

Automatically responding to customers

Master Thesis

Rik Huijzer

Supervisor:
dr. N. Yakovets

Examination committee:
dr. N. Yakovets
dr. G.H.L. Fletcher
dr. J. Vanschoren

Eindhoven, January 2019

Abstract

In the last few years artificial intelligence has considerably advanced the natural language processing (NLP) field. The graduation company is interested in seeing whether this can be used to automate customer support. NLP has evolved to contain many tasks. Intent classification is used to classify the intent of a sentence like ‘what is the weather in London tomorrow?’. The intent for this sentence could be ‘get_weather’. Named-entity recognition (NER) aims to extract information from subparts of the sentence. For example, ‘London’ is a location and ‘tomorrow’ is a date. Intents and entities are used by chatbots to understand the text written by users. This has caused intent classification and NER to have the following practical constraints. The text should be analysed in real-time and training data consists of a few dozen training examples. The latter makes it an interesting problem from a machine learning perspective.

Multiple systems and services provide intent classification and NER. Accuracy of classification differs per system. Higher accuracy means responding correctly to customer utterances more often. Many systems claim to make the fewest mistakes during classification when comparing their system to others. To validate this a benchmarking tool is created. This tool is aimed on creating comparisons in such a way that users can easily run new or re-run existing evaluations. The code can be extended to allow comparison of more datasets and systems.

To improve the accuracy of intent classification and NER, deep learning architectures for NLP have been investigated. New accuracy records are set every few months for various NLP tasks. One of the most promising systems at the time of writing is considered. This system, Google BERT, uses context from both sides of some word to predict the meaning of the word. For example, the meaning of the word ‘bank’ differs in the sentences ‘river bank’ and ‘bank account’. BERT has shown state-of-the-art results for eleven NLP tasks. An attempt is made to apply the model to intent classification. Compared to baseline models applied by industry, this obtained significant increases in running time, but not in accuracy. A second attempt trained the system jointly on intent classification and NER. BERT is well-suited for joint training, because it uses context in all hidden layers of the network. Information from the intent classification task is used, in all layers, for making NER predictions and vice versa. It is shown that joint training with BERT is feasible, and can be used to lower training time when compared to separate training of BERT. Future work is needed to see whether the improvements in accuracy are significant.

Preface

This master thesis is aimed to graduate the Computer Science and Engineering master at Eindhoven University of Technology. The project is carried out at Devhouse Spindle in Groningen. The thesis concludes my education in Eindhoven. I consider choosing this study as one of the best decisions of my life. Getting good grades for assignments and exams was tough, but grading was fair. Teachers were always busy, but also always willing to help. I would like to thank Nikolay Yakovets for guiding the project. I appreciate the fact that Nikolay allowed me to choose a subject that does not directly aid him in his research. Furthermore, I would like to thank Devhouse Spindle for allowing me to graduate at their company. Little could be improved about the workplace, colleagues and office in general. I would like to thank my graduation tutor Jukka Koivunen for the guidance, and helping me fix any bugs related to object-oriented code. Thanks to Jukka and Hylke for giving feedback on the thesis. Thanks to all colleagues for the relaxed work environment, the advices and for helping me fix bugs I would have been stuck on for many hours.

In general I would like to thank Rodger, Erik, Mark, Leon, Stan, Astrid, Sjors, Nol, Tom, Suzanne, Reggie, Imke, the student ice hockey team and the student rugby team, to help me getting through the study. Finally, special thanks go out to my parents and siblings for motivating me to aim high.

Rik Huijzer
1 February 2019

Contents

Contents	vii
List of abbreviations	ix
1 Introduction	1
1.1 Thesis context	1
1.2 Problem description	1
1.3 Project goal and outline	2
2 Preliminaries	5
2.1 Natural language processing	5
2.1.1 Language model	6
2.1.2 Machine translation	6
2.1.3 Natural language understanding	8
2.1.4 F score	8
2.2 Deep learning	9
2.2.1 Vanishing gradient problem	9
2.2.2 Recurrent neural networks	9
2.2.3 Gated recurrent units and LSTMs	10
2.2.4 Bidirectional recurrent neural networks	11
2.2.5 Convolutional neural networks	12
2.2.6 ELMo	12
2.2.7 Transformers	12
3 Benchmarking	15
3.1 Datasets	15
3.1.1 Format	15
3.1.2 Available datasets	16
3.2 Systems	17
3.2.1 Rasa	17
3.2.2 DeepPavlov	18
3.2.3 Cloud services	18
3.3 Tool and results	19
3.3.1 Overview	19
3.3.2 Benchmark results	21
3.4 Observations	22
3.4.1 Benchmarking system	22
3.4.2 Methodology	23

4	Improving accuracy	25
4.1	Search	25
4.1.1	Duplicate finding	25
4.1.2	Using data	25
4.1.3	Kaggle	26
4.1.4	Meta-learning	26
4.1.5	Embeddings	26
4.2	BERT	27
4.2.1	Model description	27
4.2.2	Training	28
4.2.3	Joint training	29
4.2.4	BERT joint training	30
4.3	Results	31
4.4	Implementation improvements	32
5	Conclusions	35
	Bibliography	37
	Appendix	45
A	bench usage	45
B	Notes on functional programming in Python	46
B.1	Mapping to functions	46
B.2	NamedTuple	47
B.3	Function caching	47
B.4	Lazy evaluation	48
C	Lazy evaluation demonstration	49
C.1	Benefits	49
C.2	Code	49
C.3	Output	50
D	Intent F1 score calculations	52
E	improv	54
E.1	Usage	54
E.2	TPU and the Estimator API	54
E.3	Additional experiment	55
F	Runs	56

List of abbreviations

Abbreviation	Phrase	Section
AI	Artificial intelligence	
API	Application programming interface	
ASIC	Application-specific integrated circuit	
BRNN	Bidirectional recurrent neural network	2.2.4
BERT	Bidirectional Encoder Representations from Transformers	4.2
CNN	Convolutional neural network	2.2.5
ELMo	Embeddings from Language Models	2.2.6
F ₁	F ₁ score	2.1.4
GB	Gigabyte	
GPU	Graphic processing unit	
GRU	Gated recurrent unit	2.2.3
LSTM	Long short-term memory	2.2.3
LUIS	(Microsoft) Language Understanding Intelligent Service	3.2.3
MB	Megabyte	
NER	Named-entity recognition	2.1.3
NLP	Natural language processing	2.1
RAM	Random-access memory	
SQuAD	The Stanford Question Answering Dataset	4.2.1
SOTA	State-of-the-art	
TPU	Tensor processing unit	4.2.2

Chapter 1

Introduction

This master thesis is the final report of the graduation project for the Computer Science and Engineering master at Eindhoven University of Technology (TU/e). The project is carried out at Devhouse Spindle in Groningen. In this chapter the research questions are presented. Section 1.1 explains the context for the research problem. The problem and the research questions are discussed in Section 1.2. Resulting goals from these questions are listed in Section 1.3 as is an outline for the rest of the thesis.

1.1 Thesis context

Spindle is interested in automatically responding to customer questions. The field working on text interpretation is natural language processing (NLP). It is expected that automated responses to text are feasible by recent examples in artificial intelligence. Smartphones include speech recognition, allowing users to control the device using speech [64]. For example, users can obtain information about the weather or the age of a specific president by talking to the device. Self-driving car technology created by Waymo “has driven over 10 million miles of real-world roads” [98]. At the same time Google AlphaGo has learnt how to beat the best Go players in the world [37] and another Google team is working on Duplex [62]. The goal of Duplex is to fill in missing information from Google sites by calling people and asking for the information. Demonstrations by Google show that it is difficult for unsuspecting humans to tell the difference between Duplex and a real human.

Remarkable about the context is the speed in which the field evolves. Picking up a book as recent as 2010 will not list many common practises applied nowadays. The newer NLP approaches use deep learning (to do NLP from scratch) as introduced by Collobert et al. [21] in 2011. In 2013 vector representations for words became dense by the introduction of word2vec [70]. These dense vectors have been “producing superior results on various NLP tasks” [105].

The fast pace of the field and the popularity of machine learning causes this thesis to be in certain aspects unconventional. The benefit of deep learning is that everyone can reproduce and improve results given some time and a (cloud) computer. The fast pace and machine learning cause an atypical high amount of references to respectively arXiv publications and blogposts and websites. It is tried to regard these sources with more than usual skepticism.

1.2 Problem description

Spindle would want to see a system that is automatically trained on various sets of a few dozen text documents to answer customer questions. Research in NLP has focused their efforts on various tasks. To solve the problem defined above, we require a task which is able to interpret questions in real-time using little training data. Natural language understanding (NLU), which maps text to its meaning [52], is working on interpretation of text. Specifically, chatbots are using intent

classification to understand the intention of a user when the user utters some sentence [8, 13, 108]. The IBM sales department claims that Autodesk using chatbots cut down their resolution time “from 1.5 days to 5.4 minutes for most inquiries” [46].

Various parties run benchmarks and use this to draw conclusions about the best performing system. Issues can be pointed out which question the validity of these conclusions. A methodology and three datasets for benchmarking intent classification and named-entity recognition (NER) are created and published by Braun et al. [10]. NER aims to detect information like dates, locations and person names from text. The paper compares accuracy for Microsoft LUIS [69], IBM Watson Conversation (now Watson Assistant [3]), Api.ai (now Google DialogFlow [39]), Wit.ai [101], Amazon Lex [63] and Rasa [8]. When knowing that the field is rapidly advancing [105] it becomes clear that the scores from this paper are outdated. Snips [85] show that they outperform the competition by a large margin [23]. The competition consists of Api.ai, Wit.ai, Luis.ai and Amazon Alexa (now Amazon Lex [63]). Their small benchmark tests the systems against 70 queries per intent on their own dataset. Snips claim to score 79% accuracy, while the second best scores 73%. Also, via sentence examples Snips show that some named-entities are classified incorrectly by systems other than Snips. Although the authors “guarantee transparency” about the benchmark [24], the dataset could still be cherry-picked. DeepPavlov [13] reports another high score for intent classification. It is based on the Snips dataset [24] and compared against Api.ai, Watson Conversation, Microsoft LUIS, Wit.ai, Snips, Recast.ai (now SAP Conversational AI [1]) and Amazon Lex. Their model uses embeddings trained on the DSTC2 dataset [6, 5]. DSTC2 contains communications with users or ‘calls’ [42]. The dataset includes roughly 500 dialogs with 7.88 turns on average in each condition for 6 conditions [42], hence about 20.000 turns or utterances. Knowing that the focus for Snips also lies in interpretation of voice commands [85] it is expected that the model created by DeepPavlov does not obtain state-of-the-art results for other datasets. Botfuel [9] claims to be ‘on par’ with the competition [92]. This is based on runs on the same datasets as Braun et al. [10]. Botfuel shows it is one percent lower than Watson, equals LUIS and outperforms DialogFlow, Rasa, Snips and Recast.ai. The score for Rasa matches the score listed by Braun et al. [10]. This means that Botfuel has compared their system against an old version of Rasa. These observations give rise to the following research question.

RQ1. Can an open-source benchmarking tool for NLU systems and services be created?

An interesting problem from an academic point of view is increasing accuracy. The second research question aims to do that.

RQ2. Can the classification accuracy for NLU systems and services be increased?

For the graduation company Dutch datasets would match their use-case, however the focus in NLP research is on English datasets [14, 105]. To be able to compare our results this thesis will also focus on English datasets. It is expected that the knowledge from answering this question can be transferred to Dutch datasets since modern multilingual models exist [88, 89, 32]. The first and second research question are discussed respectively in Chapter 3 and 4.

1.3 Project goal and outline

Answering RQ1 means writing code. Even if the answer is negative, it will have provided a baseline to use when answering RQ2. The aim of the software is to be used by others, so it should be easy to run. Software extensions should also be possible. The goal related to the first research question is as follows.

RG1. Develop an open-source reproducible tool for benchmarking of NLU systems and services.

The goal related to the second research question is stated ambitiously.

RG2. Improve the accuracy of NLU systems and services.

The first research question is discussed in Chapter 3. NLP and NLU are introduced in detail in Chapter 2, specifically in Section 2.1. The second research question is discussed in Chapter 4. To be able to compare the introduced deep learning model with existing models one need to know about the existing models. Well-known models in NLP are explained in Chapter 2, specifically in Section 2.2.

Chapter 2

Preliminaries

Both research questions are related to natural language processing, which is introduced in Section 2.1. Deep learning forms the basis for state-of-the-art systems in the field. Well-known deep learning models and concepts for NLP are introduced in Section 2.2.

2.1 Natural language processing

Natural language processing (NLP) aims to read text or listen to speech and extract meaningful information from it. One deviation from this definition is natural language generation, where text is generated. The difficulty of NLP is that humans for every word activate “a cascade of semantically related concepts, relevant episodes, and sensory experiences” [14]. These activations may also differ per context. For example, the activations for “I want more money.” depend on whether the sentence is uttered by a child or employee. The field is often associated with artificial intelligence (AI). AI is defined as “a system’s ability to correctly interpret external data, to learn from such data, and to use those learnings to achieve specific goals and tasks through flexible adaptation” by [54]. By this definition modern NLP approaches are AI, specifically artificial narrow intelligence [54].

The field is divided in tasks. Some well-known tasks are:

- Translating texts between (human) languages is called **machine translation**.
- **Question answering** which is NLP on question sentences only to point to an answer in text or generate an answer sentence.
- Classifying words or parts of sentences is done by **part-of-speech tagging** (for example, nouns and verbs) and **named-entity recognition** (for example, dates and locations).
- **Optical character recognition** attempts to recognize characters in images and can use NLP knowledge to improve accuracy.
- Finding co-references like ‘house’ and ‘it’ in the sentence “The house is white, and it is located on a hill.” is done using **coreference resolution**.
- NLP is not limited to text, because it includes **speech recognition** which transforms speech to text.
- **Entailment classification** contains examples such as “People formed a line at the end of Pennsylvania Avenue.” which is contained in (logically implied by) “At the other end of Pennsylvania Avenue, people began to line up for a White House tour.” [100]
- Determining whether two sentences have the same meaning is called **semantic text similarity**.

- **Sentence classification** is the broad tasks of classifying a sentence (for example, the sentiment or intention of user)

Intent classification, machine translation and named-entity recognition are discussed further in Section 2.1.2 and 2.1.3.

2.1.1 Language model

To be able to explain more complex language representations using neural network architectures we first take a look at a simple statistical language model. Language models try to capture the grammar of a language. Capturing grammar can be done by using probabilities for sentences or probabilities of upcoming words given a part of a sentence. It is based on the assumption that grammatically correct sentences occur more often than incorrect sentences. The following three tasks and examples demonstrate the usefulness of probabilities.

Task	Example
Spell correction	$P(\text{my car broke}) > P(\text{my car boke})$
Machine translation	$P(\text{the green house}) > P(\text{the house green})$
Speech recognition	$P(\text{the red car}) > P(\text{she read ar})$

In the spell correction example it is much more likely that the third word should be ‘broke’ than ‘boke’. This is used by automatic spell checkers to correct mistakes. In machine translation knowledge about the expected order of words is used to improve translation accuracy. The speech recognition example shows that ambiguity introduced by phonetic similarity can be solved by choosing the most likely phrase.

A simple approach is to use counters to calculate the sentence probabilities. This makes use of the chain rule or ‘general product rule’. Let W denote a sentence, or equivalently, sequence of words. For the probability of the sentence we have $P(W) = P(w_1, w_2, \dots, w_n)$. The probability of the upcoming word w_i is $P(w_i) = P(w_i | w_1, w_2, \dots, w_{i-1})$. Using the chain rule we can, for three variables, state that $P(w_3, w_2, w_1) = P(w_3 | w_2, w_1) \cdot P(w_2 | w_1) \cdot P(w_1)$. This pattern scales to any number of variables. To get the probability for the sentence ‘the car broke’ we rewrite it as follows.

$$P(\text{the car broke}) = P(\text{the}) \cdot P(\text{car} | \text{the}) \cdot P(\text{broke} | \text{the car})$$

Each term can be rewritten to a combination of counters, for example: $P(\text{broke} | \text{the car}) = \text{COUNT}(\text{the car}) / \text{COUNT}(\text{broke})$. This does not scale well. The counters for each word and each pair of words (each combination of two words) are feasible. However, when doing this for combinations of four words and up the number of counters to keep track of becomes too large.

To reduce the number of counters an approximation defined by Markov is used. Markov states that only looking at a fixed number of previous words gives an approximation. Considering only one previous word for our example we get the following.

$$P(\text{broke} | \text{the car}) \approx P(\text{broke} | \text{car})$$

This is called the bigram model. When using this to generate sentences it becomes clear that bigrams do not have enough information. Take the generated sentence “I cannot betray a trust of them.” [59]. Each pair of sequential words is correct, while the sentence as a whole is not. To improve this systems can be created which use n-grams, where a larger n gives a more accurate language model. Although n-grams offer good performance for certain cases they are in practise not able to capture long-distance dependencies in texts [68].

