-2019]. 2016. URL: <https://futurism.com/images/the-history-of-chatbots-infographic>.
- [15] Glosser.ca. *File:Colored neural network.svg*. [Online; accessed 26-May-2019]. 2013. URL: https://commons.wikimedia.org/wiki/File:Colored_neural_network.svg.
- [16] Google. *Google Assistant*. [Online; accessed 26-May-2019]. 2012. URL: <https://assistant.google.com>.
- [17] Sepp Hochreiter and Jürgen Schmidhuber. "Long Short-Term Memory". In: *Neural Comput.* 9.8 (Nov. 1997), pp. 1735–1780. ISSN: 0899-7667. DOI: 10.1162/neco.1997.9.8.1735. URL: <http://dx.doi.org/10.1162/neco.1997.9.8.1735>.
- [18] IBM. *Watson*. [Online; accessed 26-May-2019]. 2007. URL: [https://researcher.watson.ibm.com/researcher/files/us-heq/W\(3\)%20INTRODUCTION%2006177724.pdf](https://researcher.watson.ibm.com/researcher/files/us-heq/W(3)%20INTRODUCTION%2006177724.pdf).
- [19] Ken Colby. *PARRY: Paranoia mental hospital patient*. [Online; accessed 26-May-2019]. 1995. URL: <https://www.cs.cmu.edu/afs/cs/project/ai-repository/ai/areas/classics/parry/0.html>.
- [20] Ryan Kiros et al. "Skip-Thought Vectors". In: *arXiv e-prints*, arXiv:1506.06726 (2015), arXiv:1506.06726. arXiv: 1506.06726 [cs.CL].
- [21] Jey Han Lau and Timothy Baldwin. "An Empirical Evaluation of doc2vec with Practical Insights into Document Embedding Generation". In: *Proceedings of the 1st Workshop on Representation Learning for NLP*. 2016, 78–86. URL: <http://www.aclweb.org/anthology/D14-1162>.
- [22] Quoc V. Le and Tomas Mikolov. "Distributed Representations of Sentences and Documents". In: *arXiv e-prints*, arXiv:1405.4053 (2014), arXiv:1405.4053. arXiv: 1405.4053 [cs.CL].
- [23] Michal Wallace & George Dunlop. *Eliza, the Rogerian Therapist*. [Online; accessed 26-May-2019]. 1999. URL: <http://psych.fullerton.edu/mbirnbaum/psych101/Eliza.htm>.
- [24] Microsoft. *Cortana*. [Online; accessed 26-May-2019]. 2014. URL: <https://www.microsoft.com/en-us/research/group/cortana-research/>.
- [25] Tomas Mikolov et al. "Efficient Estimation of Word Representations in Vector Space". In: *arXiv e-prints*, arXiv:1301.3781 (2013), arXiv:1301.3781. arXiv: 1301.3781 [cs.CL].
- [26] Jacob Goldberger Oren Melamud and Ido Dagan. "context2vec: Learning Generic Context Embedding with Bidirectional LSTM". In: *Proceedings of the 20th SIGNLL Conference on Computational Natural Language Learning (CoNLL)*. 2016, 51–61. URL: <https://www.aclweb.org/anthology/K16-1006>.

- [27] Hamid Palangi et al. “Deep Sentence Embedding Using Long Short-Term Memory Networks: Analysis and Application to Information Retrieval”. In: *arXiv e-prints*, arXiv:1502.06922 (2015), arXiv:1502.06922. arXiv: 1502.06922 [cs.CL].
- [28] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. “GloVe: Global Vectors for Word Representation”. In: *Empirical Methods in Natural Language Processing (EMNLP)*. 2014, pp. 1532–1543. URL: <http://www.aclweb.org/anthology/D14-1162>.
- [29] Peter Lee. *Learning from Tay's introduction*. [Online; accessed 26-May-2019]. 2016. URL: <https://blogs.microsoft.com/blog/2016/03/25/learning-tays-introduction/>.
- [30] Radim Řehůřek. *RaRe Consulting*. [Online; accessed 26-May-2019]. 2011. URL: <https://radimrehurek.com>.
- [31] Radim Řehůřek. *RaRe Technologies*. [Online; accessed 26-May-2019]. 2013. URL: <https://rare-technologies.com>.
- [32] Radim Řehůřek and Petr Sojka. “Software Framework for Topic Modelling with Large Corpora”. English. In: *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*. <http://is.muni.cz/publication/884893/en>. Valletta, Malta: ELRA, May 2010, pp. 45–50.
- [33] Richard Wallace. *A.L.I.C.E. A.I Foundation*. [Online; accessed 26-May-2019]. 1995. URL: <http://www.alicebot.org/>.
- [34] Robert Hoffer, Timothy Kay and Peter Levitan. *SmarterChild*. [Online; accessed 26-May-2019]. 2001. URL: <https://patents.google.com/patent/EP1767015A1/da>.
- [35] Rollo Carpenter. *Jabberwacky*. [Online; accessed 26-May-2019]. 1997. URL: <http://www.jabberwacky.com>.
- [36] Ashish Vaswani et al. “Attention Is All You Need”. In: *arXiv e-prints*, arXiv:1706.03762 (2017), arXiv:1706.03762. arXiv: 1706.03762 [cs.CL].

Appendix

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

.0.1 Appendix

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

