

IP9

Daniel Tiefenauer

November 22, 2018

Abstract

tbd.

Contents

1	Introduction	1
1.1	Scope and overall goal	1
1.2	Chosen approach and previous work	1
1.2.1	Previous results and problems	1
1.3	Goal of this project	2
2	Training an Recurrent Neural Network (RNN) for Automatic Speech Recognition (ASR)	4
2.1	Related research	4
2.2	<i>DeepSpeech</i> : A reference model	4
2.3	Exploiting the <i>DeepSpeech</i> model	4
2.4	A simpler model	5
3	Integrating a Language Model	7
3.1	Measuring and improving the performance of a Speech-To-Text (STT) engine	7
3.2	Language Models in ASR	7
3.3	A simple spell checker	8
3.3.1	Reducing the vocabulary size	9
3.4	Further thoughts and considerations	9
4	Plotting a learning curve	10
4.1	Previous corpora and their problems	10
4.2	The <i>CommonVoice</i> (CommonVoice (CV)) Corpus	10
4.3	Plotting the learning curve	11
4.3.1	Decoder dimension	11
4.3.2	Language Model (LM) dimension	12
4.4	Results and interpretation	13
4.5	Regularization	15
4.6	Final thoughts and considerations	15
5	Measuring the performance of the pipeline	16
5.1	The quality of alignments	16
5.2	Test and results	17
5.3	Pipeline performance using an ASR model for a different language	18
6	Forced Alignment for other languages	19
6.1	Inferring German transcripts	19
6.2	Data augmentation	19
6.3	Creating a Language Model for German	20
6.4	n-Gram Language Models	20
6.4.1	Perplexity, discount and smoothing	20
6.4.2	Kneser-Ney Smoothing	21
6.5	Creating a raw text corpus	21
6.6	Training the LM	22
6.6.1	Data structures	22
6.6.2	Quantization	22
6.6.3	Pointer Compression	22
6.6.4	Building the model	23
6.7	Evaluating the LM	23
6.7.1	Extrinsic and intrinsic evaluation	23
6.7.2	Evaluation of KenLM	23

6.7.3	Evaluation of the German Wikipedia LM	23
6.7.4	Evaluation 1: Comparing scores of randomized sentences	24
6.7.5	Experiment 2: Word predictor	25
6.8	Results	25
7	Conclusion	26
8	Ehrlichkeitserklärung	31

1 Introduction

This report documents the progress of the project *Speech-To-Text Engine for Forced Alignment*, my master thesis at University of Applied Sciences (FHNW) (referred to as *IP9*). Some preliminary work has been done in a previous project (referred to as *IP8*). The overall goal, project situation and some background information are described in detail in the project report for *IP8* and shall not be repeated here. Only a quick recap of the relevant terms and aspects is given as far as they are relevant for the understanding of this document.

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 the textual 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 *Forced Alignment (FA)*. 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 pursued in this project.

1.2 Chosen approach and previous work

The approach chosen for 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 by means of Local Sequence Alignment (LSA).

All these parts were treated as individual stages of a pipeline:

- **VAD:** the audio was split into non-silent parts
- **ASR:** each part was transcribed resulting in a partial transcript, which can contain transcription errors
- **LSA:** each partial transcript was localized within the original transcript

Since the quality of the ASR stage has an imminent impact on the subsequent LSA stage, the quality of the alignments depends heavily on the quality of the partial transcripts. This makes the ASR stage the crucial stage of the pipeline. 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 and problems

For the VAD stage, an implementation¹ of *WebRTC*² was used. This implementation which has proved 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.

For the ASR stage on the other hand, no RNN could be trained that was capable of transcribing the audio segments with a quality that was high enough for the LSA stage. The main problems were the lack of readily available training data in high quality, very long training times and therefore very long feedback cycles. Because the ASR stage is at the heart of the pipeline, the self-trained model was replaced by Google Cloud

¹<https://github.com/wiseman/py-webrtcvad>

²<https://webrtc.org/>

Speech (GCS)³, which was embedded into the pipeline over an API. This step however made the pipeline dependent on a commercial product, whose inner workings remain unknown and who cannot be tuned to the project's needs. Furthermore, although the quality of the transcriptions produced by GCS is very good, it might be an overkill for the purpose of this project and the API calls are subject to charges.

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. Some experiments were made to find out what features can be used to train a RNN for the ASR stage. The features considered 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 *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 lies on the ASR stage. Because the pipeline should become language-agnostic and self-contained, a RNN must be trained that can be used in this stage in the pipeline. Such a RNN could be a simplified variant of the *DeepSpeech* model. The idea behind this are the following two points:

- A simpler model will not be able to produce transcripts with the same quality as complex models like *DeepSpeech* or GCS. But since the superordinated goal of this project is FA and not Speech Recognition, this may not be required in the first place. The RNN only needs to be *good enough* for the downstream LSA stage
- On the other hand, a simpler model will probably need less training data to learn. This might open up new possibilities to align audio and text in other languages, for which there is not a third-party solution yet, such as Swiss German.