2.1.2 Machine translation

Machine translation became known to the public by introduction of Google Translate in 2006. The system was based on a statistical language model (as explained in Section 2.1.1) and not on a rule-based model. Rule-based means formulating linguistic rules which is “a difficult job and

requires a linguistically trained staff” [90]. In an attempt to visualise the progress made in the field we consider one example translation through time as recorded by Manning and Socher [68]. This ‘one sentence benchmark’ contains one Chinese to English example which is compared with Google Translate output. The correct translation for the example is:

In 1519, six hundred Spaniards landed in Mexico to conquer the Aztec Empire with a population of a few million. They lost two thirds of their soldiers in the first clash.

Google Translate returned the following translations.

2009 1519 600 Spaniards landed in Mexico, millions of people to conquer the Aztec empire, the first two-thirds of soldiers against their loss.

2011 1519 600 Spaniards landed in Mexico, millions of people to conquer the Aztec empire, the initial loss of soldiers, two thirds of their encounters.

2013 1519 600 Spaniards landed in Mexico to conquer the Aztec empire, hundreds of millions of people, the initial confrontation loss of soldiers two-thirds.

2015 1519 600 Spaniards landed in Mexico, millions of people to conquer the Aztec empire, the first two-thirds of the loss of soldiers they clash.

2017 In 1519, 600 Spaniards landed in Mexico, to conquer the millions of people of the Aztec empire, the first confrontation they killed two-thirds.

One important concept in machine translation is alignment. Alignment refers to the fact that words in different languages tend to be located at similar parts in the sentence. Consider the sentences

“well i think if we can make it at eight on both days”

and

“ja ich denke wenn wir das hinkriegen an beiden tagen acht uhr”.

The first five words of the sentences are perfectly aligned. The last five words of the sentences are not. Alignment is visualised in Figure 2.1. Non-perfect alignment can also be observed from the fact that the number of words in English sentence is higher.

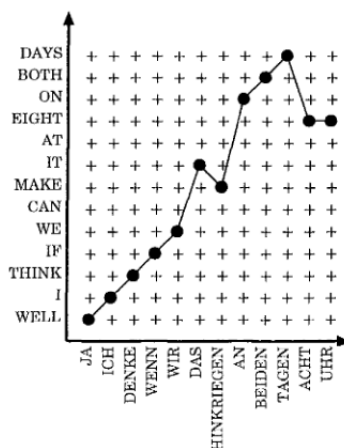


Figure 2.1: Word alignment for a German-English sentence pair [96, Figure 1].

2.1.3 Natural language understanding

Conversational agent or dialogue systems aim to communicate with humans using natural language. Consensus is not clear on whether a chatbot is synonymous to conversational agent. Some argue it is [47], and some argue that conversational agents are more sophisticated than bots which can chat. This thesis will consider the words to be synonymous. Conversational agents can be used to replace graphical user interfaces [11] or to act as a social companion [35, 108]. Interpreting user utterances gives rise to a new task, often denoted as natural language understanding (NLU) [48, 10, 104]. NLU extracts information from user sentences. The meaning for the entire sentence, or the *intention* of the user when uttering a sentence, is defined as the intent. Information from one or more sequential words is defined as named-entities (for example, dates and locations). Consider the following sentence:

“I would like to book a ticket to London tomorrow.”

In this sentence the intent of the user is to book a ticket. Often the chatbot needs to know more than just the intent. For this book ticket example the system needs to know the destination of the user and when the user wants to arrive. Named-entity recognition (NER) can be used to find this information. A named-entity classifier can be trained to classify London as a destination and tomorrow as a date. Most systems allow entities to be defined by examples and regular expressions. The examples can be used for keyword matching by a simple system. More sophisticated systems use machine learning and language knowledge to not only find exact (keyword) matches, but also texts similar to the examples.

2.1.4 F score

A common way to calculate system performance for NLU is the F score, or specifically the F_1 score. It is based on the confusion matrix, see Table 2.1.

Class \ Recognized	as Positive	as Negative
Positive	tp	fn
Negative	fp	tn

Table 2.1: Confusion matrix for binary classification [86, Table 1].

The confusion matrix can be used to define precision π and recall (or sensitivity) ρ as [86].

$$\pi = \frac{tp}{tp + fp}, \quad \rho = \frac{tp}{tp + fn}$$

Then, the F score is defined as [26]

$$F_\beta = \frac{(\beta^2 + 1)\pi\rho}{\beta^2\pi + \rho}$$

where $\beta = 1$ to obtain the F_1 score.

The metric has evolved to be used for three different averages, namely micro, macro and weighted.

- Let y be the set of *predicted* (sample, label) pairs,
- \hat{y} be the set of *true* or *gold standard* (sample, label) pairs,
- L the set of labels,
- y_l the subset of y with label l , or formally: $y_l := \{ (s, l') \in y \mid l' = l \}$,
- $P(A, B) := \frac{|A \cap B|}{|A|}$,

- $R(A, B) := \frac{|A \cap B|}{|B|}$ (where $R(A, B) := 0$ and $P(A, B) := 0$ for $B = \emptyset$), and
- $F_\beta(A, B) := (1 + \beta^2) \frac{P(A, B) \times R(A, B)}{\beta^2 P(A, B) + R(A, B)}.$

Then the metrics are defined as [84]:

$$\begin{aligned} F_{\beta-micro} &= F_\beta(y, \hat{y}) \\ F_{\beta-macro} &= \frac{1}{|L|} \sum_{l \in L} F_\beta(y_l, \hat{y}_l) \\ F_{\beta-weighted} &= \frac{1}{\sum_{l \in L} |\hat{y}_l|} \sum_{l \in L} |\hat{y}_l| F_\beta(y_l, \hat{y}_l) \end{aligned}$$

In this thesis the implementation by Scikit-learn version 0.20.0 [82] is used where possible.

2.2 Deep learning

As with many computer science subfields deep learning has outperformed manual feature engineering on many NLP tasks. For the last few years the best performing NLP systems have been neural networks [105]. This section will provide a basic overview of the most important model architectures for NLP.

2.2.1 Vanishing gradient problem

The vanishing gradient problem, as recognized by Hochreiter et al. [45], is one of the main issues for training neural networks. The backpropagation algorithm updates the weights in the network by sending information from the output layer in the direction of the input layer. The problem relates to vanishing and exploding updates to weights in the layers further from the output layer. These extreme updates are caused by the backpropagation algorithm. Consider some weight w which is near the input layer of some network. This weight is updated by the following partial derivative $w = d_1 \cdot d_2 \cdot \dots \cdot d_i$. Here i corresponds to the number of layers in the network. These partial derivatives d_x can become small ($0 \leq d_x < 1$) or large ($1 \ll d_x$). Typically the learning rate has an order of magnitude of 1e-3. When one derivative is small the weight update when multiplied by learning rate might become very small. Conversely when one derivative is big the update might become very big. Since the weight is near the input layer, i is big. Thus, the chance that weights further away from the output layer vanish or explode increase by increment of i . Exploding gradients can be solved by putting a threshold on the update [70]. Vanishing gradients are more difficult since it is unknown whether an update is small due to vanishing gradients or due to the fact that the weight should not be changed according to the loss function. Effectively the vanishing gradient problem causes layers which are have a long distance from the output layer to stop learning.

2.2.2 Recurrent neural networks

A recurrent neural network (RNN) is an extension on neural networks which allows the network to have memory. This is effective for problems where the data is sequential. In a RNN the information from the state of the network is passed to the next state, see Figure 2.2. This state contains information from previous states, we call this a ‘summary’, denoted W in the image. On the right side the model is unfolded in time. The unfolded representation shows the states of the network and the flow of information for three consecutive points in time. To know its history the neural network in current state s_t obtains the summary W from the previous state s_{t-1} . Suppose we are training a RNN in an unsupervised way. Hence, the model trains on unlabeled real-world

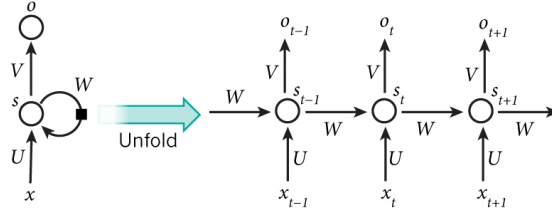


Figure 2.2: Recurrent neural network unfolded in time [61, Figure 5].

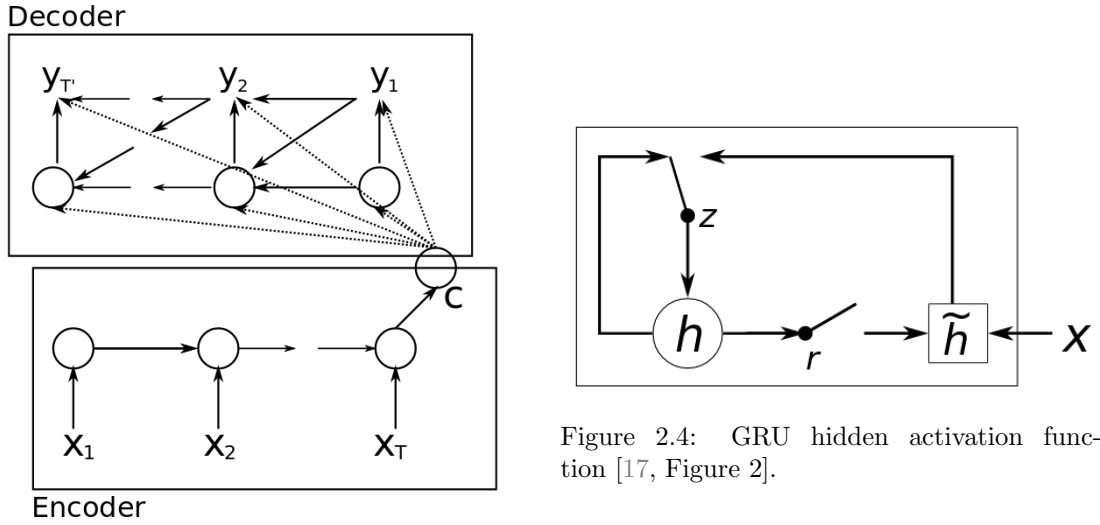


Figure 2.3: RNN Encoder-decoder [17, Figure 1].

texts. For each step it is asked to predict the output word o_t based only on all k previous words $x_{t-1}, x_{t-2}, \dots, x_{t-k}$. The prediction is compared to the correct word, if these are not equal the loss is backpropagated. The backpropagation is then able to ‘change history’ to improve prediction accuracy.

The benefit of this architecture over n-grams, as presented in Section 2.1.1, is that the information is compressed inside the neural network. Also, there is a certain sense of importance since the weights are not uniformly updated for all previous states (words). Take, for example, ‘the car, which is green, broke’. For this sentence the word ‘broke’ can more easily be predicted based on ‘the car’ than on ‘which is green’. It is found that RNNs are able to capture this importance [68].

In practice RNNs are not using the complete history. The cause for this are vanishing gradients. “In practise gradients vanish and simple RNNs become like 7-grams for most words” [68].

2.2.3 Gated recurrent units and LSTMs

Basic RNNs do not yet provide a way to translate languages. Machine translation requires to convert some sentence from a source language to a target language. To this end a RNN encoder-decoder has been proposed by Cho et al. [17].

Similar to the basic RNN the decoder at some time has access to only the previous state, see figure 2.3. The encoder takes the source sentence of variable length T and maps it to a fixed length vector. The decoder in turn maps the fixed length vector to the target sentence of length T' . To do this the network reads all words x_1, x_2, \dots, x_T until an end of line token is received. At that point the entire sentence is captured in the hidden state C . The decoder starts generating

words $y_1, y_2, \dots, y_{T'}$ based on the previous decoder states and C . This is visualised by the arrows. The authors of the paper recognize that this approach has the same limitations as the basic RNN. Typically words in translation sentence pairs are aligned, as described in Section 2.1.2. For example, when generating the first word y_1 the algorithm mainly needs info from x_1 , but has more recently seen the next words in the sequence $x_2, x_3, \dots, x_{T'}$. Vanishing gradients will cause the network to forget its far history, so based on the 7-grams claim RNNs are only effective for translating sentences shorter than 7 words. To solve this the authors [17] have introduced gates to RNNs. To improve this gated recurrent architectures such as long short-term memory networks (LSTMs) [44] and gated recurrent units (GRUs) [17] have been developed.

GRUs have gates which automatically learn to open and close for some hidden state. This can be visualised by looking at the information which is passed through the states. The information is captured in a matrix. In a RNN the information in the entire matrix is updated in each step. GRUs learn to read and write more selectively, see figure 2.5. For some point in time the update consist of reading a specific subset and writing a specific subset of the matrix. In effect the network learns to look only at the word it needs [68]. For example when translating a sentence from German to English it will look at the verb in German to come up with the verb in English. Writing, in effect, lets the model allocate specific parts in the matrix for specific parts of speech (for example, nouns and verbs).

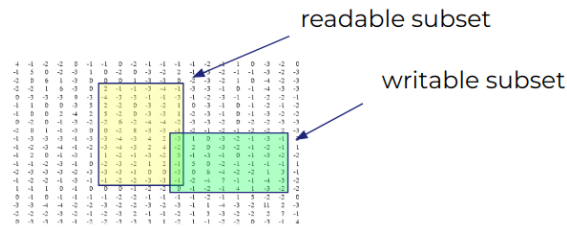


Figure 2.5: Simplistic visualisation for updating the hidden state in a GRU.

LSTMs [44] are similar to GRUs. An LSTM does not only contain update and reset gates but also uses a internal memory state. In practise LSTMs take longer to converge than GRUs [19], but remember around 100 words where GRUs remember around 40 words [68].

2.2.4 Bidirectional recurrent neural networks

All recurrent architectures described above use only the information on the left of some word to predict it. Take for example the following sentences containing the polyseme ‘drinking’.

Peter was drinking after a workout.

Peter was drinking excessively.

The meaning of the word ‘drinking’ changes after reading the next words in the sentence. To take this into account bidirectional recurrent neural networks (BRNN) have been developed by Schuster and Paliwal [80]. A BRNN contains two separate RNNs as depicted in figure 2.6. The paper only considers RNNs, but the method can be applied to gated recurrent models as well [76]. One RNN goes through elements of the sequence from left to right and the other in the reverse direction. Training can be done by showing the network all words except for one. Both networks learn to predict the next word given a sequence of words. Calculating the loss is done by taking the average of the predictions of both RNNs. To reduce required computational power one simplification is used. Suppose we want to learn from the word at location k , x_k , and stepped through many states to reach s_{k-1} and s'_{k+1} . Here we let the model make a prediction for y_k . Then we update the weights and, assuming we go forward, now want to learn from x_{k+1} . The RNN in state s_{k-1} takes one step forward, but the RNN in state s'_{k+1} has to restart from the last word in the sequence. To

solve this both RNNs make one prediction for each word and the answers of both for the entire sequence are batched and used to update the weights.

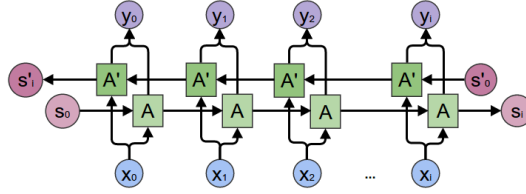


Figure 2.6: Bidirectional RNN [74].

2.2.5 Convolutional neural networks

Convolutional neural networks (CNNs) [60] are an extension to vanilla neural networks which allow the model to exploit spatial locality. This makes the networks highly successful for image classification. During image classification CNNs learn a hierarchical representation of the input, according to Conneau et al. [22]. The authors claim that such an hierarchical representation would have benefits for NLP over sequential models. One of the benefits is that the depth of the tokens, which can be chosen as words but also characters, varies for the tokens in the sentence. In CNNs this depth would be constant, which mitigates the problem of forgetting tokens seen in a less recent past.

A review by Young et al. [105] argues that CNNs are suited for specific NLP tasks, but when data is scarce they are not effective. Foundational problems with CNNs are that they are unable to model long-distance contextual information and to preserve sequential order [105]. The former implies that CNNs are not well suited for question answering tasks such as SQuAD [78].

2.2.6 ELMo

Word embeddings generated by well-known models such as Word2vec [70] and GloVe [75] do not take context into account when determining word representations. For ELMo “each token is assigned a representation that is a function of the entire input sentence” [76]. This is done by using a bidirectional LSTM. Word embeddings are used only to map words to vector representations. To improve accuracy further, compared to traditional embeddings, the authors advise to use ‘the deep internals of the network’ [76]. These internals can be used for downstream models, also known as transfer learning. For example, the authors show that word sense disambiguation tasks are captured by higher level LSTM states. Part-of-speech tagging or named entity recognition are captured by lower level states. ELMo is not the first system to use context, but was obtaining state-of-the-art empirical results on multiple non-trivial tasks at the time of publication. Another reason for the good results is that the system is character based. Word based systems cannot generate an embedding for a word they have not seen during training (out-of-vocabulary tokens). In character based systems morphological clues can be used to guess the meaning of the out-of-vocabulary words. The system has quickly become very popular. Reasons for this seem to be the high accuracy, that the system generalizes well, and that it is integrated into the AllenNLP open-source NLP library created by Gardner et al. [36].

2.2.7 Transformers

The main issue in the recurrent approaches is that distant information needs to pass through all the intermediate states. In the basic RNN for each state all the information is updated, causing less recent information to gradually disappear. Gated recurrent architectures (GRUs and LSTMs) reduce this problem by being more selective when viewing or changing information. Transformer

networks allow the model to look at previous inputs instead of previous states [95]. For example, suppose we are translating ‘impounded’ in the following sentence pair from the WMT’14 English-German dataset.

