

# IP9

Daniel Tiefenauer

September 28, 2018

## **Abstract**

tbd.

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Scope and overall goal . . . . .	1
1.2	Chosen approach and previous work . . . . .	1
1.2.1	Previous results . . . . .	1
1.3	Goal of this project . . . . .	2
<b>2</b>	<b>Training an Recurrent Neural Network (RNN) for Automatic Speech Recognition (ASR)</b>	<b>3</b>
2.1	Exploiting the <i>DeepSpeech</i> model . . . . .	3
2.2	A simpler <i>DeepSpeech</i> model . . . . .	3
<b>3</b>	<b>The importance of Language Model (LM)s for ASR</b>	<b>5</b>
<b>4</b>	<b>Plotting a learning curve</b>	<b>7</b>
<b>5</b>	<b>Creating a Language Model for German</b>	<b>10</b>
5.1	n-Gram Language Models . . . . .	10
5.2	Creating a raw text corpus and training the model . . . . .	10
5.3	Evaluating the model . . . . .	11
<b>6</b>	<b>Measuring the performance of the pipeline</b>	<b>13</b>
<b>7</b>	<b>Ehrlichkeitserklärung</b>	<b>17</b>

# 1 Introduction

This report documents the progress of the project *Forced Alignment with a Recurrent Neural Network*. The project serves as a master thesis at University of Applied Sciences (FHNW) (IP9). Some preliminary work has been done in a previous project (IP8). The overall goal, project situation and background information are described in detail in the project report for IP8 and are not repeated here. Instead, a quick recap is given for the sake of completeness of this documentation.

## 1.1 Scope and overall goal

*ReadyLingua* is a Switzerland based company that develops tools and produces content for language learning. Some of this content consists of audio/video data with an accompanying transcript. The overall goal is to enrich this data with temporal information, so that for each part of the transcript the corresponding point in the audio/video data can be found. This process is called *acfa*. An *InnoSuisse* project was started in 2018 to research how this could be achieved. The *InnoSuisse* project foresees three different approaches, one of which is followed in this project.

## 1.2 Chosen approach and previous work

The approach chosen in this project is based on speech pauses, which can be detected using *Voice Activity Detection (VAD)*. The utterances in between are transcribed using *ASR*, for which a *RNN* is used. The resulting partial transcripts contain the desired temporal information and can be matched up with the full transcript with a process called *Local Sequence Alignment (LSA)*.

In the IP8 project, VAD, ASR and LSA were treated as part of a pipeline which split a given audio file into individual utterances, transcribe them and localize them in the original transcript. Since the quality of the ASR stage has an imminent impact on the subsequent LSA stage, the quality of the alignments is heavily dependent on the quality of the partial transcripts. However, ASR is highly prone to external influences like background noise, properties of the speaker (gender, speaking rate, pitch, loudness). Apart from that, language is inherently ambiguous (e.g. accents), inconsistent (e.g. linguistic subtleties like homonyms or homophones) and messy (stuttering, unwanted repetitions, mispronunciation).

### 1.2.1 Previous results

For the VAD stage, an implementation of WebRTC was used which was shown to be capable of detecting utterances with very high accuracy within reasonable time. For the LSA stage a combination of the Smith-Waterman algorithm and the Levenshtein distance was used. This combination included tunable parameters and proved to be able to be able to localize partial transcript within the full transcript pretty well, provided the similarity between actual and predicted text was high enough.

Because the final pipeline should be language-agnostic, the IP8 project proposed the use of *DeepSpeech* for the ASR stage, which uses Connectionist Temporal Classification (CTC) (Graves, Fernández, and Gomez 2006) as its cost function. It included some experiments on what features could be used to train a RNN for the ASR stage. Possible features were raw power-spectrograms (as stipulated by the *DeepSpeech* paper), Mel-Spectrograms and Mel-Frequency Cepstral Coefficients (MFCC). It was found that training on MFCC features would probably require the least amount of training data because. An RNN using a simplified version of the *DeepSpeech* architecture was trained on data from the *LibriSpeech* project (containing only English samples). However, developing a fully-fledged ASR system is extremely time-consuming and could not be done within the project time. For that reason a state-of-the-art *Speech-To-Text (STT)* engine (*Google Cloud Speech*) was embedded in the pipeline as the ASR stage. Using this engine, the pipeline was able to produce very good (although not perfect) transcripts for the individual utterances. Therefore the chosen approach was validated and the pipeline could shown to be generally functional.

### 1.3 Goal of this project

In this project, the chosen pipelined approach shall further be refined. Because the VAD and the LSA stage already work pretty well, the focus in this project is on the ASR stage. A RNN is trained that can be used as for this stage in the pipeline. Because the superordinated goal of this project is Forced Alignment (FA) and not Speech Recognition, this RNN only needs to be *good enough* for the downstream LSA stage. By exploring various combinations of properties of the network or the data as well as varying amounts of training data, the conditions should be researched under which such a network could be trained. Concretely, the following aspects shall be examined more closely:

- **How does the quality of the simplified *DeepSpeech*-RNN change with increasing training data?**  
By plotting the learning curve we should be able to see whether the RNN is able to learn something useful at all and also get some intuition about how much training data is needed to get reasonably accurate partial transcripts.
- **How does the quality of the partial transcripts change when using synthesized training data?**  
Neural Network usually require large amounts of training data and often improve with increasing size of the training set. However, labelled training data is usually scarce and expensive to acquire. For the purpose of Forced Alignment however, synthesized training data can be easily obtained by adding some distortion to the original signal (reverb, change of pitch, change of tempo).
- **How does the quality of the partial transcript change when integrating a LM?** STT-engines traditionally use a LM that models the probabilities of characters, words or sentences. A LM can help producing valid transcripts by mapping transcripts (that may sound similar to what was actually said) to orthographically correct sentences.
- **How can we assess the quality of the alignments?** This should give us some insight about how the quality of the alignment changes with varying capability of the STT-engine and what quality of transcripts is required.