The motorcycle was seized and impounded for three months.

Das Motorrad wurde sichergestellt und für drei Monate beschlagnahmt.

Suppose the system has correctly predicted the German sentence up to and including ‘Monate’. The next step is to predict ‘beschlagnahmt’. To do this the system needs mainly information about the word ‘impounded’. Gated recurrent architectures learn to look at the previous state in such a way that the attention is focused on ‘impounded’. This requires the information of the word to not be overwritten during execution.

Transformers evade this overwriting problem by allowing the system to see all d previous words, where d is 1024 for the biggest model. The only thing the transformer then needs to learn is where to focus its attention. The information of all the previous words is stored in an array of word vectors (a tensor). To apply focus to parts of this tensor the model learns to put a mask over the tensor. In the mask hidden items are multiplied by infinity [95]. One drawback of this architecture is the required computational power. Suppose we only need one word from the d previous words. The mask will hide $d - 1$ words. This still requires to multiply the masked word vectors by infinity. Google argues that this is not really an issue since matrix multiplication code is highly optimized and graphic processing units (GPUs) and tensor processing units (TPUs) exist. So, the model can relate any dependency in constant time when the range is lower than d . This in contrast to recurrent layers which are in linear time. When the sequence length is greater than d , computations will require more than linear time [95].

Another benefit of the transformers are that self-attention visualisations can more easily be done than in recurrent architectures. By self-attention the authors refer to attention that is used to generate context aware word representations. An example of a transformer model correctly applying coreference resolution is shown in figure 2.7.

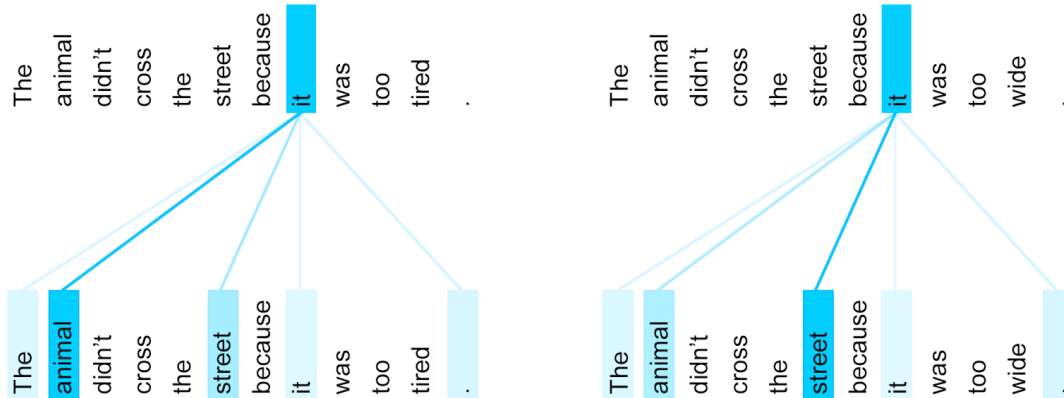


Figure 2.7: Example of transformer encoder self-attention distribution [93].

Chapter 3

Benchmarking

This chapter aims to answer the first research question. The goal is to create a reproducible benchmark tool, simply called BENCH. Datasets used by the benchmark tool are described in Section 3.1. The benchmarked systems are described in Section 3.2. BENCH is described in Section 3.3 as are results obtained by running the benchmarks. It has been observed that BENCH needs further requirements to be more useful, these requirements are described in Section 3.4. Notes on using the tool are presented in Appendix 3.3.

3.1 Datasets

Trained models are considered black boxes at the time of writing. To verify performance of a model data is required. This is fed to the model and results are measured. This section will describe the datasets used for benchmarking.

3.1.1 Format

The dataset format needs to be able to specify a sentence annotation and subsentence annotation. One often used dataset for subsentence, or token, classification is CoNLL-2003 [91]. It uses the NER task definition as described by Chinchor et al. [16]. The definition uses tags to specify entities, for example:

```
<B_ENAMEX TYPE="PERSON">bill<E_ENAMEX> and <B_ENAMEX TYPE="PERSON">
susan jones<E_ENAMEX>
```

Note that this constraints the text since angled brackets (‘<’, ‘>’) cannot be used without escaping the brackets. A less verbose and non text constraining annotation standard is the BIOES-style standard. The origin is unclear, but adaptation is done by at least Stanford as seen in the GloVe paper [75]. Here sentences are annotated as follows.

```
I B-Person
and O
John B-Person
Doe I-Person
worked O
yesterday B-Date
. O
```

where B, I and O respectively mean begin, intermediate and ‘empty’ annotation. Note that other annotations, like part-of-speech tagging, are possible by adding another column of tokens. A

benefit of this annotation is that measuring performance can be done by looking at each token annotation separately. In a tag example like above it is unclear how cases where the classifier is partly correct should be solved. Suppose only ‘susan’ is classified as person and not her last name ‘jones’. Metrics now have to decide how to handle this partially correct situation. In the BIO2 annotation standard token classifications can only be correct or incorrect. One drawback is that the standard is not easy to read for humans. A more readable format is the Rasa Markdown format. Here it constraints the text by using square (‘[’, ‘]’) and round bracket (‘(’, ‘)’) symbols to denote annotations. For example:

[I](person) and [John Doe](person) worked [yesterday](date).

Unlike BIO2 this standard does not easily allow to also specify other annotations, for example part of speech tagging. One could argue that the readability of a Markdown format and the versatility of the BIO standard show that there is no single best approach.

A combination of sentence annotations and token annotations is not supported by the standards described above. For BIO one could track the sentence annotations in a separate file or put it before or after the sentence. The former adds duplicate information while the latter makes the file incompatible with the standard. For Rasa one could change it to a tab separated file and put the Markdown and sentence annotation in separate columns. This is very readable and compact, but means transforming the multi-token annotations to separate token annotations for easier validation. Datasets which combine sentence annotations with token annotations seem to take yet another approach. They use json to store any information they have. These formats should allow for easier parsing, but can not easily be read by humans. For example one dataset annotates entities as follows.

```
"entities": [{
  "entity": "StationDest",
  "start": 4,
  "stop": 4,
  "text": "marienplatz"
}]
```

This entity belongs to the sentence “i want to go marienplatz”. ‘start’ and ‘stop’ here assume the sentence to be tokenized using a WordPunctTokenizer [7] having regexp `\w+|[\^\w\s]+`. Drawbacks are that the entity text is duplicated and that the datasets are hard to read for humans. The sentences have to be manually tokenized for verification and the number of lines of the dataset is an order of magnitude higher than the Markdown format.

3.1.2 Available datasets

Three datasets for intent classification and entity recognition are created and made publicly available by Braun et al. [10]. The paper and its Github version use different names, this thesis will stick to WebApplications, AskUbuntu and Chatbot. WebApplications and AskUbuntu are obtained by pulling questions from StackExchange¹. For example, they respectively contain “How can I delete my [Hunch](WebService) account?” and “How to install a [Brother MFC-5890CN](Printer) network printer?”. Intents for these examples are respectively ‘Delete Account’ and ‘Setup Printer’. StackExchange datasets are labeled using Amazon Mechanical Turk. The Chatbot dataset is based on a Telegram chatbot in production use. This dataset contains sentences like “when is the [next](Criterion) [train](Vehicle) in [munchen freiheit](StationStart)?” having intent ‘Departure-Time’. Labeling for Chatbot is done by the authors of the paper.

Snips [85] is a company which provides software to locally run a voice assistant. They have shared some of the data generated by their users as well as results for their benchmarks [24]. The

¹<https://stackexchange.com>

Dataset	Train	Test	Intents	Entities
WebApplications	30	54	7	1
AskUbuntu	53	109	4	3
Chatbot	100	106	2	5
Snips2017	2100	700	7	unknown

Table 3.1: Number of labeled train and test sentences and unique intents and entities per dataset

incentive for sharing these datasets seems to be showing that their system performs better than other systems. Two datasets have been published by Snips. The thesis has only used the 2017 version and not the 2016 version. The 2017 version will from now on be referred to as Snips2017. Sentences in this dataset are typically short. They utter some command to the system, for example for the intent ‘PlayMusic’: “i want to listen to [Say It Again](track) by [Blackstratblues](artist)”.

Full datasets can be inspected at Github². Summary statistics for these datasets are listed in Table 3.1. In this table ‘None’ is not counted as an intent. The reason for specifying this is that falling back to *null* or some intent during unsure predictions result in different scores for most metrics. F₁ score calculations, for example, do not ignore *nulls* or ‘None’, but instead consider them as a separate group. Information about the unique number of entities for Snips2017 is not specified by the dataset authors.

3.2 Systems

Two open-source systems and some cloud services are considered for the benchmark. The open-source systems are described in Section 3.2.1 and 3.2.2. Cloud services are described in Section 3.2.3.

3.2.1 Rasa

Rasa [8] is an open-source system allowing users to build conversational agents. The system consists of two parts, namely RASA_NLU and RASA_CORE. The former classifies sentences and subsentences. To train the system users can specify (hierarchical) intents, synonyms and regular expressions. Hierarchical intents is a recent addition which allows the system to extract multiple intents from a sentence. For example, it can extract ‘hi+book.ticket’ from “Good morning. Can I order a ticket to London please?”. Regular expressions can, for example, be used to detect numbers and dates in fixed formatting. The system is actively used in production. As a result the code is well documented and stable.

RASA_CORE aims to handle dialogue management. This is an extension on the classifiers of RASA_NLU which aims to understand text in context. Also, it can be used to specify conversation flow. This part remains one of the most difficult problems for conversational agents. Humans tend to switch rapidly between topics in conversations. For example, suppose one ticket order conversation flow contains six questions to be answered by the customer. Customers expect to be able to switch topic during each one of these questions and then return to the flow. Enabling this behaviour via state machines or flowcharts is cumbersome, because the number of transitions grows quickly. One of the Rasa solutions is applying machine learning to let developers train dialog flows interactively.

Rasa can be used via the API or via the ‘Python API’. The Python API is the most efficient. Here users install RASA_NLU in their programming environment and call functions directly. Depending on the used configuration a selection of dependencies have to be installed. The regular API advises use a Docker container. This is less efficient, but more modular and does not require to install dependencies. Containers are published to Docker Hub by Rasa. Users can pull these for free and use the newest stable configuration of choice.

²https://github.com/rikhuijzer/nlu_datasets

Configurations are defined as a pipeline. Pipelines specify what components should be applied to sentences and in what order. Typical pipelines contain at least a tokenizer followed by some classifier. Pipelines are meant to be modified easily and are specified using configuration files. In practise default pipelines often suffice for end-users. A back-end refers to the used intent classifier in some pipeline, for example ‘tensorflow’. Three back-ends are offered by Rasa via Docker Hub, namely RASA-MITIE, RASA-SPACY and RASA-TENSORFLOW. RASA-MITIE is the oldest and depreciated. Training MITIE (<https://github.com/mit-nlp/MITIE>) takes at least a few minutes for small datasets. This is caused by the fact that it is tuning hyperparameters during training. On two different computers used for this thesis the MITIE Docker Hub image occasionally hangs on various datasets. RASA-SPACY is the successor of MITIE and based on spaCy [87]. In 2015 spaCy was in the top 1% for accuracy and the fastest syntactic parser [18]. spaCy (and by that RASA-SPACY) uses a language model to parse text. It includes seven language models for which English, Spanish and French include word vectors. The multilingual model supports only named entities. Unlike the other two back-ends RASA-TENSORFLOW is not based on a pre-trained language model. This is like classifying sentences in an unfamiliar language (say Chinese) after only seeing some examples. Rasa advises to use this back-end when training data contains more than 1000 training examples. The benefit of RASA-TENSORFLOW is that the back-end is language independent and can handle domain specific data and hierarchical intents.

3.2.2 DeepPavlov

DeepPavlov [13] is similar to Rasa. Unlike Rasa, DeepPavlov aims to aid researchers in development of new techniques for conversational agents. Being a newer system than Rasa and aimed at researchers the system is not yet production ready. The system does only provide a Python API, requiring Python 3.6. One claimed benefit of the system is that they do not export machine learning components from other systems. A reason why the system is not well suited for production is that pipelines can download information. This means that a generic Docker needs to download many megabytes of data for each time the Docker is started. Manually defining new Dockers holding this information is possible, but does require some knowledge about Docker and some time to set it up. For users who want to use only a few training examples pre-trained models are necessary. DeepPavlov by default includes DSTC 2, Wikipedia, Reddit, RuWiki+Lenta and ELMo embeddings.

3.2.3 Cloud services

Cloud service providers and various small companies (start-ups) provide APIs for conversational agents. Functionality differs per provider, but in the basis they all offer the same features. Naming conventions do not seem to exist. (For example, Rasa calls intent classification training examples *utterances*, while Google Dialogflow calls them *training phrases*.) Via the web interface or API examples can be sent to a server and the system can be configured. Configurations specify the example utterances, dialog flows, how to classify entities (used for slot filling) and input language. At some point the server will use the provided examples to train the model. This takes a few seconds. An extension on the intent classification and slot filling described above is using knowledge bases. When looking at IBM Watson a document needs to be uploaded and annotated by humans. This is an example of a structured knowledge base.

In the period that IBM Watson [3] won the Jeopardy quiz [43] a lot of math and reasoning was required to create NLP systems. Nowadays, training a competitive neural network for natural language processing is relatively easy. It takes a PhD candidate a few months [68]. This results in a lot of companies providing natural language processing services. An in-depth analysis of all services is left out. The following is a non-exhaustive list. Some offer full conversational agent capabilities, while others focus on natural language understanding.

Watson Assistant (<https://www.ibm.com>) is a conversational agent by IBM.

Dialogflow (<https://dialogflow.com>) is a conversational agent by Google.

Lex (<https://aws.amazon.com/lex>) is a variant created by Amazon.

LUIS (<https://www.luis.ai>) is the conversational agent created by Microsoft.

wit.ai (<https://wit.ai>) can be used for chatbots and is acquired by Facebook. The system is free to use, but Wit is allowed to use the data sent to their servers.

Deep Text (<https://deeptext.ir>) provides sentiment analysis, text classification, named-entity recognition and more.

Lexalytics (<https://www.lexalytics.com>) provides categorization, named entity recognition, sentiment analysis and more.

Pat (<https://pat.ai>) has as goal to humanize AI and provides some conversational agent services.

kore.ai (<https://kore.ai>) focus lies on intent classification and entity extraction with as goal to replace graphical user interfaces with chatbots.

Sixteen more are listed by Dale [25].

3.3 Tool and results

The benchmarking tool is called BENCH and available on Github³. Its goal is to be a reproducible benchmarking tool for intent classification and named-entity recognition. Reproducible means that anyone can clone and run the code to reproduce the results presented in this thesis. The code is written in Python, since it is the default choice for machine learning. Python is conceived as a object-oriented language. Over time it has included more and more functional programming ideas. The code in this project will aim to be adhering to functional programming. Reasons are pedagogic value, improved modularity, expressiveness, ease of testing, and brevity. Some general notes on functional programming in Python are listed in Appendix B.

The functional programming constraints for the project are that we do not define any new classes. Specifically, we do not use the class keyword. Exceptions being NamedTuples and Enums. The code prefers returning iterators over collections, the reason for this is explained in Appendix C. A final remark is about the imports. When importing, an attempt is made to explicitly import using ‘from <module> import <class>’. When more implicit imports are used ‘import <module>’ this is can have multiple causes. It is either caused by the appearance of circular imports, by the fact that some names are too common or to avoid reader confusion. An example for the latter are the types defined in SRC.TYP. The names are quite generic and could cause name clashing or confusion when imported explicitly.

3.3.1 Overview

Since the code does not contain classes, the high-level overview is simply a tree-like structure. This is analogous with a book, where subsections are contained in sections and sections are contained in chapters. In the code small functions are called by larger functions and these larger functions are called by even larger functions. For an overview this idea can be generalized to modules. An overview for the modules of BENCH is roughly as follows. The ‘.py’ suffix is omitted for all elements in the tree. Tests are also omitted.

- BENCH
 - SRC.UTILS
 - SRC.TYP

³<https://github.com/rikhuijzer/bench>

- SRC.DATASET
 - * SRC.DATASETS.CORPORA
 - * SRC.DATASETS.SNIPS
- SRC.SYSTEM
 - * SRC.SYSTEMS.AMAZON_LEX
 - * SRC.SYSTEMS.DEEPPAVLOV
 - * SRC.SYSTEMS.DIALOGFLOW
 - * SRC.SYSTEMS.MOCK
 - * SRC.SYSTEMS.RASA
 - * SRC.SYSTEMS.WATSON
- SRC.EVALUATE
 - * SRC.RESULTS

Some generic functions are listed in SRC.UTILS and used through the entire project. The project makes use of type hints as introduced in Python 2.7. All NamedTuples or ‘types’ are defined in SRC.TYP. An overview of the most important types is depicted in Figure 3.1. These types also contains enumerables or ‘Enums’. These are used in cases where function behaviour depends on some parameter having a fixed set of options. Alternatively one could use strings for these cases depending on user-preference. Notable is the usage of SYSTEM, and by that the usage of CORPUS, in QUERY, SYSTEMCORPUS, CLASSIFICATION and F1SCORE. This is caused by the fact that external systems (for example, DialogFlow) have a certain state which needs to be passed through many functions. This state could be that the system has not yet seen the dataset, resulting in CORPUS.EMPTY, or the system has trained on AskUbuntu, which results in CORPUS.ASKUBUNTU. In, for example, CLASSIFICATION this is used to let some evaluation function know context for an input sentence. This context includes from what dataset the sentence came and what system has classified the sentence.

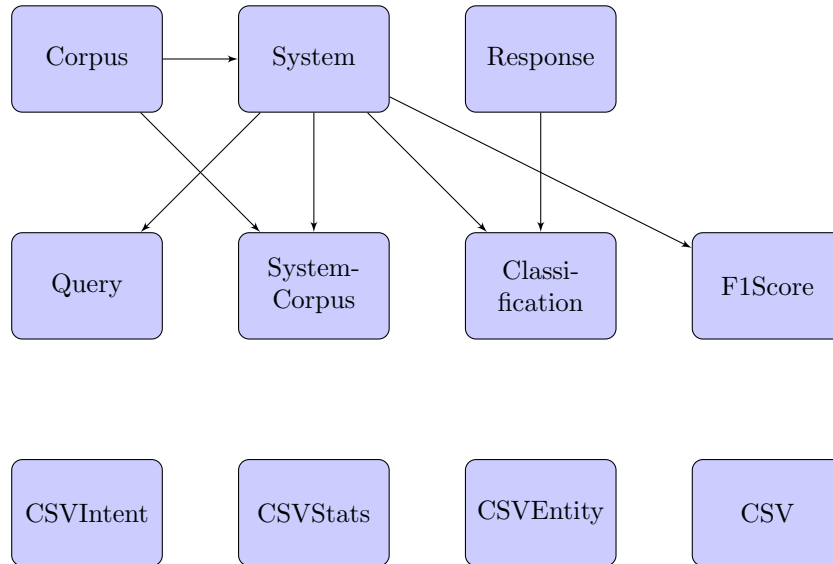


Figure 3.1: Overview of most important type classes (NamedTuples and Enums) and their relations in BENCH. Here a line from A to B means that B is a type which includes A (B extends A).

The real work of the project is done by SRC.DATASET, SRC.SYSTEM and SRC.EVALUATE, as shown in Figure 3.2. ‘Dataset’ takes input files and converts them to an internal representation as defined by SRC.TYP. Input files here denote the original dataset files as created by the

dataset publishers. For the internal representation a Rasa Message⁴ is used. The benefit of this is that it avoids defining the same structure and that it can be used in combination with Rasa code. For example, `SRC.DATASET.CONVERT_MESSAGE_TO_ANNOTATED_STR` uses Rasa code to print the internal data representation as a sentence in Markdown format (Section 3.1.1). Next, the data reaches `SRC.SYSTEM`. Here it is passed to the system under consideration, either in training or prediction mode. For the predictions this is done by finding out which function can convert `SRC.TYP.QUERY` to `SRC.TYP.RESPONSE`. When, for example, Rasa is under consideration the function `SRC.SYSTEMS.RASA.GET_RESPONSE` is called. DeepPavlov would be handled by `SRC.SYSTEMS.DEEPPAVLOV.GET_RESPONSE`. PyCharm is known to have the best type inference for Python. The IDE is not yet able to infer function type for a function mapping, even when all functions have the same input and output type. A workaround is to manually define the type of the function returned by the mapping as `func: Callable[[tp.Query], tp.Response] = ...`. `SRC.EVALUATE` takes all responses `TP.RESPONSE` evaluates the performance of the system under consideration. Printing F_1 score is a matter of three functions and about a dozen lines of code. At one point more advanced logging has been included which is responsible for the other 12 functions and 110 lines of code.

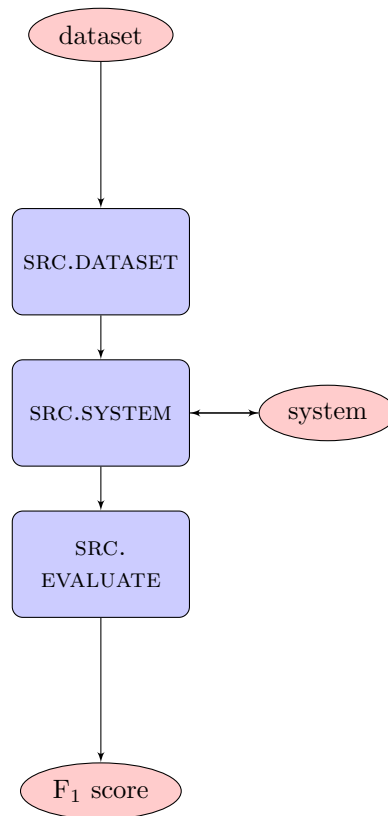


Figure 3.2: Diagram showing the dataflow in BENCH.

3.3.2 Benchmark results

This section presents the benchmark results for intent classification using F_1 scoring with micro averaging. F_1 formulas are listed in Section 2.1.4. An explanation for the differences in averages for the F_1 score is presented in Section 3.4.2. Here micro F_1 scores are used to allow comparing results with Braun et al. [10]. The results are listed in Table 3.2.

⁴`RASA_NLU.TRAINING_DATA.MESSAGE.MESSAGE`

System	Source	AskUbuntu	Chatbot	WebApplications
Rasa:0.5-mitie	see Table D.1	0.862	0.981	0.746
Microsoft LUIS	see Table D.3	0.899	0.981	0.814
Watson Conversation (2017)	see Table D.2	0.917	0.972	0.831
Rasa:0.13.7-mitie	BENCH	0.881		0.763
Rasa:0.13.8-spacy	BENCH	0.853	0.981	0.627
Watson Conversation (2018)	BENCH	0.881	0.934	0.831
Dialogflow	BENCH	0.879	0.986	0.830

Table 3.2: Micro F_1 scores for intent classification. One score is missing due to a bug in BENCH.

The paper remarks that “For our two corpora, LUIS showed the best results, however, the open source alternative RASA could achieve similar results” [10]. When considering only intents this does not hold. Watson Conversation has very similar results, and in fact slightly higher scores on two out of three datasets. The MITIE back-end outperforms the spaCy back-end in terms of accuracy. This would not support the choice of Rasa to depreciate MITIE. It is expected to be caused by the facts that training MITIE takes more time than spaCy and MITIE tends to freeze during training. Interesting to see is that the accuracy for Watson Conversation has dropped. The cause can only be guessed since IBM does not provide information about the Watson back-end. It could be that the calculations for BENCH and the paper differ. Alternatively it could be that the back-end for Watson has changed. The datasets under consideration are small, so it might be that Watson has chosen a back-end better suited for large datasets. Note that IBM is aimed at large companies. These companies have the resources for creating lots of training examples.

3.4 Observations

BENCH requires some further improvements, as explained in Section 3.4.1. Tool design was guided by the methodology as presented by Braun et al. [10]. Section 3.4.2 points out some observations which could improve the the proposed methodology.

3.4.1 Benchmarking system

While developing the benchmarking system, it is observed that the current implementation requires further work. The observations are as follows. Few datasets are publicly available. Current datasets are either small (AskUbuntu, Chatbot, WebApplications) or domain specific (Snips2017). Rasa, for example, uses “a dozen different datasets” [71] for verification. One possible explanation for not sharing these datasets is the sensitive nature of natural language. Another is that they simply not share it to have a competitive advantage. We also observe that dependencies require the product to be continuously maintained. The dependencies are mainly in the form of APIs. APIs are meant to be used by software and will therefore not change often. However, eventually they will change or stop functioning. So, eventually the benchmarking software needs to be updated. Another problem is that the services which offer APIs are not free to use. For each system to be evaluated we need a separate API key. The owner of the benchmarking tool can decide to offer paid keys or let users set keys manually. The former might be possible by negotiations with NLU service providers which are incentivized by selling their product. The latter requires users to have an account for each service. A closed-source solution is Intento⁵. One can send data to the site via their API and they will run a benchmark on various services for a given task. Their ‘catalog’ contains machine translation, intent detection, sentiment analysis, text classification, dictionaries, image tagging, optical character recognition and speech-to-text. The final observation for BENCH is as follows. When choosing a system not only the performance matters. One reason for this is that accuracies are volatile. If companies would base their decision solely on ‘the highest accuracy’

⁵<https://inten.to>

then they would need to change system each month. Since companies can not spend their time constantly switching systems they should also take other factors into account. These factors can include pricing, memory usage, classification speed, privacy (whether open-source) and in-house API preferences.

In summary, we define the following requirements.

- The tool should be continuously maintained,
- offer an API key for each service, or let users add their own keys,
- report more information than just accuracy statistics,
- report evaluation statistics over multiple training runs, and
- include bigger and more varied datasets.

3.4.2 Methodology

Creating a benchmarking tool has resulted in more insight into intent classification. This helped in identifying improvements for the methodology proposed by Braun et al. [10]. As discussed in Section 3.1.2 falling back to ‘None’ or a random intent changes F_1 score. In the paper Chatbot does not have a ‘None’ intent, while WebApplications and AskUbuntu do. Furthermore, drawn conclusions about some system being more accurate than others seems insubstantial. The conclusion that accuracy of some system depends on the domain seems convincing, but is poorly grounded. Reason for this is that both conclusions are based on the F_1 score.

In this paper, the F_1 score is calculated using micro F_1 score. Such a score does not take classes of different size, so called class imbalances, into account. This is combined with a situation where intents and entities are given the same weight. For WebApplications there are in total 74 labeled intents and 151 labeled entities. AskUbuntu contains 128 labeled intents and 123 labeled entities. So, when using micro F_1 on AskUbuntu the score is based somewhat equally on intents and entities. For WebApplications the score is based for about one thirds on intents and two thirds on entities. This could mean that some system has scored significantly better than others simply because it labels entities in WebApplications particularly well. Another observation is that users interested in either intent or entity classification are not well informed. Better seems to be using macro or weighted F_1 and reporting separate intent and entity scores. Macro averaging (and by that weighted averaging) is better suited to “get a sense of effectiveness on small classes” [81]. Rasa, for example, uses the weighted average in their evaluation according to their code on Github.

Another reason for caution with regard to the presented F_1 scores is the probabilistic nature of neural networks. Although inference (classification) is deterministic, training is not. During training models often start with random weights. Random initializations can move into different local minima for the same training data. This could change the inference results. During benchmarking this effect has been observed for Rasa using the spaCy back-end. According to a mail from the main author Rasa 0.5 with the MITIE back-end is used for the results described in the paper. The MITIE back-end has not shown to change accuracy after re-training the model. Microsoft LUIS and Google Dialogflow also did not show a change in accuracy after re-training. So, it could be that all the systems in the benchmark were deterministic. Still, the problem of not considering the possibility of changing accuracy persists for the methodology.

Chapter 4

Improving accuracy

The goal is to improve the classification accuracy for natural language understanding, specifically intent classification. A search is conducted in Section 4.1 to find ways to improve the accuracy. BERT is deemed to be the most promising and is discussed in Section 4.2. The section about BERT describes the model and how it is implemented for this thesis. Accuracy scores for BERT are listed in Section 4.3 and compared to a baseline.

4.1 Search

The field of NLP is rapidly evolving due to the introduction of deep learning [14]. Systems which obtain state of the art (SOTA) accuracy are often surpassed within a few months. A recent example of this is ELMo [76] as published in March 2018 which has been surpassed [105] by BERT [32] in October 2018. A search is conducted to improve the accuracy of the existing systems. The research for this part has not been systematic. The method for finding an improvement is based on coming up with ‘novel’ approaches to improve accuracy. After having such an ‘novel’ approach the literature is consulted. This search method relies on the assumption that papers have done their research and will provide proper related work. This section will explain the considered ideas and related literature.

4.1.1 Duplicate finding

A large part of communications with customers consist of answering questions. Some questions will be duplicates, or in other words, some questions will have been asked and answered before. Finding duplicate questions is the same as finding clusters in the data. Another approach could be based on template responses used by customer support teams. A classifier could be trained to come up with these template responses.

Another approach to find duplicates is using semantic text similarity. It is is a NLP tasks focusing on finding sentences (or texts) having the same meaning and a recurring task in the SemEval workshop. Systems in this field obtain impressive results, however it would not help with the chosen task of intent classification. Also, intent classification and NER extract more information from the sentence than only finding whether there is a similar sentence.

4.1.2 Using data

According to Warden [97] it is more effective to get more training data than to apply better models and algorithms. For conversational agents in company settings it is easy to get raw data. This shifts the problem to applying to automatic data wrangling. Learning automatically from users is applied by some Microsoft chatbots. Microsoft Tay famously started to learn offensive language from users and has been shut down as a result. In China Microsoft has had a more successful release of XiaoIce. XiaoIce is optimized for “long-term user engagement” [108]. Engagement is

achieved by establishing an emotional connection with the user. This system which is created by a research lab and being used by 660 million users does not automatically use the data to learn. The authors manually optimize the engagement of the system. From this it is concluded that reinforcement learning and automatic data wrangling are not yet feasible approaches to increase accuracy.

4.1.3 Kaggle

Kaggle¹ is a well-known site in the machine learning domain. On this site a framework exists where datasets can be published. The site, along other things, shows statistics, a comments section and a scoreboard. It is famous for hosting ‘competitions’, where the person or team obtaining the highest accuracy for some task gets prize money from the dataset hoster. Kaggle provides a way for machine learning enthusiasts to communicate. People who obtain top 20 high scores on difficult tasks tend to explain their pipeline in a blogpost. Research papers tend to focus on designing the best deep learning architectures. The Kaggle explanations are valuable sources for learning how to get the most out of the architectures.

One such post [57] uses three embeddings, namely Glove, FastText and Paragram. The author argues that “there is a good chance that they [the embeddings] capture different type of information from the data”. This method is called boosting. Predictions from the embeddings are combined by taking the average score. A threshold is set to remove answers where the model is unsure. This method could be used to improve performance for natural language understanding. Running three systems in parallel does increase the training time, but the difference is not too large. It would be interesting to test whether averaging can be replaced by a more involved calculation. Meta-algorithms such as boosting, bagging and stacking are not investigated further since the improvement is expected to be insignificant.

4.1.4 Meta-learning

Meta-learning is “is the science of systematically observing how different machine learning approaches perform on a wide range of learning tasks, and then learning from this experience, or meta-data, to learn new tasks much faster than otherwise possible” Vanschoren [94]. Few-shot learning aims to learn useful representations from a few examples. In practise most intent classification systems use few examples, so few-shot learning is interesting to the research question. This was also concluded by the IBM T. J. Watson Research Center [106]. The authors show that their system outperforms other few-shot learning approaches. They do not compare their system against natural language understanding solutions and conclude that their research should be applied to other few-shot learning tasks. This implies that natural language understanding specific systems obtain higher accuracies. Automatically tuning hyperparameters as done in TensorFlow’s AutoML is based on the work by Andrychowicz et al. [2]. Industry claim that AutoML obtains 95% of the accuracy of hand-tuning hyperparameters. Another problem is that it does not scale well [49]. Transfer learning approaches like MAML [34] and Reptile [72] could be useful for intent classification as well. Different domains require different models. Reptile seems interesting to be used to train a model on one domain and then be able to easily switch the model to other domains. This would introduce a lot of complexity in the code. More convenient would be using a model which works on all domains.

4.1.5 Embeddings

Embeddings capture knowledge about language and use that for downstream tasks. There appears to be a consensus about the timeline of embeddings evolution. GloVe [75] was superseded by FastText [50]. The Facebook FastText embedding is aimed to be quick, allowing it to be used as a baseline. With 157 languages (<https://fasttext.cc/>) it is a mutli-lingual model. Another

¹<https://www.kaggle.com>

state-of-the-art and easy to implement embedding is the universal sentence encoder [15]. The word ‘universal’ denotes that the system has used a supervised training task which has been chosen such that the embeddings generalize to downstream tasks. Not only Google, but also Microsoft research is working on multi-task learning [89]. These embeddings are not enough to improve on existing systems, since Rasa is using the universal sentence encoder [99]. One step further would be to let the model decide what embedding it wants to use [56]. A caveat is the fact that one then needs to implement multiple embeddings (even when the model decides that only one embedding should be used).

4.2 BERT

At the start of October 2018 Google published their NLP model called Bidirectional Encoder Representations from Transformers (BERT) [32]. The authors show it is able to score state-of-the-art (SOTA) results for eleven NLP tasks. A comparison by Young et al. [105] shows ELMo [76] outperforms various SOTA models on six distinct non-trivial NLP tasks. The comparison [105] continues by showing that BERT gets higher accuracy scores than ELMo for all six tasks. This by transitivity means that BERT obtains the highest accuracy scores at the time of writing. BERT being SOTA is also supported by a maintained scoreboard for the Stanford Question Answering (SQuAD) dataset [79].

4.2.1 Model description

Results are obtained for a wide range of tasks presented by various datasets. These tasks include entailment classification, semantic text similarity, sentence classification and question answering.