The answers to above questions should help in estimating the effort to create a generic solution (required minimum amount of training data, architecture, etc.). <sup>1</sup>

---

<sup>1</sup>Because ASR is highly dependent on the language that should be recognized, a different STT system has to be trained for each language.

## 2 Training an RNN for ASR

As stated above, this project does not aim at training a state of the art STT engine. Because the Smith Waterman (SW) algorithm used for local alignment is tolerant to a certain amount of errors in the transcripts, the RNN need only be *good enough* for the task at hand (FA). If such a network can be trained under the given circumstances it could be used in the ASR stage of the pipeline. The pipeline would then become self-contained and would not be dependent on a commercial solution that cannot be tuned and whose inner workings are unknown. Furthermore such a RNN would open up to recognizing languages for which there is not a third-party solution yet, such as Swiss German.

### 2.1 Exploiting the *DeepSpeech* model

A Neural Network (NN) that had quite an impact on ASR was *DeepSpeech* (Hannun et al. 2014) which reached recognition rates near-par to human performance, despite using a comparably simpler than traditional speech systems. Because the relation between audio signal and text was learned end-to-end (E2E) the network was also pretty robust to distortions like background noise or speaker variation. An open source implementation of a *DeepSpeech* model is available from Mozilla <sup>2</sup>. Since this implementation uses a LM, the quality of the model is measured as the percentage of misspelled or wrong words (called Word Error Rate (WER)) or as the edit distance (also called Levenshtein distance or Label Error Rate (LER)). A pre-trained model for inference of English transcript can be downloaded, which achieves a WER of just 6.5%, which is close to what a human is able to recognize (Morais 2017).

A model could be trained by providing training-, validation- and test-data for an arbitrary language (e.g. from the *ReadyLingua* corpus). However, this is not the preferred procedure for this project for various reasons:

1. The *DeepSpeech* implementation was designed for ASR. In such settings a low WER is desirable. But this is not the main focus for this project. As a result, the architecture of the Mozilla implementation might be overly complicated for this project, although it might make sense for pure ASR tasks.
2. The problem with above point is that more complex models usually require more training data. However, as for any neural network, the limiting factor for training a RNN is often the lack of enough high quality training data. This becomes especially important when recordings in a minority language should be aligned.
3. The implementation requires an (optional) LM, which is tightly integrated with the training process which might not be available for the target languages.

For these reasons, the RNN architecture of the *DeepSpeech* model was used as a basis for a simplified version, which should (hopefully) require less training data in order to converge and still produce partial transcriptions that can be locally aligned.

### 2.2 A simpler *DeepSpeech* model

An implementation of the RNN used for STT in the previous IP8 project was done in Python using Keras<sup>3</sup>. This model is further referred to as *previous model*. Unfortunately, this model did not perform well, i.e. it was not able to learn how to infer a transcript from a given sequence of feature vectors from a spectrogram. Furthermore, performance was a big issue, although the RNN used a simpler architecture and no computational power was needed to query a LM. Training on aligned speech segments from the *LibriSpeech* corpus was not possible within project time because it would have taken approximately two months when using a single acGPU. This duration is at least consistent with the experience made by the Machine Learning team at Mozilla

---

<sup>2</sup><https://github.com/mozilla/DeepSpeech>

<sup>3</sup><https://keras.io>

Research, which used a cluster of 16 GPUs that required about a week (Morais 2017) to train a variant <sup>4</sup> of the RNN originally proposed in (Graves, Fernández, and Gomez 2006).

For this project, the previous model was closer examined to find out what works best and to help the model converge. After some changes to the previous model architecture for this project, a new model was obtained which was able to learn something meaningful, i.e. it started to infer transcripts that – although still not perfect – resembled the ground truth. This model is further referred to as *new model*. Summarized, the following changes were made to the previous model architecture to constitute the new model:

- **Optimizer:** The new model uses Stochastic Gradient Descent (SGD) instead of Adam. Adam was used in the previous model because it was assumed to be generally superior to other other optimizers in most cases.
- **RNN cell type:** While Long Short Term Memory (LSTM) cells were used for the recurrent layer in the previous model, this was changed back to a simple RNN layer without gates as they appear in Gated Recurrent Unit (GRU) or LSTM cells in the new model. This also corresponded to the original *DeepSpeech* model architecture, which did not use sophisticated gated cells in their recurrent layer for computational reasons.
- **number of features:** While the use of MFCC as features was examined in the previous model, the number of features was set to 13, a value which is found often used in acoustic modelling. The Mozilla implementation of *DeepSpeech* however doubled this number to 26, which is also what is used in the new model. In spite of the increase in the number of features, this value is still much smaller than the 160 filter banks used in the original *DeepSpeech* model. The amount of training data is therefore still expected to be smaller than in the original model.

The final model used in this project was therefore a variant of the original *DeepSpeech* model with the following simplifications:

- **optional LM:** While a key part in the original *DeepSpeech* model, the inclusion of a LM is optional in the simplified architecture. This is on account of the fact that such a LM needs to be trained first and to assess the influence of a trained LM (see the following sections).
- No convolution in first layer
- LSTM instead of SimpleRNN
- **smaller alphabet:** The Mozilla implementation uses an alphabet of 29 characters ( $a, b, c, \dots, z, space, apostrophe, blank$ ), which is also what is proposed in the *DeepSpeech* paper. This is due to the fact that apostrophes are frequently found in English word tokens (like "don't" or "isn't"). For simplification and also because a model should be trainable for any language, this character was dropped from the alphabet used in this project and also removed from the transcripts of all training samples by the normalization process. It is expected that missing apostrophes can be fixed by post-processing the inferred transcript with a LM.
- **no context:** The *DeepSpeech* paper proposes using combining each feature vector  $x_t$  (a frame in the spectrogram) with  $C \in \{5, 7, 9\}$  context frames. This context frame was dropped to keep the number of features in the input layer small. As a result, the first layer in the model only depends on the 26 features of the feature vector  $x_t$ .
- **no convolution in input layer:** The *DeepSpeech* paper proposes a series of optimization to reduce computational cost. Among these optimization is a convolution over time in the input layer with by striding with step size 2. Because the context frame was dropped in this project, the striding was also not applied in order not to lose the information from the intermediate frames.