The paper describes three reasons for the good results on the GLUE, MultiNLI and SQuAD datasets. One being that they pre-train the model and let users fine-tune it on their downstream task [32]. Fine-tuning starts with a model for which all the layers are initialized based upon a pre-trained model [41]. Then the output layer is replaced by a task specific output layer, hence the number of labels equals the classes in the downstream dataset [41]. This can also be denoted as transfer learning [20]. The output layer uses dropout during training as described in the CREATE_MODEL function in RUN_CLASSIFIER.PY [30]. Basically, pre-training learns a language model which is used for the downstream task. Another is that transformer models parallelize better than recurrent architectures [95]. This allowed the BERT researchers to train a model having 340 million parameters (BERT_{LARGE}). Lastly, the model is presented as being ‘deeply bidirectional’. The bidirectionality allows the model to use context from both sides to determine the meaning of a word. Deep bidirectionality denotes that the model uses left and right context in all layers of the model. This is visualised and compared to ELMo [76] and OpenAI GPT [77] in Figure 4.1.

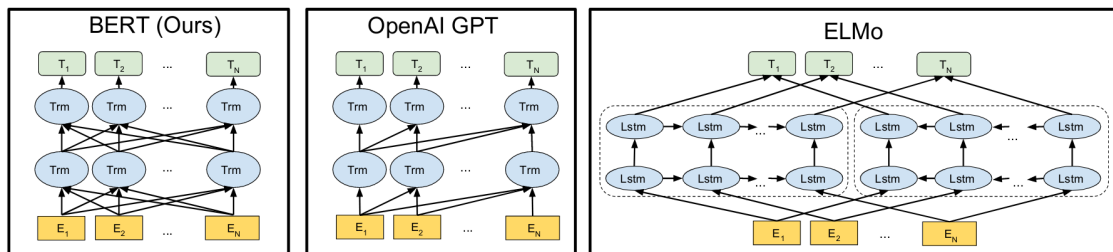


Figure 4.1: Comparison of flow of information in the layers of various recent pre-training model architectures [32, Figure 1]. Note that “only BERT representations are jointly conditioned on both left and right context in all layers” [32].



Figure 4.2: TensorBoard visualisation of ‘summaries’ for fine-tuning pre-trained models on AskUbuntu. Here the horizontal axis denotes the number of steps. Accuracy and loss for the model trained on the test set is evaluated at fixed steps on the evaluation set. The plots show that BERT is in some cases able to converge to a correct local minimum after 100 steps (3200 examples shown to model). The scores in these images should not be compared to the scores in Section 3.3.2 since the metrics and the data split differ.

Another benefit of BERT is that they provide a wide range of pre-trained models. The basic models presented in the BERT paper are BERT_{BASE} and BERT_{LARGE}. BERT_{BASE} is “chosen to have an identical model size as OpenAI GPT for comparison purposes” [32]. BERT_{LARGE} obtains higher accuracy on most tasks and has 340 million parameters in total. Compared to BERT_{BASE} this is an increase from 110 million to 340 million parameters. The BERT Github repository [28] lists some more models, namely uncased and cased variants for BERT_{BASE} and BERT_{LARGE}. In general uncased models suffice, but for certain tasks (for example, NER) performance can be increased by using a cased model [28]. Also, they provide BERT_{MULTILINGUAL} and BERT_{CHINESE}. The multilingual model is trained on the 100 languages having the most Wikipedia pages [31].

4.2.2 Training

From now on training is used to denote fine-tuning of the model. Training the general language model on some downstream task is presented as being inexpensive [28]. Relative to the pre-training it is. Experiments show that fine-tuning with default hyperparameters will run out of RAM on a 16 GB RAM machine. Lowering the batch size reduces the memory usage, but running a few training steps still takes at least a few hours.

To train the model on some tasks it is advised to run “a few epochs” [28]. Based on the example code provided by Google researchers the number of epochs is 3 and the number of training examples is about 1000 [4]. So, it is advised to show the system 3000 examples. For our smaller datasets of around 50 examples this means running $3000/50 = 60$ epochs. When measuring the training time in steps it means running $3000/16 \approx 188$ steps for a batch size of 16. Preliminary experiments on the AskUbuntu dataset (having 53 training examples) with a batch size of 32 confirm this estimate, see Figure 4.2. The images show that the system does not converge smoothly, and can even have a sudden drop in performance. One possible explanation for the performance drop is that the model moved into a non-generalizing local minimum.

The results are interesting because it shows that the model is able to learn something even for a dataset with only tens of training examples. Training the model for 5 steps or 80 examples takes at least a few hours on a modern computer. Interpolation suggests that training 188 steps will take at least 36 hours. This is impractical when doing experiments.

According to the paper the benefit of the Transformer models is that they are highly parallelizable. Training BERT consist mainly of matrix multiplications [27]. These can be done quickly and efficiently on graphic processing units (GPUs) and tensor processing units (TPUs). The latter are ASICs created by Google specifically to do machine learning inference [51] and contain 64 GB of RAM [28]. When using the TensorFlow implementation of BERT, GPUs with 16 GB of RAM are required [28]. GPU optimizations are available in the PyTorch implementation provided by Wolf et al. [102], but PyTorch does not support TPUs at the time of writing. Prices for these GPUs are at least a few thousand euros, which means most users and companies resort to cloud services. Google Colab Google [38] provides free access to a GPU and TPU instance. Code which uses Google Colab for BERT is based on an example implementation provided by Google [4].

Using Colab is a compromise between usability and costs. The costs are limited to the use of some storage in a Google Cloud Bucket. Usability is hindered by the usual constraints of online Jupyter Notebook editors, for example no unittests, no autocomple and poor Git integration. To overcome these issues most of the code is written and tested locally and pushed to a Github repository called IMPROV, see Appendix E. In the Colab the code is then pulled from the repository and main functions are called. Using Colab has benefits as well. Hyperparameters and output are visible in one document and can easily be modified in the Notebook, this eases verification. Reproducibility is possible by opening the Notebook and running all cells. The first cell will ask to link the Colab to a Google account, make sure this account has access to a Google Cloud Bucket.

The plots in Figure 4.2 are created using the default TensorFlow visualisation tool TensorBoard. Generating these plots can be done by specifying a model and metrics using the TensorFlow Estimator API. The plots will not be generated for the rest of the runs for reasons explained in Section E.2. For the rest of this document all results are for the BERT_{LARGE} model since BERT_{BASE} is only created for a fair comparison with OpenAI GPT [32].

4.2.3 Joint training

One reason why neural networks are obtaining the best results for many fields is because networks are now deep. Deep networks have more layers and can therefore learn more complex tasks. One application of this is adding a larger portion of the pipeline to the model. For example, the code by Keller and Bocklisch [55] for the default pipeline for RASA-SPACY, as introduced in Section 3.2.1, contains the following steps. In these steps a featurizer denotes a system component which transforms text to vector representations [12].

1. Tokenization which splits texts up in tokens.
2. Regular expression based intent and entity featurizer (for example able to featurize phone numbers),
3. Intent featurizer based on spaCy [87].
4. Stanford Named Entity Recognizer based on conditional random fields [33].
5. NER synonym detection.
6. Intent classification based on scikit-learn [83].

For this pipeline the Stanford Named Entity Recognizer and scikit-learn classify separately. Preferably one would have one model which could learn to do the entire pipeline, also known as an end-to-end model. End-to-end models have two benefits. Firstly, an end-to-end model avoids feature engineering and data pre-processing [67]. Secondly, end-to-end models can obtain higher accuracies because (semi-)optimal features are found automatically.

That the combination improves independent models has been shown by Ma et al. [66] and Zhang et al. [107]. The results for the former are obtained by using a LSTM network. The latter introduces an algorithm to combine hidden states from an LSTM. They show this for the more general problem of sequential labeling and classification. Intuitively the improvement was to be

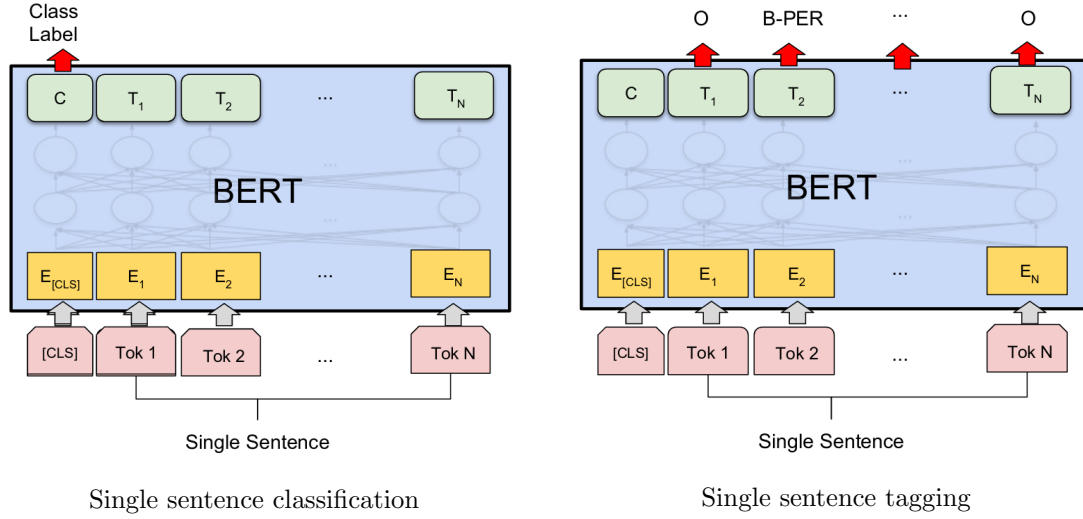


Figure 4.3: Two of the four single sentence tasks presented in the BERT publication [32, Figure 3].

expected for the following reason. Suppose we are trying to classify a dataset which contains the sentence:

“I would like to book a ticket to London tomorrow.”

The sentence has intent ‘BookFlight’. Training the model could be simplified by providing the sentence classifier with:

“I would like to book a ticket to <location> <date>.”

Now the model does not have to learn to classify sentences while also learning that London is a location and that tomorrow is a date.

Note that an end-to-end model is preferred over two separate models. At the time of writing NER classifiers do not obtain perfect accuracy. This means that some classifications will be incorrect. The example from above could instead be converted to:

“I would like to book a <date> to <location> tomorrow.”

This could make the intent classifier drop in accuracy. In an ideal end-to-end model incorrect NER classifications would be less of an issue. The model would learn to ignore the named entity recognition if it would not increase accuracy.

4.2.4 BERT joint training

That joint training BERT is possible can be observed from Figure 4.3.

Let $A = A_1, A_2, \dots, A_n$ denote the layer which is depicted below C, T_1, \dots, T_n , and let $B = B_1, B_2, \dots, B_n$ denote the layer below A . Let s denote the number of tokens for some input sentence. By default the max sequence length for the model is set to 128. For each sentence the sequence length is padded to this max sequence length. When predicting a ‘class label’ C will only be based on A_1 which is based on B_1, B_2, \dots, B_s . A_2, A_3, \dots, A_n are not used. When predicting entities only A_2, A_3, \dots, A_s are used. It seems that a joint model is possible by providing the model with a combination of these two. Consider the NER as a base model and suppose we add some input to C . Now when predicting C the model is expected to learn to look at input from A . For this it can use entity information from A_2, A_3, \dots, A_s . To also learn non-trivial patterns in

non-entity words in the sentence it can use A_2, A_3, \dots, A_n . Typically sentences are much shorter than 128 tokens so enough space should be available in A_2, A_3, \dots, A_n . To allow for more space the max sequence length can be increased, this will increase training and inference time.

To do this the input for the model has been changed from:

```
text: ['how', 'do', 'i', 'disable', 'the', 'spam', 'filter', 'in', 'gmail', '?']
true: ['0', '0', '0', '0', '0', '0', '0', '0', '0', 'B-WebService', '0']
```

to

```
text: ['INTENT', 'how', 'do', 'i', 'disable', 'the', 'spam', 'filter', 'in',
      'gmail', '?']
true: ['FilterSpam', '0', '0', '0', '0', '0', '0', '0', '0', '0', 'B-WebService', '0']
```

where ‘text’ is passed to C, T_1, \dots, T_s and ‘true’ to $E_{[CLS]}, E_1, \dots, E_n$ during training. The BERT tokenizer splits words which are not listed in the vocabulary corresponding to a pre-trained model. INTENT is capitalized to force it to be out-of-vocabulary, it is converted by BERT to [UNK]. The goal of this is to avoid overriding the default interpretation the pre-trained model has for ‘intent’ (or any other uncased token we choose).

The NER loss function can be applied to the joint training without change. This is counter-intuitive because input examples have one position for the class C and s positions for the entities T_1, T_2, \dots, T_s . Typically sentences have around 10 tokens after tokenization. So, the loss is based on around one intent position and ten entities. This means that the model will learn entity recognition much quicker than intent classification. For situations where the model is able to learn the entities this is not expected to affect the performance of the intent classification significantly. It is expected that the only significant difference of this change is in the difficulty of the loss function.

4.3 Results

Experiments are conducted on the AskUbuntu, Webapplications, Chatbot and SNIPS2017 dataset as introduced in Section 3.1.2. Comparisons are made for a fixed number of steps (or equivalently epochs). The reason for this is that intermediate results are not easily reported for the BERT model as explained in Section 4.2.2. The number of steps to be used for training is based on a guess on what should be enough. Validating whether the model should have been trained for more epochs can be done by looking at the loss reported in Appendix F for different number of steps. For each dataset various runs for BERT are executed. During one run only intents are shown to the system and accuracy is measured for intents. Another run only shows entities and measures intents. A third run shows the system intents and entities and measures both. These methods are denoted as separate or joint. It is expected that the joint training increases accuracy for both intent and entity classifications. The reason for this is that the model sees more varied data and hence should be able to more easily find a good internal representation of the data. Results are listed in Table 4.1. Further experiments on the near zero score on Snips2017 separate intent are located in Section E.3.

Note that separate training consists of two runs, and hence ran twice as many epochs. This seems fair, since intent or entity improvements which require twice as many training steps are not interesting for expensive models such as BERT. As a baseline Rasa 0.13.8 with the Spacy pipeline is used. At the time of writing only intent classification is implemented using the benchmark code, so entity scores are missing. Omitting scores for other systems has been deliberate. The table is merely meant to support that joint training is feasible. A final remark is that the scores have been rounded to two decimals. The number of epochs is calculated by taking number of training steps times training batch size and dividing by number of training examples. The training time for 600 steps is around 10 minutes, for 6000 steps it is around 20 minutes. The reason for the relatively small increase in training time is that most time is spent on training preparations and

transferring model checkpoints between TPU and cloud storage.

Dataset	Steps	Batch size	Epochs	Method	Intent	Entity
WebApps				Rasa	0.67 ± 0.04	
	600 (twice)	32	640 (twice)	separate	0.72 ± 0.03	0.81 ± 0.01
	600	32	640	joint	0.76 ± 0.07	0.82 ± 0.01
AskUbuntu				Rasa	0.84 ± 0.00	
	600 (twice)	32	362 (twice)	separate	0.82 ± 0.05	0.81 ± 0.01
	600	32	362	joint	0.87 ± 0.01	0.83 ± 0.00
Chatbot				Rasa	0.98 ± 0.00	
	600 (twice)	32	192 (twice)	separate	0.84 ± 0.21	0.76 ± 0.00
	600	32	192	joint	0.98 ± 0.00	0.79 ± 0.00
Snips2017				Rasa	0.99 ± 0.00	
	6000 (twice)	32	91 (twice)	separate	0.04 ± 0.00	0.84 ± 0.00
	6000	32	91	joint	0.98 ± 0.02	0.86 ± 0.00

Table 4.1: Weighted F_1 accuracy scores (mean \pm standard deviation, over three runs) for separate and joint training on four datasets. The NER accuracy calculation is based on the FULLTOKENIZER [29]. Details are listed in Appendix F.

From the results it can be concluded that joint training is feasible. Joint training increases the named-entity recognition accuracy for each dataset. For intent classification joint training significantly increases accuracy compared to separate training. Compared to the RASA-SPACY baseline it performs better or similar on all the small datasets (WebApplications, AskUbuntu, Chatbot). The accuracy is slightly lower for Snips2017. The model accuracy will be near zero when training on intents separately for SNIPS2017 and the model accuracy will vary for Chatbot. For Snips2017 it has been found that lowering the step size does not solve this problem, as listed in Table F.4. One reason for the poor performance on intent classification only is expected to be that the model get stuck in an incorrect local minimum. This might be explained by the fact that joint training examples are more varied. A typical sentence contains 12 tokens. Then a joint training batch of size 32 will contain about 3 tokens related to intents and 29 related to entities. For an intent training batch of size 32 it will contain 32 tokens related to intents. Hence, the data for the joint training is much more complex. This seems to indicate that the joint training forces the model to learn a more complex representation.

An important thing to note about the results is that the datasets are very small. One would expect that the large BERT model is better suited for datasets which contain more training examples. Furthermore, the experiments are based on a basic implementation. For intent the model as defined by Devlin [28] is used. For named-entity recognition and joint training the model by Kaiyinzhou [53] is used. Not only the batch size but also other hyperparameters can be tuned for better results. Training has used a fixed number of steps or epochs. It might be that more epochs give higher accuracies. On the other hand it might also be that less epochs correspond to similar accuracies in less training time. Other interesting hyperparameters are *max_seq_length* and *learning_rate*. Lowering the former to the expected maximum number of tokens in sentences reduces training and inference time.

Observe that joint training generalizes to sequential labeling and sentence classification. In other words, any combination of tasks where sentences are classified as well as parts of the sentence and these tasks are correlated. The tasks, of course, are expected to be correlated for any real-world NLP dataset. Also, it seems that joint training can be used for multiple sequential labeling tasks (for example, NER and part-of-speech tagging) combined with a sentence classification task. Further work is needed to validate this.

4.4 Implementation improvements

The current implementation is a proof-of-concept. The loss function for joint training needs to be reconsidered. Currently the importance of intent and all named-entities is the same, meaning one error in a intent or in a named-entity weights the same. It might be better for the model to have a higher importance for the intents, since the intent contains information about the entire sentence. Further tests are needed to see whether the increase in accuracy is significant. Comparing NER with a baseline is currently not possible using BENCH since the F_1 score calculation is based on the tokenizer used by BERT. The BERT tokenizer (FULLTOKENIZER) is based on the vocabulary file corresponding to the pre-trained BERT model. So, when trying to compare the score with, for example, Rasa one has to use the same tokenizer and vocabulary file to calculate the Rasa score. Better would be to use a more standard tokenizer as provided by NLTK [7]. Another thing needed to determine the significance of the accuracy improvements is more and more varied datasets. No dataset of around a few hundred distinct training examples has been evaluated, while this seems realistic for production use-cases. One could merge AskUbuntu, WebApplications and Chatbot, but this would be an unrealistic dataset. It would be a dataset where two thirds is related to software questions and one third is related to public transport information. Tests should also be conducted using the other models such as BERT_{MULTILINGUAL} and BERT_{BASE} (cased). The authors [32] use English models to get SOTA results, this suggests that BERT_{MULTILINGUAL} performs worse. It is unknown how big the difference between the two is. Another improvement would be setting up fine-tuning such that TensorBoard visualisations (like Figure 4.2) are possible. This allows users to get a better insight into the model training. Lastly, the code is in need of refactoring and validation. For example, the NER code has a distinction between X and O. O is used by the BIO annotation standard and X is used for padding (since the sentence is shorter than the maximum sequence length). The model in its current implementation does use X to classify tokens. Accuracy can be improved by simply converting all X in the output to O.

Less important improvements are using conditional code to fix I-ENTITY statements appearing without a B-ENTITY. This constraint is enforced by the BIO annotation standard, where ‘intermediate’ statements are preceded by ‘begin’ statements. Also, for production use it might be beneficial to do some hyperparameter tuning.

Chapter 5

Conclusions

Recent well-known artificial intelligence (AI) examples, such as Apple Siri, suggest that difficult natural language processing (NLP) tasks can be automated. The graduation company is interested to see whether this can be used to automate customer support. To this end various NLP tasks have been considered. Eventually it is decided that intent classification and named-entity recognition are an interesting combination of tasks for the graduation company. This combination of tasks is called natural language understanding (NLU) and often used in chatbots to respond to users in real-time. While reading about this task it was found that various parties run benchmarks and draw conclusions. For each party an issue which affects the validity is found. This gives rise to the following research question and goal.

RQ1. Can an open-source benchmarking tool for NLU systems and services be created?

RG1. Develop an open-source reproducible tool for benchmarking of NLU systems and services.

The answer for the first research question is that it is possible. However, stronger requirements are needed to make the tool more useful. The tool is required to be continuously maintained to adapt to changing APIs for NLU services and software. It would require to offer an API key, or let users add their own keys. Furthermore it should offer better metrics. Deciding on a product depends not only on accuracy, but also depends on at least pricing, memory usage and running time. The system should include more metrics to reflect these properties. Currently the accuracy metric is based on the last run. Training is probabilistic, meaning that accuracy may vary for each training run. Reporting only the last run can lead to incorrect conclusions, so the program should execute multiple runs and summarize the results. Lastly, the used datasets are small or domain specific. More and bigger datasets would allow for more statistical powerful conclusions.

Next, the tool (and knowledge obtained by creating the tool) is used to work on the following research question and goal.

RQ2. Can the classification accuracy for NLU systems and services be increased?

RG2. Improve the accuracy of NLU systems and services.

The search for improvement has considered increasing the amount of training data and using new meta-learning algorithms and embeddings. A recent model called Google BERT [32] is expected to be the most likely candidate for increasing accuracy. The model provides a pre-trained checkpoint which has ‘learned’ about language by reading large amounts of text. The pre-trained checkpoint can then be fine-tuned on some specific NLP task using transfer learning. It is a big model, meaning that fine-tuning takes around 1,5 days on a modern computer and a few minutes on a high-end GPU. Experiments on intent classification datasets show non-significant improve-

ments in accuracy. To improve accuracy further the model has been jointly trained on intent classification and named-entity recognition. The benefit is that named-entity information can be used to determine the intent and vice versa. The Google model is a good candidate for jointly training because it, unlike other recent models, uses left and right context in all layers of the network. BERT has obtained state-of-the-art results in a wide range of tasks including named-entity recognition. These two facts imply that jointly training BERT should obtain state-of-the-art results on the joint intent classification and NER task. Basic experiments are conducted in which training BERT separately is compared to training it jointly. The experiments show that jointly training is possible and compared to separate training obtains higher accuracies for intent classification and NER while requiring fewer training steps. Compared to a baseline the intent accuracy is equal or higher for datasets with around 100 training examples. Future work is needed to see whether the model implementation can be improved and whether the improvements in accuracy are significant.

Bibliography

- [1] SAP Conversational AI. Build great bots in minutes. <https://cai.tools.sap>, 2019.
- [2] Marcin Andrychowicz, Misha Denil, Sergio Gomez, Matthew W Hoffman, David Pfau, Tom Schaul, Brendan Shillingford, and Nando De Freitas. Learning to learn by gradient descent by gradient descent. In *Advances in Neural Information Processing Systems*, pages 3981–3989, 2016.
- [3] IBM Watson Assistant. Improve customer and employee experiences with AI. <https://www.ibm.com/cloud/watson-assistant>, 2019.
- [4] Sourabh Bajaj. BERT finetuning tasks in 5 minutes with Cloud TPU. https://colab.research.google.com/github/tensorflow/tpu/blob/master/tools/colab/bert_finetuning_with_cloud_tpus.ipynb, 2018. Accessed: 2018-12-27.
- [5] Dilyara Baymurzina, Aleksey Lymar, and Alexey Sorokin. intents_snips.json. https://github.com/deepmipt/DeepPavlov/blob/0.1.5.1/deeppavlov/configs/classifiers/intents_snips.json, 2019.
- [6] Dilyara Baymurzina, Aleksey Lymar, Mary Trofimova, Yura Kuratov, and Nikolay Bushkov. classifiers.rst. <https://github.com/deepmipt/DeepPavlov/blob/0.1.5.1/docs/components/classifiers.rst>, 2019.
- [7] Steven Bird and Edward Loper. NLTK: the natural language toolkit. In *Proceedings of the ACL 2004 on Interactive poster and demonstration sessions*, page 31. Association for Computational Linguistics, 2004.
- [8] Tom Bocklisch, Joey Faulker, Nick Pawlowski, and Alan Nichol. Rasa: Open source language understanding and dialogue management. *arXiv preprint arXiv:1712.05181*, 2017.
- [9] Botfuel.io. Build enterprise grade chatbots fueled by conversational AI. <https://www.botfuel.io>, 2019.
- [10] Daniel Braun, Adrian Hernandez-Mendez, Florian Matthes, and Manfred Langen. Evaluating natural language understanding services for conversational question answering systems. In *Proceedings of the 18th Annual SIGdial Meeting on Discourse and Dialogue*, pages 174–185, 2017.
- [11] Susan E Brennan. Conversation as direct manipulation: An iconoclastic view. 1990.
- [12] Jake D Brutlag and Christopher Meek. Challenges of the email domain for text classification. In *ICML*, pages 103–110, 2000.
- [13] Mikhail Burtsev, Alexander Seliverstov, Rafael Airapetyan, Mikhail Arkhipov, Dilyara Baymurzina, Nikolay Bushkov, Olga Gureenkova, Taras Khakhulin, Yuri Kuratov, Denis Kuznetsov, et al. DeepPavlov: Open-source library for dialogue systems. *Proceedings of ACL 2018, System Demonstrations*, pages 122–127, 2018.

- [14] Erik Cambria and Bebo White. Jumping NLP curves: A review of natural language processing research. *IEEE Computational intelligence magazine*, 9(2):48–57, 2014.
- [15] Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua, Nicole Limtiaco, Rhomni St John, Noah Constant, Mario Guajardo-Cespedes, Steve Yuan, Chris Tar, et al. Universal sentence encoder. *arXiv preprint arXiv:1803.11175*, 2018.
- [16] Nancy Chinchor, Erica Brown, Lisa Ferro, and Patty Robinson. 1999 named entity recognition task definition. *MITRE and SAIC*, 1999.
- [17] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using RNN encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.
- [18] Jinho D Choi, Joel Tetreault, and Amanda Stent. It depends: Dependency parser comparison using a web-based evaluation tool. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, volume 1, pages 387–396, 2015.
- [19] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*, 2014.
- [20] Dan C Cireşan, Ueli Meier, and Jürgen Schmidhuber. Transfer learning for Latin and Chinese characters with deep neural networks. In *Neural Networks (IJCNN), The 2012 International Joint Conference on*, pages 1–6. IEEE, 2012.
- [21] Ronan Collobert, Jason Weston, Léon Bottou, Michael Karlen, Koray Kavukcuoglu, and Pavel Kuksa. Natural language processing (almost) from scratch. *Journal of Machine Learning Research*, 12(Aug):2493–2537, 2011.
- [22] Alexis Conneau, Holger Schwenk, Loic Barrault, and Yann Lecun. Very deep convolutional networks for natural language processing. *arXiv preprint*, 2016.
- [23] Alice Coucke. Benchmarking natural language understanding systems: Google, Facebook, Microsoft, Amazon, and Snips. <https://medium.com/snips-ai/2b8ddcf9fb19>, 2017. Accessed: 2019-01-18.
- [24] Alice Coucke. 2017-06-custom-intent-engines. <https://github.com/snipsco/nlu-benchmark/tree/master/2017-06-custom-intent-engines>, 2017. Commit ab53441 accessed: 2019-01-18.
- [25] Robert Dale. Text analytics APIs, part 2: The smaller players. *Natural Language Engineering*, 24(5):797–803, 2018.
- [26] Franca Debole and Fabrizio Sebastiani. Supervised term weighting for automated text categorization. In *Text mining and its applications*, pages 81–97. Springer, 2004.
- [27] Tim Dettmers. TPUs vs GPUs for transformers (BERT). <http://timdettmers.com/2018/10/17/tpus-vs-gpus-for-transformers-bert>, 2018. Accessed: 2018-12-27.
- [28] Jacob Devlin. bert - TensorFlow code and pre-trained models for BERT. <https://github.com/google-research/bert>, 2018.
- [29] Jacob Devlin. tokenization.py. <https://github.com/google-research/bert/blob/master/tokenization.py>, 2018.
- [30] Jacob Devlin. run_classifier.py. https://github.com/google-research/bert/blob/master/run_classifier.py, 2019. Accessed: 2019-01-21.

- [31] Jacob Devlin. multilingual.md. <https://github.com/google-research/bert/blob/master/multilingual.md>, 2019. Accessed: 2019-01-21.
- [32] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [33] Jenny Rose Finkel, Trond Grenager, and Christopher Manning. Incorporating non-local information into information extraction systems by Gibbs sampling. In *Proceedings of the 43rd annual meeting on association for computational linguistics*, pages 363–370. Association for Computational Linguistics, 2005.
- [34] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. *arXiv preprint arXiv:1703.03400*, 2017.
- [35] Kathleen Kara Fitzpatrick, Alison Darcy, and Molly Vierhile. Delivering cognitive behavior therapy to young adults with symptoms of depression and anxiety using a fully automated conversational agent (Woebot): a randomized controlled trial. *JMIR mental health*, 4(2), 2017.
- [36] Matt Gardner, Joel Grus, Mark Neumann, Oyvind Tafjord, Pradeep Dasigi, Nelson F. Liu, Matthew Peters, Michael Schmitz, and Luke S. Zettlemoyer. AllenNLP: A deep semantic natural language processing platform. In *ACL workshop for NLP Open Source Software*, 2018.
- [37] Elizabeth Gibney. Google AI algorithm masters ancient game of Go. *Nature News*, 529 (7587):445, 2016.
- [38] Google. Colab. <https://colab.research.google.com>, 2019.
- [39] Google. Dialogflow. <https://dialogflow.com>, 2019.
- [40] Google. Using TPUs. https://www.tensorflow.org/guide/using_tpu, 2019. Accessed: 2019-01-24.
- [41] Yanming Guo, Yu Liu, Ard Oerlemans, Songyang Lao, Song Wu, and Michael S Lew. Deep learning for visual understanding: A review. *Neurocomputing*, 187:27–48, 2016.
- [42] Matthew Henderson, Blaise Thomson, and Jason D Williams. The second dialog state tracking challenge. In *Proceedings of the 15th Annual Meeting of the Special Interest Group on Discourse and Dialogue (SIGDIAL)*, pages 263–272, 2014.
- [43] Rob High. The era of cognitive systems: An inside look at IBM Watson and how it works. *IBM Corporation, Redbooks*, 2012.
- [44] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [45] Sepp Hochreiter, Yoshua Bengio, Paolo Frasconi, and Jürgen Schmidhuber. Gradient flow in recurrent nets: the difficulty of learning long-term dependencies, 2001.
- [46] IBM. Autodesk inc. <https://www.ibm.com/case-studies/autodesk-inc>, 2018. Accessed: 2019-01-16.
- [47] HN Io and CB Lee. Chatbots and conversational agents: A bibliometric analysis. In *Industrial Engineering and Engineering Management (IEEM), 2017 IEEE International Conference on*, pages 215–219. IEEE, 2017.
- [48] Aaron Jaech, Larry Heck, and Mari Ostendorf. Domain adaptation of recurrent neural networks for natural language understanding. *arXiv preprint arXiv:1604.00117*, 2016.

- [49] Llion Jones. Learning to learn by gradient descent by gradient descent. <https://hackernoon.com/learning-to-learn-by-gradient-descent-by-gradient-descent-4da2273d64f2>, 2017. Accessed: 2019-01-14.
- [50] Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. Bag of tricks for efficient text classification. *arXiv preprint arXiv:1607.01759*, 2016.
- [51] Norman P Jouppi, Cliff Young, Nishant Patil, David Patterson, Gaurav Agrawal, Raminder Bajwa, Sarah Bates, Suresh Bhatia, Nan Boden, Al Borchers, et al. In-datacenter performance analysis of a tensor processing unit. In *Computer Architecture (ISCA), 2017 ACM/IEEE 44th Annual International Symposium on*, pages 1–12. IEEE, 2017.
- [52] Dan Jurafsky and James H Martin. *Speech and language processing*, volume 3. Pearson London, 2014.
- [53] Kaiyinzhou. BERT-NER. <https://github.com/kyzhouhau/BERT-NER>, 2018.
- [54] Andreas Kaplan and Michael Haenlein. Siri, Siri, in my hand: Whos the fairest in the land? on the interpretations, illustrations, and implications of artificial intelligence. *Business Horizons*, 62(1):15–25, 2019.
- [55] Caleb M Keller and Tom Bocklisch. config_spacy_duckling.yml. https://github.com/RasaHQ/rasa_nlu/blob/master/sample_configs/config_spacy_duckling.yml, 2018. Accessed: 2018-12-24.
- [56] Douwe Kiela, Chaghan Wang, and Kyunghyun Cho. Dynamic meta-embeddings for improved sentence representations. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1466–1477, 2018.
- [57] Sudalairaj Kumar. A look at different embeddings! <https://www.kaggle.com/sudalairajkumar/a-look-at-different-embeddings/notebook>, 2018. Accessed: 2018-11-11.
- [58] Lak Lakshmanan. How to write a custom Estimator model for the Cloud TPU. <https://medium.com/tensorflow/7d8bd9068c26>, 2018. Accessed: 2019-01-16.
- [59] Irene Langkilde and Kevin Knight. The practical value of n-grams is in generation. *Natural Language Generation*, 1998.
- [60] Yann LeCun, Yoshua Bengio, et al. Convolutional networks for images, speech, and time series. *The handbook of brain theory and neural networks*, 3361(10):1995, 1995.
- [61] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436, 2015.
- [62] Yanviv Leviathan and Yossi Matias. Google Duplex: An AI system for accomplishing real world tasks over the phone. *Google AI Blog*, 2018.
- [63] Amazon Lex. Build conversational bots. <https://aws.amazon.com/lex>, 2019.
- [64] Gustavo López, Luis Quesada, and Luis A Guerrero. Alexa vs. Siri vs. Cortana vs. Google Assistant: a comparison of speech-based natural user interfaces. In *International Conference on Applied Human Factors and Ergonomics*, pages 241–250. Springer, 2017.
- [65] Steven Lott. *Functional Python Programming*. Packt Publishing Ltd, 2015.
- [66] Mingbo Ma, Kai Zhao, Liang Huang, Bing Xiang, and Bowen Zhou. Jointly trained sequential labeling and classification by sparse attention neural networks. *arXiv preprint arXiv:1709.10191*, 2017.