<sup>4</sup>the variant used MFCC as features whereas the original paper proposed raw spectrograms

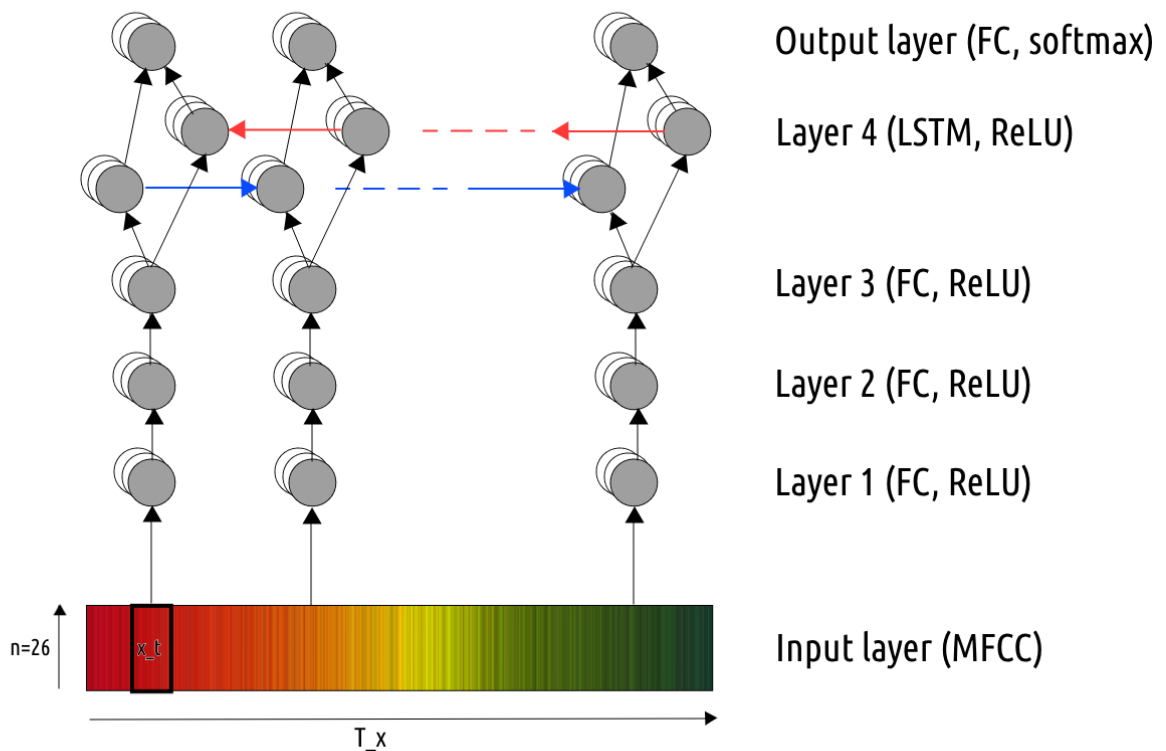


Figure 1: architecture of the simplified model

Figure 1 shows the architecture proposed in the *DeepSpeech* with the changes applied for this project. It looks similar to the one shown in the paper. Note the missing context frame, the use of MFCC features and LSTM cells in the recurrent layer.

### 3 The importance of LMs for ASR

Although CTC is the cost that is optimized during training, the usual metrics to evaluate an ASR system are WER and LER. The LER is defined as the mean normalized edit distance ( $ed(a, b)$  i.e. the number of insertions, deletions and changes required to produce string  $b$  from string  $a$ ) between a an inferred transcription (*prediction*) and the actual transcription (*ground truth* or *label*). It operates on character level and is sometimes also referred to as *Levenshtein Distance*. The WER builds upon the LER but operates on word level, i.e. it represents the number of words in a inferred transcription, that need to be inserted, deleted or changed in order to arrive at the ground truth.

If a single evaluation metric is required, the WER is often the better choice because it is more related to the way humans would assess the quality of a transcription: A transcription which might sound correct when read out loud but is full of spelling mistakes is not a good transcription. A LM can help inferring orthographically correct words from sequences of characters detected by CTC and hence decrease the WER. Therefore, by using a LM the quality of transcriptions improves perceptibly, as the following example shows:

A LM models probabilities of sequences of words. Therefore a LM can be used to assign a probability to a sentence. A simple LM that does that is the  $n$ -gram LM.  $n$ -grams are overlapping tuples of words whose



transcript	value	LER
actual transcript	and i put the vice president in charge of mission control	1.00
inference without LM	ii put he bice president in charge of mission control	0.11
inference with LM	i put the vice president in charge of mission control	0.08

Table 1: examples of inferred transcripts with pre-trained DeepSpeech model with and without LM  
(sample 20161203potusweeklyaddress from the ReadyLingua corpus)

probability can be approximated by training on massive text corpora. Although a lot of research has been made in the field of using NN for language modelling (like for machine translation),  $n$ -grams LM are still widely used and often the right tool for many tasks (Jurafsky and Martin 2019), because they are faster to train and require significantly less training data.

The Mozilla implementation includes an  $n$ -Gram LM using *KenLM*. The LM is queried while decoding the numeric matrices produced by CTC using *Beam Search* or *Best-Path* decoding. It uses a *trie* and precompiled custom implementations in C of *TensorFlow*-operations to maximize performance and dedicated weights for the influence the number of valid words and the LM itself on the inferred transcription. It is therefore deeply baked in with the decoding process.