-
- [67] Xuezhe Ma and Eduard Hovy. End-to-end sequence labeling via bi-directional LSTM-CNNs-CRF. *arXiv preprint arXiv:1603.01354*, 2016.
 - [68] Christopher Manning and Richard Socher. Natural language processing with deep learning. Lecture Notes Stanford University School of Engineering, 2017.
 - [69] Microsoft. Language understanding (LUIS). <https://www.luis.ai>, 2019.
 - [70] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119, 2013.
 - [71] Alan Nichol. How to write a custom Estimator model for the Cloud TPU. <https://medium.com/rasa-blog/6daf794efcd8>, 2018. Accessed: 2019-01-25.
 - [72] Alex Nichol and John Schulman. Reptile: a scalable metalearning algorithm. *arXiv preprint arXiv:1803.02999*, 2018.
 - [73] Taiichi Ohno. *Toyota production system: beyond large-scale production*. crc Press, 1988.
 - [74] Chris Olah. colah’s blog. <http://colah.github.io>, 2005. Accessed: 2018-12-10.
 - [75] Jeffrey Pennington, Richard Socher, and Christopher Manning. GloVe: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543, 2014.
 - [76] Matthew E Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep contextualized word representations. *arXiv preprint arXiv:1802.05365*, 2018.
 - [77] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training.
 - [78] Pranav Rajpurkar, Robin Jia, and Percy Liang. Know what you don’t know: Unanswerable questions for SQuAD. *arXiv preprint arXiv:1806.03822*, 2018.
 - [79] Pranav Rajpurkar, Robin Jia, and Percy Liang. The stanford question answering dataset (SQuAD) explorer. <https://rajpurkar.github.io/SQuAD-explorer>, 2019.
 - [80] Mike Schuster and Kuldip K Paliwal. Bidirectional recurrent neural networks. *IEEE Transactions on Signal Processing*, 45(11):2673–2681, 1997.
 - [81] Hinrich Schütze, Christopher Manning, and Prabhakar Raghavan. *Introduction to information retrieval*, volume 39. Cambridge University Press, 2008.
 - [82] Scikit-learn. classification.py. <https://github.com/scikit-learn/scikit-learn/blob/0.20.X/sklearn/metrics/classification.py>, 2018.
 - [83] Scikit-learn. Machine learning in Python. <https://scikit-learn.org>, 2019.
 - [84] Scikit-learn. Precision, recall and F-measures. https://scikit-learn.org/stable/modules/model_evaluation.html#precision-recall-f-measure-metrics, 2019. Accessed: 2019-01-30.
 - [85] SNIPS. Using voice to make technology disappear. <https://snips.ai>, 2018.
 - [86] Marina Sokolova, Nathalie Japkowicz, and Stan Szpakowicz. Beyond accuracy, F-score and ROC: a family of discriminant measures for performance evaluation. In *Australasian joint conference on artificial intelligence*, pages 1015–1021. Springer, 2006.

- [87] spaCy. Industrial-strength natural language processing in Python. <https://spacy.io/>, 2019.
- [88] spaCy. Models overview. <https://spacy.io/models>, 2019. Accessed: 2019-01-16.
- [89] Sandeep Subramanian, Adam Trischler, Yoshua Bengio, and Christopher J Pal. Learning general purpose distributed sentence representations via large scale multi-task learning. *arXiv preprint arXiv:1804.00079*, 2018.
- [90] Eiichiro Sumita and Hitoshi Iida. Experiments and prospects of example-based machine translation. In *Proceedings of the 29th annual meeting on Association for Computational Linguistics*, pages 185–192. Association for Computational Linguistics, 1991.
- [91] Erik Tjong Kim Sang and Fien De Meulder. Introduction to the CoNLL-2003 shared task: Language-independent named entity recognition. In *Proceedings of the seventh conference on Natural language learning at HLT-NAACL 2003-Volume 4*, pages 142–147. Association for Computational Linguistics, 2003.
- [92] Nguyen Trong Canh. Benchmarking intent classification services - June 2018. <https://medium.com/botfuel/eb8684a1e55f>, 2018. Accessed: 2019-01-18.
- [93] Jakob Uszkoreit. Transformer: A novel neural network architecture for language understanding. <https://ai.googleblog.com/2017/08/transformer-novel-neural-network.html>, 2017. Accessed: 2018-12-10.
- [94] Joaquin Vanschoren. Meta-learning: A survey. *arXiv preprint arXiv:1810.03548*, 2018.
- [95] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, pages 5998–6008, 2017.
- [96] Stephan Vogel, Hermann Ney, and Christoph Tillmann. HMM-based word alignment in statistical translation. In *Proceedings of the 16th conference on Computational linguistics-Volume 2*, pages 836–841. Association for Computational Linguistics, 1996.
- [97] Peter Warden. Why you need to improve your training data, and how to do it. <https://petewarden.com/2018/05/28/why-you-need-to-improve-your-training-data-and-how-to-do-it>, 2018. Accessed: 2019-01-12.
- [98] Waymo. Waymo safety report. <https://waymo.com/safety>, 2018. Accessed: 2019-01-16.
- [99] Georg Wiese. Enhancing intent classification with the universal sentence encoder. <https://scalableminds.com/blog/MachineLearning/2018/08/rasa-universal-sentence-encoder>, 2018. Accessed: 2018-12-03.
- [100] Adina Williams, Nikita Nangia, and Samuel R. Bowman. A broad-coverage challenge corpus for sentence understanding through inference. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, 2018.
- [101] Wit.ai. Natural language for developers. <https://wit.ai>, 2019.
- [102] Thomas Wolf, Victor Sanh, Gregory Chatel, and Tim Rault. pytorch-pretrained-BERT. <https://github.com/huggingface/pytorch-pretrained-BERT>, 2018. Accessed: 2018-12-27.
- [103] Tom Wolf, Victor Sanh, and Tim Rault. bert - TensorFlow code and pre-trained models for BERT. <https://github.com/huggingface/pytorch-pretrained-BERT>, 2018.

- [104] Xuesong Yang, Yun-Nung Chen, Dilek Hakkani-Tür, Paul Crook, Xiujun Li, Jianfeng Gao, and Li Deng. End-to-end joint learning of natural language understanding and dialogue manager. In *Acoustics, Speech and Signal Processing (ICASSP), 2017 IEEE International Conference on*, pages 5690–5694. IEEE, 2017.
- [105] Tom Young, Devamanyu Hazarika, Soujanya Poria, and Erik Cambria. Recent trends in deep learning based natural language processing. *ieee Computational intelligence magazine*, 13(3):55–75, 2018.
- [106] Mo Yu, Xiaoxiao Guo, Jinfeng Yi, Shiyu Chang, Saloni Potdar, Yu Cheng, Gerald Tesauro, Haoyu Wang, and Bowen Zhou. Diverse few-shot text classification with multiple metrics. *arXiv preprint arXiv:1805.07513*, 2018.
- [107] Chenwei Zhang, Yaliang Li, Nan Du, Wei Fan, and Philip S Yu. Joint slot filling and intent detection via capsule neural networks. *arXiv preprint arXiv:1812.09471*, 2018.
- [108] Li Zhou, Jianfeng Gao, Di Li, and Heung-Yeung Shum. The design and implementation of XiaoIce, an empathetic social chatbot. *arXiv preprint arXiv:1812.08989*, 2018.

Appendix A

bench usage

Installation is similar to other Python projects. Pull the code and in a terminal set the current working directory to the project folder. Install the required pip packages by running the following command.

```
pip install -r requirements.txt
```

If one wants to check accuracy for an open-source system then run the following command.

```
docker-compose up
```

DOCKER-COMPOSE will read 'docker-compose.yml' and use that information to spin up various Docker containers. All dockers listed in the file are available from Docker Hub. This avoids having to build Dockers manually. DeepPavlov has been removed from the configuration file, since it was found to be unstable, see Section 3.2.2.

After the set-up the program can be executed by running 'bench.py'. To change on which system the benchmarking occurs, replace the first parameter in the GET_SYSTEM_CORPUS call. The prefix is used to determine which system is being tested. Possible prefix options are 'mock', 'rasa', 'deeppavlov', 'lex' and 'dialogflow'. Rasa and DeepPavlov will use the complete string to find a matching port from 'docker-compose.yml'. So, based on the Docker configuration one can also specify 'rasa-tensorflow', 'rasa-spacy' or 'rasa-mitie'. The corpus (dataset) to run the bench on is specified by an enumerable, see SRC.TYP.CORPUS for possible options. When running the script in a modern IDE autocomplete will suggest the possible corpora. Slightly more convenient would be to have the script take input arguments using SYS.ARGV. After setting the two parameters the script can be executed and will display all predictions as well as micro, macro and weighted F_1 scores. The predictions and F_1 scores will also be written to files, see the 'results' folder.

Appendix B

Notes on functional programming in Python

Python is not a pure functional language. However more and more constructs of functional programming are being added to the language each year. This appendix will explain some functional ideas used in the code, as presented by Lott [65]. Higher-order functions take or return functions, this is used to replace the factory design pattern as explained in Section B.1. Keeping track of state without a class results in function signatures to contain many parameters, these can be handled by using NamedTuples, see Section B.2. Another benefit of classes is that data can be stored, used for example in caching. A convenient solution for caching is described in Section B.3. Collections of data are typically transformed via loops. Here each loop will transform the entire collection and move to the next transformation. Lazy evaluation, as described in Section B.4, uses a more efficient way.

B.1 Mapping to functions

In code we often have a function which calls other functions depending on some conditionals. For example in ‘system.py’ the factory design pattern is replaced by a more functional design. In this design ‘system.py’ behaves like a super and delegates the work based on what system we currently interested in. We give an example for two systems. The delegation could be done via conditional statements.

```
if 'mock' in system.name:
    response = src.systems.mock.get_response(tp.Query(system, message.text))
elif 'rasa' in system.name:
    response = src.systems.rasa.get_response(tp.Query(system, message.text))
elif ...
```

This introduces a lot of code duplication. Therefore a dict is created.

```
get_intent_systems = {
    'mock': src.systems.mock.get_response,
    'rasa': src.systems.rasa.get_response,
    ...
}
```

Now we can just get the correct function from the dict and call it.

```
func: Callable[[tp.Query], tp.Response] =
    get_substring_match(get_intent_systems, system.name)
query = tp.Query(system, message.text)
response = func(query)
```


Note that `GET_SUBSTRING_MATCH()` implements the substring matching used in the conditional code (IF 'MOCK' IN `SYSTEM.NAME`:). Since the code can return any of the functions contained in the mapping they should all have the same signature and output. The used IDE (PyCharm 2018.2.4) is not able to check this. Therefore, functions from the mapping `FUNC` get a type hint. This allows the IDE to check types again and it allows developers to see what signature should be used for all the functions in the mapping.

B.2 NamedTuple

Pure functions by definition cannot rely on information stored somewhere in the system. We provide one example from the code where this created a problem and how this can be solved using `NamedTuples`.

The benchmarking tools communicates with a system called Rasa. Rasa starts in a default, untrained, state. To measure its performance we train Rasa and then send many sentences to the system. In general one prefers to functions should be as generic as possible. It makes sense to have one function which takes some sentence, sends it to Rasa to be classified and returns all information from the response. To avoid re-training Rasa for each system we have to remember whether Rasa is already trained. Passing a flag 'retrain' to the system is insufficient, since the function does not know where Rasa should train on. To make it all work we need the following parameters:

- `SENTENCE`: The sentence text.
- `SENTENCE_CORPUS`: The corpus the sentence is taken from.
- `SYSTEM_NAME`: Used to call the function which can train the specific system we are interested in.
- `SYSTEM_KNOWLEDGE`: Used in combination with `SENTENCE_CORPUS` to determine whether we need to re-train.
- `SYSTEM_DATA`: In specific cases even more information is needed.

re-training the system to check whether its outputs differ.

When this function has decided to train the system the `system.knowledge` changes. So as output we need to return “

Since Python 3.5 a `NamedTuple` with type hints is available.

To allow for better type checking and reduce the number of function parameters use is made of `TYPING.NAMEDTUPLES`.

B.3 Function caching

Functions can be cached using `FUNCTOOLS.LRU_CACHE`. This is mainly used for reducing the number of filesystem operations. A typical example is as follows. Suppose we write some text to a file iteratively by calling `WRITE` multiple times. Since we try to avoid storing a state `WRITE` does not know whether the file already exists. To solve we can do two things. The first option is passing parameters telling the function whether the file already exists. This is cumbersome, since this state needs to be passed through all the functions to the function which is calling the loop over `WRITE`. This can be done directly by calling `WRITE` or indirectly by calling some other function. The second option is defining a function to create a file if it does not yet exists `CREATE_FILE`. We call this function every time `WRITE` is called. This does mean that the filesystem is accessed to check the folder each time `WRITE` is called. To avoid all those filesystem operations `CREATE_FILE` can be decorated using `FUNCTOOLS.LRU_CACHE`. Now on all but the first calls to `CREATE_FILE` just query memory.

There is one caveat with this using function caching. Make sure to not try to mimic state. In other words the program should not change behaviour if the cache is removed. Reason for this is that any state introduced via the cache is similar to creating functions with side-effects but without all the constructs from object-oriented programming.

B.4 Lazy evaluation

By default Python is not interested in performance and advises to use a list for every collection. However, lists are mutable and therefore not suitable for hashing. Since hashing is not possible any function taking lists as input is not suitable for function caching.

Also, in many cases the list might not be the final structure we need. Consider the following use cases where the output of type list is used:

- Only unique values are required, so the list is casted to a set.
- Only whether some value satisfies P is required.
- The x first elements are required.
- Only the values satisfying x are required.
- Only an output which is transformed is required.

Considering all these use cases it makes more sense to return an iterator by default instead of a collection. One practical example for the bench project which supports this notion is using an iterator on classification requests.

Suppose we want to measure the performance of some cloud service. Suppose we wrote some code which takes a sentence from some corpus and performs the following operations on this sentence:

1. Send the sentence to some cloud service.
2. Transforms the response to the pieces of information we need.
3. Store this information.

Suppose one of the last two operations contains a mistake causing the program to crash. When not using an iterator all sentences will have been sent to the cloud service after the first operation. Since the post-processing did not succeed we did not obtain results and need to redo this operation. In effect the programming error caused us to waste about as many API calls as there are sentences in the corpus we are testing. This is a problem since the API calls cost money and take time to execute.

To solve this use lazy evaluation. For example, functions supporting lazy evaluation in Python are MAP, FILTER, REDUCE and ANY. Another benefit for using iterators is that it improves modularity and, once used to the paradigm, readability. Take the following typical Python code.

```
my_list = []
for item in some_iterable:
    updated_item = g(f(item))
    my_list.append(updated_item)
```

In this code some iterable is read and the transformation F and G are applied to each item in the iterable. The same code can be rewritten to use MAP as follows.

```
def transform(item: SomeType) -> OtherType:
    return g(f(item))

my_iterable = map(transform, some_iterable)
```

Appendix C

Lazy evaluation demonstration

This appendix demonstrates the effect of using iterators instead of regular collections. The code demonstrates this by processing some fictional raw materials to a chair. The first is function called FORD is similar to a Ford factory around 1915. Here each part of the assembly line just keeps producing items as long as there is input coming in. After a while the other parts of the assembly line start processing the items and discover a fault in the items. One problem of this way of working is that the factory now has a pile of incorrect items in their stock.

The second function called TOYOTA is similar to a Toyota factory after 1960. Here just-in-time (JIT) manufacturing is used as developed by Toyota [73]. Each item is processed only when the next step in the process makes a request for this item.

C.1 Benefits

Using JIT makes sense in computer programs for the following reasons. It saves memory. In each step in the process we only store one intermediate result instead of all intermediate results.

It can detect bugs earlier. Suppose you got a combination of processing steps, lets call them F and G and you apply them to 100 items. In F we send some object to a system and get a response. In H we store the response of this API call. Suppose there is a bug in H, lets say the file name is incorrect. Suppose this is not covered in the tests and we decide to run our program to get all the results we want. Using an approach similar to FORD the program crashes after doing 100 executions of F and G. This means that the program executed 100 API calls. Using TOYOTA the program crashes after just one API call. Here FORD has in essence wasted 99 API calls.

It does not make assumptions for the caller. Suppose some function K returns an iterable and is called by L. The function L can now decide how it wants to use the iterable. For example it can be casted to unique values via SET or it can partly be evaluated by using ANY.