A simpler variant that is used in this project is to infer the transcriptions first with *Beam Search* or *Best-Path* decoding using the standard tools provided by Keras. The inferred transcriptions are then post-processed by running it through some sort of spell-checking, which is done as follows:

- split the sentence into words
- for each word  $w_i$  in the sentence check the spelling by generating the set  $C_i$  of possible corrections by looking it up in  $V$ , the vocabulary of the LM, as follows:

- if  $w_i \in V$  its spelling is already correct and  $w_i$  is kept as the only possible correction, i.e.

$$C_i = C_j^0 = \{w_i\}$$

- if  $w_i \notin V$  generate  $C_i^1$  as the set of all possible words  $w_i^1$  with  $ed(w_i, w_i^1) = 1$ . This is the combined set of all possible words with one character inserted, deleted or replaced. Keep the words from this combined set that appear in  $V$ , i.e.

$$C_i = C_i^1 = \{w_i^1 \mid (w_i, w_i^1) = 1 \wedge w_i^1 \in V\}$$

- if  $C_i^1 = \emptyset$  generate  $C_i^2$  as the set of all possible words  $w_i^2$  with  $ed(w_i, w_i^2) = 2$ .  $C_i^2$  can be recursively calculated from  $C_i^1$ . Again only keep the words that appear in  $V$ , i.e.

$$C_i = C_i^2 = \{w_i^2 \mid ed(w_i, w_i^2) = 2 \wedge w_i^2 \in V\}$$

- if  $C_i^2 = \emptyset$  keep  $w_i$  as the only word, accepting that it might be either misspelled, a wrong word, gibberish or simply has never been seen by the LM, i.e.

$$C_i = C_i^{>2} = \{w_i\}$$

- for each possible spelling in  $C_i$  build the set  $P$  of all possible 2-grams with the possible spellings in the next word as the cartesian product of all words, i.e.

$$P = \{(w_j, w_{j+1}) \mid w_j \in C_j \wedge w_{j+1} \in C_{j+1}\}, \quad C_j \in \{C_i^0, C_i^1, C_i^2, C_i^{>2}\}$$

- score each 2-gram calculating the log-based probability using a pre-trained 2-gram-LM

## 4 Plotting a learning curve

The total length of all transcribed speech segments in the *LibriSpeech* corpus is roughly 47 days (1141 hours). Training on all these samples was not feasible within project time. However, we can get an estimation of the expected learning progress by plotting a *learning curve*. This was done using transcribed audio segments from the *LibriSpeech* corpus. Training was done on exponentially increasing amounts of training data (1, 10, 100 and 1000 minutes of transcribed audio). To assess the impact of using a LM, the learning curve was plotted twice: Figure 2 shows the learning curve for WER and LER without using a LM. Figure 3 shows the learning curve for training with a 4-gram LM.