C.2 Code

```
from typing import List, NamedTuple
""" See README.md """

Wood = NamedTuple('Wood', [('id', int)])
Chair = NamedTuple('Chair', [('id', int)])

materials = [Wood(0), Wood(1), Wood(2)]

def ford():
```

```
""" Processing all the items at once and going to the next step. """
def remove_faulty(items: List[Wood]) -> List[Wood]:
    out = []
    for material in items:
        print('inspecting {}'.format(material))
        if material.id != 1:
            out.append(material)
    return out

def process(items: List[Wood]) -> List[Chair]:
    out = []
    for material in items:
        print('processing {}'.format(material))
        out.append(Chair(material.id))
    return out

filtered = remove_faulty(materials)
processed = process(filtered)
print('Result of ford(): {}'.format(processed))

def toyota():
    """ Processing all the items one by one. """
    def is_not_faulty(material: Wood) -> bool:
        print('inspecting {}'.format(material))
        return material.id != 1

    def process(material: Wood) -> Chair:
        print('processing {}'.format(material))
        return Chair(material.id)

    filtered = filter(is_not_faulty, materials)
    processed = list(map(process, filtered))
    print('Result of toyota(): {}'.format(processed))

if __name__ == '__main__':
    ford()
    print()
    toyota()
```

C.3 Output

The output for the program is as follows.

```
inspecting Wood(id=0)
inspecting Wood(id=1)
inspecting Wood(id=2)
processing Wood(id=0)
processing Wood(id=2)
Result of ford(): [Chair(id=0), Chair(id=2)]

inspecting Wood(id=0)
```

```
processing Wood(id=0)
inspecting Wood(id=1)
inspecting Wood(id=2)
processing Wood(id=2)
Result of toyota(): [Chair(id=0), Chair(id=2)]
```

This demonstrates that iterator elements are only executed when called.

Appendix D

Intent F1 score calculations

The F_1 score calculation by [10] uses micro averaging. For two reasons this appendix focuses these calculations on intent only. The first reason is that these results are compared against benchmarks from the BENCH project in Section 3.3.2. Secondly, micro averages could be skewed when the number of intents and entities differ, as described in Section 3.4.2. It is interesting to compare the differences. This appendix lists calculations for Rasa in Table D.1, Watson Conversation in Table D.2 and Microsoft LUIS in Table D.3.

Corpus	Intent	True +	False -	False +	Prec- ision	Recall	F_1 score
Chatbot	DepartureTime	34	1	1	0.971		
	FindConnection	70	1	1	0.986	0.986	0.986
	Σ	104	2	2	0.981	0.981	0.981
WebApps	ChangePassword	4	2	0	1	0.667	0.8
	DeleteAccount	9	1	5	0.643	0.9	0.75
	DownloadVideo	0	0	1	0		
	ExportData	0	3	0		0	
	FilterSpam	13	1	0	1	0.929	0.963
	FindAlternative	15	1	8	0.652	0.938	0.769
	None	0	4	1	0	0	
	SyncAccounts	3	3	0	1	0.5	0.667
	Σ	44	15	15	0.746	0.746	0.746
AskUbuntu	MakeUpdate	34	3	2	0.944	0.919	0.931
	SetupPrinter	13	0	2	0.867	1	0.929
	ShutdownComputer	14	0	6	0.7	1	0.824
	SRecommendation	33	7	4	0.892	0.825	0.857
	None	0	5	1	0	0	
	Σ	94	15	15	0.862	0.862	0.862

Table D.1: Rasa micro F_1 score calculation.

Corpus	Intent	True +	False -	False +	Prec- ision	Recall	F ₁ score
Chatbot	DepartureTime	33	2	1	0.971	0.943	0.957
	FindConnection	70	1	2	0.972	0.986	0.979
	Σ	103	3	3	0.972	0.972	0.972
WebApps	ChangePassword	5	1	0	1	0.833	0.909
	DeleteAccount	9	1	3	0.75	0.9	0.818
	DownloadVideo	0	0	1	0		
	ExportData	2	1	2	0.5	0.667	0.572
	FilterSpam	13	1	2	0.867	0.929	0.897
	FindAlternative	15	1	1	0.938	0.938	0.938
	None	0	4	1	0	0	
	SyncAccounts	5	1	0	1	0.833	0.909
	Σ	49	10	10	0.831	0.831	0.831
AskUbuntu	MakeUpdate	37	0	4	0.902	1	0.948
	SetupPrinter	13	0	1	0.929	1	0.963
	ShutdownComputer	14	0	1	0.929	1	0.963
	SRecommendation	35	5	3	0.921	0.875	0.897
	None	1	4	1	0.5	0.2	0.286
	Σ	100	9	9	0.917	0.917	0.917

Table D.2: IBM Watson Conversation micro F₁ score calculation.

Corpus	Intent	True +	False -	False +	Prec- ision	Recall	F ₁ score
Chatbot	DepartureTime	34	1	1	0.971	0.971	0.971
	FindConnection	70	1	1	0.986	0.986	0.986
	Σ	104	2	2	0.981	0.981	0.981
WebApps	ChangePassword	3	3	0	1	0.5	0.667
	DeleteAccount	8	2	0	1	0.8	0.889
	DownloadVideo	0	0	0			
	ExportData	3	0	1	0.75	1	0.857
	FilterSpam	12	2	0	1	0.857	0.923
	FindAlternative	14	2	2	0.875	0.875	0.875
	None	3	1	8	0.273	0.75	0.4
	SyncAccounts	5	1	0	1	0.833	0.909
	Σ	48	11	11	0.814	0.814	0.814
AskUbuntu	MakeUpdate	36	1	4	0.9	0.973	0.935
	SetupPrinter	12	1	2	0.857	0.923	0.935
	ShutdownComputer	14	0	0	1	1	1
	SRecommendation	36	4	5	0.878	0.9	0.889
	None	0	5	0		0	
	Σ	98	11	11	0.899	0.899	0.899

Table D.3: Microsoft LUIS micro F₁ score calculation.

Appendix E

improv

The project aimed to improve accuracy is called IMPROV and available on Github (<https://github.com/rikhuijzer/improv>). One warning for people interested in reading or using this code is that it is in need of refactoring. The code is cloned from the Google BERT code, built on the Google TensorFlow library, as provided by the researchers [28].

Alternatively, code is available for the PyTorch library [103]. The PyTorch implementation is under active development unlike the TensorFlow implementation and includes more functionality. Features include Multi-GPU, distributed and 16-bits training. These allow the model to be trained more easily on GPUs by reducing and distributing the memory used by the model. BERT contains various models including BERT_{BASE} and BERT_{LARGE}. Since Google Colab does not provide a multi-GPU set-up, we need to use a TPU. This is not yet supported by PyTorch [103].

E.1 Usage

The IMPROV code can partially be executed on a local system. However, training the model requires at least one enterprise grade GPU. This is discussed in Section 4.2.2. GPUs and TPUs are provided for free by Google Colab [38]. Using this code means importing one of IPython Notebooks from the IMPROV repository in Colab. Hyperparameters can be set in the Notebook after which the code can be executed. The Notebook require a Google Account combined with a paid Google Cloud Bucket. The Bucket is used to store the trained checkpoints created by training the model. Newer runs listed in the ‘runs’ folder in the Github repository depend on IMPROV, NLU_DATASETS and RASA_NLU. The dependencies list the used version in the Notebook. When errors occur make sure that the correct versions are cloned or installed.

E.2 TPU and the Estimator API

Visualisation of the training is supported by TensorFlow’s module called TensorBoard. It allows users to define metrics in the model definition which are used by TensorBoard to generate advanced plots. The plots contain not only the metrics (for example, accuracy and loss), but also timestamps for each point and sliders to set area to plot. TensorBoard is included in the default models for TensorFlow and included in the ESTIMATOR API. The ESTIMATOR API is a class which can be used to define a model. One has to create an INPUT_FN and MODEL_FN which respectively define data import and the model. Training is then simply a matter of calling ESTIMATOR.TRAIN(INPUT_FN, ...). Evaluation and prediction are similar.

Unfortunately, this is not supported for TPUs: “TensorBoard summaries [as shown in Figure 4.2] are a great way to inside your model. A minimal set of basic summaries are automatically recorded by the TPUESTIMATOR, to EVENT files in the MODEL_DIR. Custom summaries, however, are currently unsupported when training on a Cloud TPU.” [40].

Two reasons for not supporting TensorBoard summaries in TPUs are likely. Firstly, it seems that the TPUs are mainly used for inference [58]. Secondly, the TPUs are a type of application-specific integrated circuit (ASIC) meaning that TensorBoard summaries are omitted due to technically issues.

E.3 Additional experiment

In Section 4.3 the intent accuracy score for Snips2017 is near zero for runs using batch size 8 and 32. During the run on batch size 8 the loss on the test set was 0.07 at step 1000¹. When comparing this to loss scores in other runs this suggests that the accuracy is not near zero. The loss is in between 1 and 5 for step 2000, 3000, \dots , 6000. This leads to the hypothesis that the model is training too much. Three runs using the same hyperparameters, but reduced to 1000 steps² indicate that the hypothesis is false. It seems that the model found a suitable local minimum by chance.

¹<https://github.com/rikhuijzer/improv/blob/master/runs/snips2017/2018-12-20snipsintentbatchsize8.ipynb>

²<https://github.com/rikhuijzer/improv/tree/master/runs/2019-01-23snips>

Appendix F

Runs

In Section 4.3 the results for the Google BERT runs are summarized. The runs on which the summary is based are listed in this appendix. For WebApplications, AskUbuntu, Chatbot and Snips2017 see respectively Table F.1, F.2, F.3 and F.4.

Method	Batch size	Run	Loss			Intent	Entity
Rasa		1				0.674	
Rasa		2				0.722	
Rasa		3				0.625	
			200 steps	400 steps	600 steps	600 steps	600 steps
intent	32	1	$5.66 \cdot 10^{-5}$	$1.11 \cdot 10^{-4}$	$8.00 \cdot 10^{-5}$	0.682	
intent	32	2	$6.15 \cdot 10^{-5}$	$4.88 \cdot 10^{-5}$	$3.85 \cdot 10^{-5}$	0.742	
intent	32	3	$1.13 \cdot 10^{-4}$	$7.04 \cdot 10^{-5}$	$5.21 \cdot 10^{-5}$	0.734	
ner	32	1	$3.21 \cdot 10^{-4}$	$2.34 \cdot 10^{-4}$	$1.20 \cdot 10^{-4}$		0.815
ner	32	2	$3.31 \cdot 10^{-4}$	$1.76 \cdot 10^{-4}$	$1.31 \cdot 10^{-4}$		0.815
ner	32	3	$4.43 \cdot 10^{-4}$	$1.94 \cdot 10^{-4}$	$1.39 \cdot 10^{-4}$		0.793
joint	32	1	$1.03 \cdot 10^1$	$3.47 \cdot 10^{-4}$	$2.12 \cdot 10^{-4}$	0.668	0.826
joint	32	2	$3.20 \cdot 10^{-3}$	$1.11 \cdot 10^{-3}$	$8.23 \cdot 10^{-4}$	0.761	0.816
joint	32	3	$6.10 \cdot 10^{-4}$	$3.92 \cdot 10^{-4}$	$2.57 \cdot 10^{-4}$	0.839	0.829

Table F.1: Training loss, intent classification score and NER score for WebApplications. Loss is determined using the training set. Both scores are weighted F1 as determined using the evaluation set. Source: <https://github.com/rikhuijzer/improv/tree/master/runs/webapplications>.

Method	Batch size	Run	Loss			Intent	Entity
Rasa		1				0.834	
Rasa		2				0.833	
Rasa		3				0.843	
			200 steps	400 steps	600 steps	600 steps	600 steps
intent	32	1	$3.51 \cdot 10^{-5}$	$6.02 \cdot 10^{-5}$	$1.93 \cdot 10^{-5}$	0.813	
intent	32	2	$2.17 \cdot 10^{-4}$	$8.67 \cdot 10^{-5}$	$1.49 \cdot 10^{-5}$	0.752	
intent	32	3	$2.88 \cdot 10^{-5}$	$3.54 \cdot 10^{-5}$	$3.28 \cdot 10^{-5}$	0.881	
ner	32	1	$1.37 \cdot 10^{-3}$	$2.82 \cdot 10^{-4}$	$2.69 \cdot 10^{-4}$		0.805
ner	32	2	$9.02 \cdot 10^0$	$2.75 \cdot 10^{-4}$	$3.07 \cdot 10^{-4}$		0.802
ner	32	3	$9.17 \cdot 10^{-4}$	$2.16 \cdot 10^{-4}$	$1.74 \cdot 10^{-4}$		0.816
joint	32	1	$7.62 \cdot 10^{-3}$	$8.80 \cdot 10^{-4}$	$7.58 \cdot 10^{-4}$	0.892	0.830
joint	32	2	$1.32 \cdot 10^1$	$4.47 \cdot 10^{-3}$	$2.05 \cdot 10^{-3}$	0.859	0.824
joint	32	3	$1.01 \cdot 10^{-3}$	$3.71 \cdot 10^{-4}$	$1.40 \cdot 10^{-3}$	0.870	0.825

Table F.2: Training loss, intent classification score and NER score for AskUbuntu. Loss is determined using the training set. Both scores are weighted F1 as determined using the evaluation set. Source: <https://github.com/rikhuijzer/improv/tree/master/runs/askubuntu>.

Method	Batch size	Run	Loss			Intent	Entity
Rasa		1				0.981	
Rasa		2				0.981	
Rasa		3				0.981	
			200 steps	400 steps	600 steps	600 steps	600 steps
intent	32	1	$4.65 \cdot 10^{-6}$	$1.10 \cdot 10^{-6}$	$1.61 \cdot 10^{-6}$	0.991	
intent	32	2	$1.30 \cdot 10^{-5}$	$1.36 \cdot 10^{-5}$	$8.11 \cdot 10^{-6}$	0.991	
intent	32	3	$9.38 \cdot 10^{-1}$	$7.96 \cdot 10^{-1}$	$5.97 \cdot 10^{-1}$	0.537	
ner	32	1	$2.01 \cdot 10^{-3}$	$8.56 \cdot 10^{-4}$	$1.45 \cdot 10^0$		0.755
ner	32	2	$1.03 \cdot 10^{-3}$	$3.65 \cdot 10^{-4}$	$3.78 \cdot 10^{-4}$		0.759
ner	32	3	$6.48 \cdot 10^{-4}$	$2.91 \cdot 10^{-4}$	$2.64 \cdot 10^{-4}$		0.760
joint	32	1	$1.46 \cdot 10^{-3}$	$4.12 \cdot 10^{-4}$	$5.07 \cdot 10^{-4}$	0.990	0.791
joint	32	2	$2.25 \cdot 10^{-3}$	$6.31 \cdot 10^{-4}$	$2.93 \cdot 10^{-4}$	0.991	0.787
joint	32	3	$2.92 \cdot 10^0$	$1.25 \cdot 10^{-3}$	$1.06 \cdot 10^{-3}$	0.971	0.783

Table F.3: Training loss, intent classification score and NER score for Chatbot. Loss is determined using the training set. Both scores are weighted F1 as determined using the evaluation set. Source: <https://github.com/rikhuijzer/improv/tree/master/runs/chatbot>.

Method	Batch size	Run	Loss			Intent	Entity
Rasa		1				0.991	
Rasa		2				0.990	
Rasa		3				0.990	
			2000 steps	4000 steps	6000 steps	6000 steps	6000 steps
intent	8	1	$2.07 \cdot 10^{-1}$	$7.80 \cdot 10^{-2}$	$2.06 \cdot 10^0$	0.036	
intent	8	2	$6.14 \cdot 10^{-1}$	$1.65 \cdot 10^0$	$1.92 \cdot 10^0$	0.036	
intent	8	3	$3.27 \cdot 10^0$	$4.60 \cdot 10^{-1}$	$2.37 \cdot 10^0$	0.036	
			500 steps	1000 steps	1500 steps	1500 steps	1500 steps
intent	32	1		$1.93 \cdot 10^1$	$2.25 \cdot 10^1$	0.036	
intent	32	2		$2.05 \cdot 10^0$	$1.98 \cdot 10^0$	0.036	
intent	32	3		$2.12 \cdot 10^0$	$2.05 \cdot 10^0$	0.036	
ner	32	1		$5.75 \cdot 10^0$	$5.95 \cdot 10^0$		0.843
ner	32	2		$1.72 \cdot 10^{-3}$	$1.47 \cdot 10^{-3}$		0.843
ner	32	3		$1.21 \cdot 10^{-3}$	$1.48 \cdot 10^{-3}$		0.840
joint	32	1		$2.47 \cdot 10^{-3}$	$1.17 \cdot 10^0$	0.979	0.859
joint	32	2		$5.50 \cdot 10^0$	$1.17 \cdot 10^0$	0.955	0.861
joint	32	3		$3.64 \cdot 10^{-3}$	$3.08 \cdot 10^{-3}$	0.996	0.861
joint*	32	1	$1.5 \cdot 10^0$	$1.4 \cdot 10^{-3}$	$4.7 \cdot 10^{-1}$	0.996	0.860
joint*	32	2	$3.9 \cdot 10^{-1}$	$4.1 \cdot 10^{-3}$	$3.3 \cdot 10^{-3}$	0.949	0.860
joint*	32	3	$1.7 \cdot 10^1$	$1.5 \cdot 10^{-3}$	$1.6 \cdot 10^{-3}$	0.994	0.863

Table F.4: *Fix: Decapitalized dataset. Training loss, intent classification score and NER score for Snips2017. Loss is determined using the training set. Both scores are weighted F1 as determined using the evaluation set. Source: <https://github.com/rikhuijzer/improv/tree/master/runs/snips2017>.