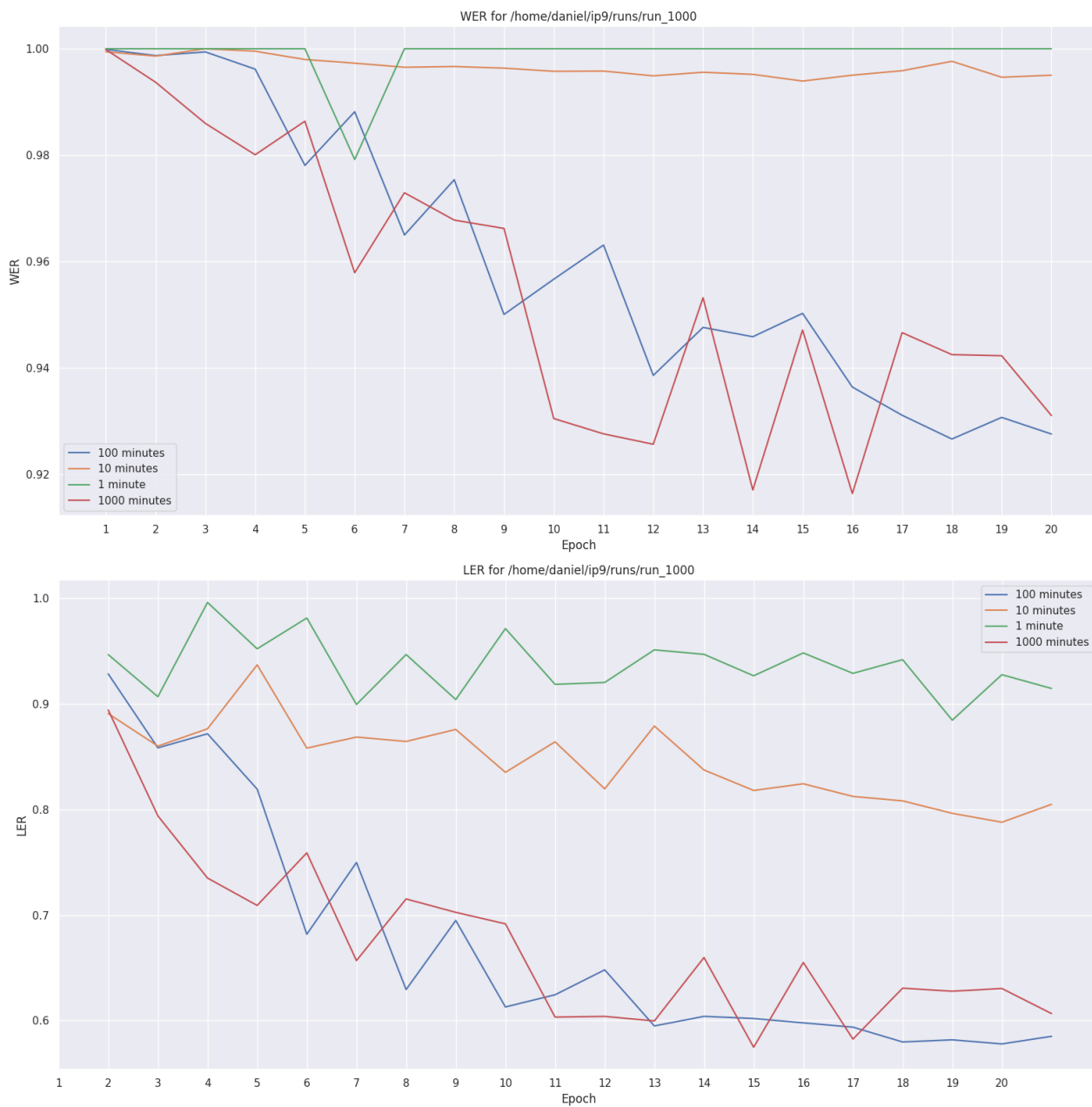


Figure 2: Learning curve for training on 1/10/100/1000 minutes of transcribed audio from ReadyLingua using a 2-gram LM and Beam-Search decoding

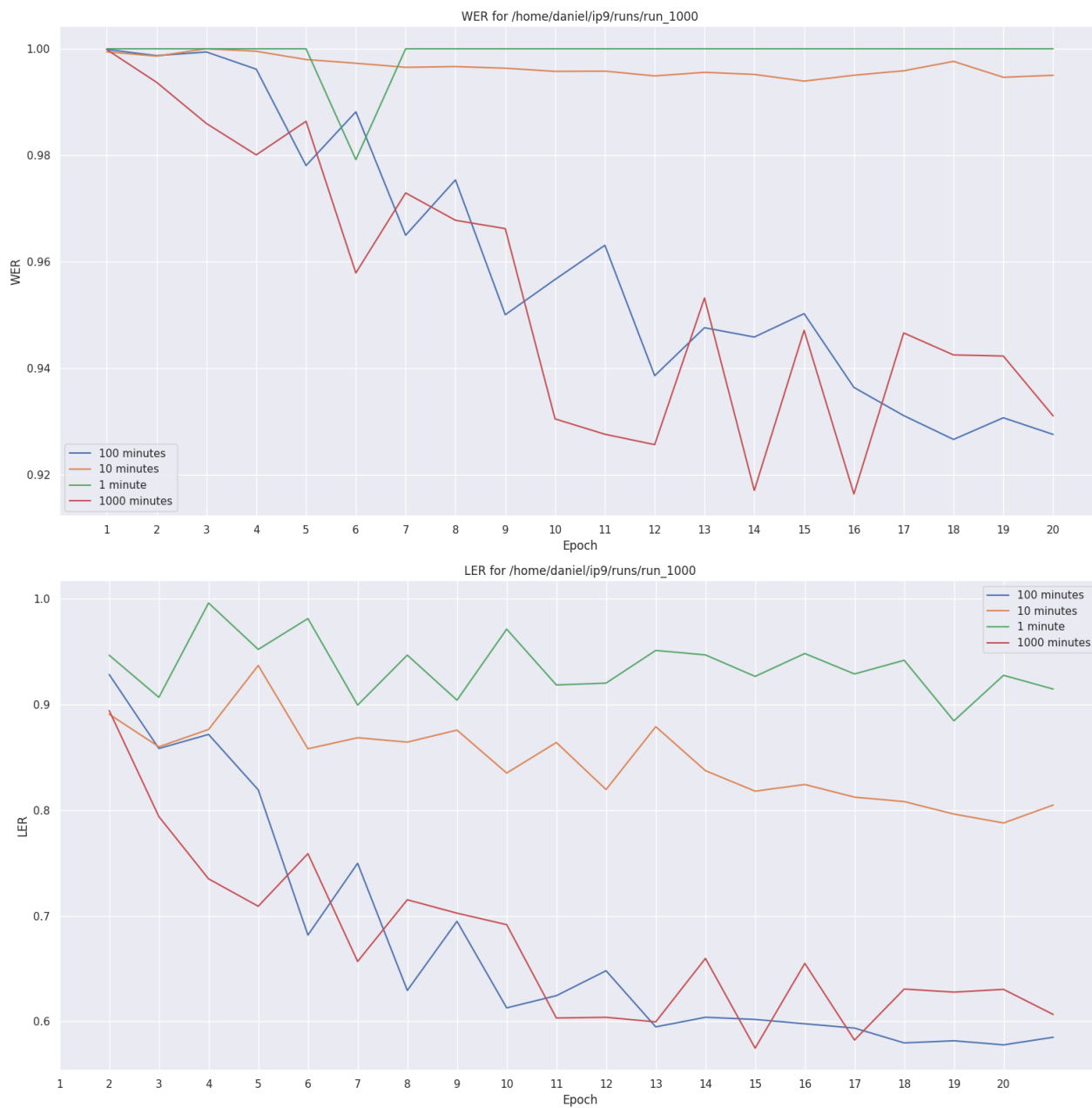


Figure 3: Learning curve for training on 1/10/100/1000 minutes of transcribed audio from ReadyLingua using a 2-gram LM and Beam-Search decoding

## 5 Creating a Language Model for German

The learning curve above was plotted using the results from training on English samples. To compare get an intuition about whether the conclusions made are transferable to training on samples in other languages, the training was done again on German samples. In order to do this, a LM for German had to be created first.

### 5.1 n-Gram Language Models

To understand the following sections better it might be helpful to get a quick recap about LM. LM are probabilistic models that model the likelihood of a given sequence of characters or words. The most widely used type for word-based models LMs are  $n$ -gram LM. However, such models can estimate probabilities only for words that appear in the vocabulary of the corpus they were trained on. All other words are Out Of Vocabulary (OOV) words with a probability of 0. Single words are 1-grams, but the same applies for  $n$ -grams of any order. But because of combinatorial explosion the  $n$ -gram, such LM suffer from sparsity with increasing order  $n$ . To handle OOV issues efficiently, a technique called *smoothing* is applied. A very rudimentary form of smoothing is *Laplace Smoothing*, which assigns a minimal count of 1 to every  $n$ -gram. All other counts are also increased by adding 1. This prevents counts of zero, which is important when calculating the perplexity of a model (because the count appears in the divisor and we cannot divide by zero) Although with Laplace Smoothing a very low probability is assigned to previously unseen  $n$ -grams (which results in a high perplexity), it performs poorly in application. A better way of smoothing is achieved using *Kneser-Ney Smoothing*.

### 5.2 Creating a raw text corpus and training the model

A LM was trained on a raw text corpus of German Wikipedia articles using KenLM (Heafield 2011). The articles were pre-processed to meet the requirements of *KenLM*. It was normalized by removing Wiki markup, punctuation and making everything lowercase. Accentuated characters (like è, é, ê, etc.) and special characters like the German ß were translated to their most similar ASCII-equivalent (e resp. ss) to account for ambiguous spelling and reduce the number of words. Umlauts (although not part of the ASCII codeset) were kept as-is because they are used very frequently in German.

Because *KenLM* expects the input as sentences (one sentence per line), the raw text was further tokenized into sentences and words using NLTK (Loper and Bird 2002). Word tokens that contain only numeric characters (such as year numbers) are changed to <unk>, a special token which is traditionally used to denote an OOV word. Although numeric tokens occur frequently in the Wikipedia articles, they are unwanted in the corpus because they do not carry any semantic meaning and because there is an infinite number of possible numbers.

The following lines are an excerpt of a article in the German Wikipedia along with its representation in the corpus.

Die Größe des Wörterbuchs hängt stark von der Sprache ab. Zum einen haben durchschnittliche deutschsprachige Sprecher mit circa 4000 Wörtern einen deutlich größeren Wortschatz als englischsprachige mit rund 800 Wörtern. Außerdem ergeben sich durch die Flexion in der deutschen Sprache in etwa zehnmal so viele Wortformen, wie in der englischen Sprache, wo nur viermal so viele Wortformen entstehen. (German Wikipedia article about Speech Recognition<sup>5</sup>)

```
1 die grösse des wörterbuchs hängt stark von der sprache ab
2 zum einen haben durchschnittliche deutschsprachige sprecher mit circa <unk> wörtern
   einen deutlich grösseren wortschatz als englischsprachige mit rund <unk> wörtern
```

<sup>5</sup><https://de.wikipedia.org/wiki/Spracherkennung>

3 ausserdem ergeben sich durch die flexion in der deutschen sprache in etwa zehnmal so viele wortformen wie in der englischen sprache wo nur viermal so viele wortformen entstehen

### Listing 1: Representation in corpus

The final corpus contained data from 2,221,101 Wikipedia articles (42,229,452 sentences, 712,167,726 words, 8,341,157 unique words). A 4-gram *KenLM* model was trained on this corpus. *KenLM* uses *Kneser-Ney Smoothing* and automatically adds the unigrams `<s>`, `</s>` and `<unk>` internally for sentence beginnings, endings and OOV words).  $n$ -grams can be represented with a tree structure <sup>6</sup>, which allows for pruning. Pruning 1-grams would help getting rid of obvious spelling mistakes and very rare tokens that only appear in very special contexts. Unfortunately, *KenLM* does not support pruning unigrams, only higher-order  $n$ -grams. Therefore the vocabulary used for training was limited to the 500,000 most frequent words not containing numbers with a minimum length of 2 characters <sup>7</sup>. This amount is about the same what was used in the *DeepSpeech* paper (Hannun et al. 2014). The words in the vocabulary make up 96.42% of the whole corpus. The most frequent word in the vocabulary is the `<num>` token (29,659,021 counts), the least frequent word is *Flachmeeres* (31 counts). Note that the latter is actually a derivate of another German word *Flachmeer*, which also appears in the vocabulary (109 counts). Having different flexions of the same word is characteristic for German texts. Handling them would require lemmatization and/or stemming the corpus, which has not been done for simplicity. It is also doubtful whether this would actually help improving the quality of inferred transcripts, since humans do not speak in lemmata or stems.

By limiting the vocabulary to the 500k most frequent words, the unigrams in the *KenLM* were artificially pruned. The number of unigrams was therefore 500,003 (one unigram for each word in the vocabulary plus one each for the `<s>`, `</s>` and `<unk>` tokens). The  $n$ -grams used to train the LM have been pruned by setting the minimal threshold for  $n$ -Grams of any order ( $n \in 1..4$ ) to 40. This is the value that Google used (reference from Jurafsky). Pruning unigrams helped getting rid of obvious spelling mistakes and very rare tokens that only appear in very special contexts (like the tokens *aaaaa* or *zzyzyxx*) (EDIT: Pruning unigrams is not supported by *KenLM*, but the vocabulary can be limited). Such words are mostly not no real German words and should therefore not be trained on. Pruning higher-order  $n$ -grams was done to increase performance (both in space and time).

## 5.3 Evaluating the model

The best way to evaluate a LM is to embed it in an application and measure how much the application improves (Jurafsky and Martin 2019). This is called *extrinsic evaluation* and has been done by comparing the learning curves with and without using a LM. However, to measure the performance of a LM independently (*intrinsic evaluation*) one would have to provide a test set containing unseen sentences and assess the scores of the LM on their  $n$ -grams. The results can then be compared to a reference LM: Whatever model produces higher probabilities (or lower perplexity) to the  $n$ -grams in the test set is deemed to perform better. However, because models can only be compared if they use the same vocabulary (Jurafsky and Martin 2019), this would require training the reference model would need to be trained on the same corpus, which can become very time consuming.

*KenLM* has been extensively compared to other LM implementations like the SRI Language Modelling Toolkit (SRILM) both in terms of speed and accuracy. It has been found to be both faster and more memory efficient (Heafield 2011) than the fastest alternative. Its low memory profile makes it runnable on a single machine, while other algorithms like *MapReduce* target clusters (Heafield et al. 2013). This was a big advantage

<sup>6</sup>note that *KenLM* offers a so called *PROBING* data structure, which is fundamentally a hash table combined with interpolation search, a more sophisticated variant of binary search, which allows for constant space complexity and linear time complexity. This does however not change the fact that  $n$ -grams can conceptually be thought as a tree of grams

<sup>7</sup>this constraint was imposed because NLTK did sometimes not tokenize abbreviations like *z.B.* correctly, which resulting in two separate tokens *z* and *b*

especially for this project. The probabilistic performance of *KenLM* has been evaluated by training a 5-gram model on a 126 billion token corpus (393 unique words) (Heafield et al. 2013). This model was embedded in some Machine Translation systems (Czech-English, French-English and Spanish-English) . Evaluation was done by calculating the BLEU score and comparing it to embeddings of other LM. *KenLM* placed first in all submissions.

Because of time constraints and because *KenLM* has already been extensively evaluated on English I resign from evaluating my German LM intrinsically, although the corpus used for training is not as big as the one used in Heafield et al. 2013. To this day *KenLM* is widely recognized as the best performing LM out there, which is emphasized by the usage of a *KenLM* model in the Mozilla implementation of *DeepSpeech*.

To still get an intuition about how well the model performs, the model's score on some test sentences were calculated. To make sure the sentences could not have been seen during training, the following set of 5 sentences of the current newspaper (date after creation of the Wikipedia dump) was used:

- *Seine Pressebeauftragte ist ratlos.*
- *Fünf Minuten später steht er im Eingang des Kulturcafés an der Zürcher Europaallee.*
- *Den Leuten wird bewusst, dass das System des Neoliberalismus nicht länger tragfähig ist.*
- *Doch daneben gibt es die beeindruckende Zahl von 30'000 Bienenarten, die man unter dem Begriff «Wildbienen» zusammenfasst.*
- *Bereits 1964 plante die US-Airline Pan American touristische Weltraumflüge für das Jahr 2000.*

The score for each sentence was calculated. Then the words of the sentences were shuffled and the score was calculated again. A good LM should calculate a higher probability for the original sentence, because the shuffled sentence is most likely to be gibberish. All sentences have been normalized the same way sentences were preprocessed for training. Table 3 shows the results of the comparison. We can see that the probabilities for the shuffled sentences are much lower than for the sentences where the words appear in the correct order.

original sentence (normalized)	score	permutation	score
seine pressebeauftragte ist ratlos	-17.58	ist ratlos pressebeauftragte seine	-21.52
fünf minuten später steht er im eingang des kulturcafes an der zürcher europaallee	-40.23	des er minuten zürcher kulturcafes steht europaallee eingang fünf im später an der	-57.69
den leuten wird bewusst dass das system des neoliberalismus nicht länger tragfähig ist	-35.52	system nicht das ist dass leuten tragfähig des neoliberalismus den bewusst länger wird	-51.27
doch daneben gibt es die beeindruckende zahl von <num> bienenarten die man unter dem begriff wildbienen zusammenfasst	-48.36	dem gibt wildbienen zahl beeindruckende doch man zusammenfasst es daneben bienenarten von die unter die <num> begriff	-75.95
bereits <num> plante die usairline pan american touristische weltraumflüge für das jahr <num>	-58.04	plante touristische für jahr pan american das bereits usairline <num> <num> weltraumflüge die	-64.02

Table 2: Comparison of log10-probabilities calculated for news sentences and a permutation of their words

## 6 Measuring the performance of the pipeline

The performance of the pipelined approach can be measured by assessing the alignments it produced. Because the number of theoretically possible alignments is massive there is no reference value to compare to. In fact, assessing the quality of an alignment is somewhat subjective because the size of the aligned chunks might be considered good by one person and bad by another person. We can however derive a few objective criteria that make up a good alignment:

1. The aligned partial transcripts should cover the whole transcript
2. The aligned partial transcripts should not overlap
3. The aligned partial transcripts should be at the correct position (i.e. they should contain the text that is actually spoken)

These criteria can be quantified with the following metrics (note the correlation<sup>8</sup>):

criterion	metric	symbol	correlation
1	length of text in ground truth that is not aligned vs. total length of the ground truth	$C$	negative
2	length of overlaps in alignments vs. total length of alignments	$O$	negative
3	average Levensthein similarity between the transcript and the text in the ground truth corresponding to its alignment	$D$	positive

Table 3: Metrics to evaluate the quality of alignments

These metrics can be weighted and reduced to a single number if necessary. Figure 4 shows an example of an alignment. The sentence was split into speech segments by VAD which were then aligned with the whole sentence (ground truth). Some speech segments were misaligned, resulting in an overlap. Some of the transcriptions contain mistakes. Note that the Levensthein Similarity ( $ls(s_1, s_2)$ ) measures the similarity of two strings  $s_1$  and  $s_2$  as the normalized edit distance and is calculated as follows:

$$ls(s_1, s_2) = 1 - \frac{ed(s_1, s_2)}{\max(len(s_1), len(s_2), len(s_3))} \quad (1)$$

For  $t_1 = i \text{ see}, t_2 = i \text{ c}, t_3 = \text{sad the blind men}, t_4 = \text{to his blind daughter}$  The metrics can be calculated as follows:

$$C = \frac{len(i \text{ see})}{len(i \text{ see } i \text{ see said the blind man to his deaf daughter})} = \frac{5}{51} = 0.098 \quad (2)$$

$$\begin{aligned} O &= \frac{len(i \text{ see})}{len(i \text{ see}) + len(i \text{ c}) + len(\text{sad the blind men}) + len(\text{to his deaf daughter})} \\ &= \frac{5}{5 + 3 + 17 + 20} = 0.11 \end{aligned} \quad (3)$$

$$\begin{aligned} D &= \frac{ls(i \text{ see}, t_1) + ls(i \text{ see}, t_2) + ls(\text{said the blind man}, t_3) + ls(\text{to his deaf daughter}, t_4)}{4} \\ &= \frac{1 + 0.4 + 0.88 + 1}{4} = 0.82 \end{aligned} \quad (4)$$

<sup>8</sup>positive correlation: higher is better, negative correlation: lower is better



«I see, I see», said the blind man to his deaf daughter.

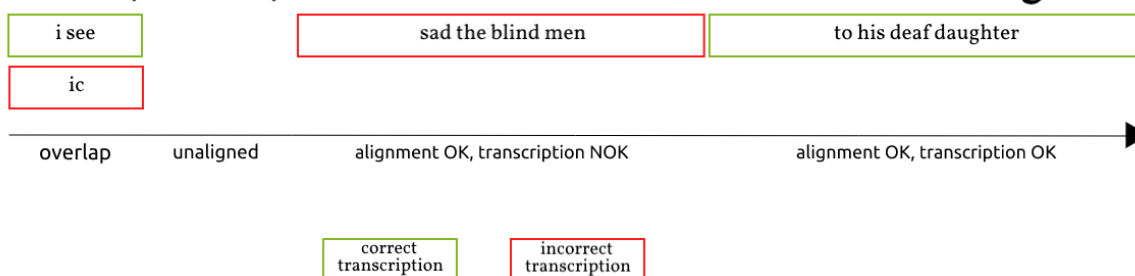


Figure 4: Example for a partially correct alignment of parts of a sentence

## List of Figures

1	architecture of the simplified model . . . . .	5
2	Learning curve for training on 1/10/100/1000 minutes of transcribed audio from <i>ReadyLingua</i> using a 2-gram LM and Beam-Search decoding . . . . .	8
3	Learning curve for training on 1/10/100/1000 minutes of transcribed audio from <i>ReadyLingua</i> using a 2-gram LM and Beam-Search decoding . . . . .	9
4	Example for a partially correct alignment of parts of a sentence . . . . .	14

## List of Tables

1	examples of inferred transcripts with pre-trained DeepSpeech model with and without LM (sample 20161203potusweeklyaddress from the <i>ReadyLingua</i> corpus . . . . .	6
2	Comparison of log10-probabilities calculated for news sentences and a permutation of their words . . . . .	12
3	Metrics to evaluate the quality of alignments . . . . .	13

## References

- Graves, Alex, Santiago Fernández, and Faustino Gomez (2006). "Connectionist temporal classification: Labelling unsegmented sequence data with recurrent neural networks". In: *In Proceedings of the International Conference on Machine Learning, ICML 2006*, pp. 369–376.
- Hannun, Awni Y. et al. (2014). "Deep Speech: Scaling up end-to-end speech recognition". In: *CoRR* abs/1412.5567. arXiv: 1412.5567. URL: <http://arxiv.org/abs/1412.5567>.
- Heafield, Kenneth (July 2011). "KenLM: Faster and Smaller Language Model Queries". In: *Proceedings of the EMNLP 2011 Sixth Workshop on Statistical Machine Translation*. Edinburgh, Scotland, United Kingdom, pp. 187–197. URL: <https://kheafeld.com/papers/avenue/kenlm.pdf>.
- Heafield, Kenneth et al. (Aug. 2013). "Scalable Modified Kneser-Ney Language Model Estimation". In: *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics*. Sofia, Bulgaria. URL: [http://kheafeld.com/professional/edinburgh/estimate%5C\\_paper.pdf](http://kheafeld.com/professional/edinburgh/estimate%5C_paper.pdf).
- Jurafsky, Daniel and James H. Martin (2019). *Speech and Language Processing (Draft of 3rd Edition)*. URL: <https://web.stanford.edu/~jurafsky/slp3/>.
- Loper, Edward and Steven Bird (2002). "NLTK: The Natural Language Toolkit". In: *In Proceedings of the ACL Workshop on Effective Tools and Methodologies for Teaching Natural Language Processing and Computational Linguistics*. Philadelphia: Association for Computational Linguistics.
- Morais, Reuben (2017). *A Journey to <10% Word Error Rate*. URL: <https://hacks.mozilla.org/2017/11/a-journey-to-10-word-error-rate> (visited on 09/14/2018).

## Acronyms used in this document

<b>ASR</b>	Automatic Speech Recognition
<b>CTC</b>	Connectionist Temporal Classification
<b>E2E</b>	end-to-end
<b>FA</b>	Forced Alignment
<b>FHNW</b>	University of Applied Sciences
<b>GPU</b>	Graphics Processing Unit
<b>GRU</b>	Gated Recurrent Unit
<b>LER</b>	Label Error Rate
<b>LM</b>	Language Model
<b>LSTM</b>	Long Short Term Memory
<b>LSA</b>	Local Sequence Alignment
<b>MFCC</b>	Mel-Frequency Cepstral Coefficients
<b>NN</b>	Neural Network
<b>RNN</b>	Recurrent Neural Network
<b>SGD</b>	Stochastic Gradient Descent
<b>STT</b>	Speech-To-Text
<b>OOV</b>	Out Of Vocabulary
<b>SRILM</b>	the SRI Language Modelling Toolkit
<b>SW</b>	Smith Waterman
<b>VAD</b>	Voice Activity Detection
<b>WER</b>	Word Error Rate

## 7 Ehrlichkeitserklärung

Hiermit erkläre ich, dass ich die vorliegende schriftliche Arbeit selbstständig und nur unter Zuhilfenahme der in den Verzeichnissen oder in den Anmerkungen genannten Quellen angefertigt habe. Ich versichere zudem, diese Arbeit nicht bereits anderweitig als Leistungsnachweis verwendet zu haben. Eine Überprüfung der Arbeit auf Plagiate unter Einsatz entsprechender Software darf vorgenommen werden.

Würenlingen, September 28, 2018

Daniel Tiefenauer