The goal of this project is therefore to make statements as to under what conditions the ASR stage can be implemented. For this, various combinations of network or data properties are explored as well as varying amounts of training data. Concretely, the following questions shall be answered:

- **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/or 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

³<https://cloud.google.com/speech-to-text/>

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.

Answering above questions should help estimating the effort to create a generic pipeline. ⁴

⁴Because 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, the focus for this project is not on 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 for any language in the ASR stage of the pipeline.

2.1 Related research

Hier etwas über Transfer Learning (z.B. wie in <https://arxiv.org/pdf/1706.00290.pdf>) und warum es nicht eingesetzt wurde (CNN anstatt RNN, Zeitaufwand). Layer Freezen bringt ausserdem offenbar auch nix. (Kunze et al. 2017)

2.2 *DeepSpeech*: A reference 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 ⁵. This implementation uses a variant of the RNN originally proposed in the *DeepSpeech* paper (Graves, Fernández, and Gomez 2006)⁶. Since the implementation also uses a LM, the quality of the model is measured as the percentage of misspelled or wrong words (referred to as 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) and serves as a reference model in this project.

2.3 Exploiting the *DeepSpeech* model

Since the goal of this project is to provide alignments for any language, one possible approach would be to train a model using the existing Mozilla implementation by providing training-, validation- and test-data for each language. However, this approach has several drawbacks:

1. The *DeepSpeech* implementation was explicitly designed for ASR. In such settings a low WER is desirable. But because accurate speech recognition is not the main concern in this project, the architecture of the Mozilla implementation might be overly complicated.
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 Mozilla 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 architecture of the RNN model from Mozilla was used only as a basis for a simplified version. This version should (hopefully) require less training data to converge and still produce partial transcriptions that can be aligned.

⁵<https://github.com/mozilla/DeepSpeech>

⁶the variant used MFCC as features whereas the original paper proposed raw spectrograms

2.4 A simpler model

The implementation of the RNN used for STT in the previous IP8 project was done in Python using Keras⁷. This model is further referred to as *previous model*. Unfortunately, the previous model did not perform very 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, even though the RNN was rather simple and no LM was used. Training on aligned speech segments from the *LibriSpeech* corpus would have taken approximately two months on a single Graphics Processing Unit (GPU). This duration is at consistent with the experience made by the Machine Learning team at Mozilla Research, which used a cluster of 16 GPUs that required about a week (Morais 2017) to train a variant. For the purpose of this project however, such a long training time was a serious impediment.

In the course of this project, the previous model was examined more closely to find out what works best and to help the model converge. A few changes were made to arrive at a new model which was able to learn something meaningful. This model is further referred to as *new model*. The new model started to infer transcripts that – although still not perfect – resembled the ground truth. The following list summarizes the differences between the previous and the new model:

- **Optimizer:** The new model uses Stochastic Gradient Descent (SGD) instead of Adam. Adam was used in the previous model because it is the Optimizer used in the Mozilla implementation of Deep Speech (DS). However, this optimizer did not seem to work for the simplified model.
- **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. This is also the number of features used in the new model. Despite 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 new model constitutes a simplified variant of the original *DeepSpeech* model with the following simplifications and changes applied:

- **Different application of LM:** In the Mozilla implementation the use of a LM is baked in with the training process, i.e. it is integrated in the decoding process. With The edit distance between prediction and ground truth is then included in the loss which is minimized. The simplified model also uses a LM, but does not include it in the training process. Instead, the LM is applied in some sort of post-processing to improve the quality of the decoded predictions.
- **Different features:** MFCC with 26 filter banks instead of Spectrogram with 161 filterbanks, because that's what the Mozilla implementation uses
- **No convolution in first layer:** Whereas Graves, Fernández, and Gomez 2006 propose a convolution over time in the first layer for performance reasons, this is not done in the simplified model.
- **LSTM instead of SimpleRNN:** Whereas Graves, Fernández, and Gomez 2006 deliberately refrain from using Long Short Term Memory (LSTM) cells for various reasons, the Mozilla implementation has shown that it is possible to implement the *DeepSpeech* model using LSTM cells. Since the simplified model is based on the Mozilla implementation, it also uses LSTM cells.
- **dynamic alphabet:** The Mozilla implementation uses an alphabet with 29 characters⁸, 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"*). The apostrophe is therefore an integral part of English words, but not for other languages. Vice versa, other languages may use a different alphabet (like

⁷<https://keras.io>

⁸*a, b, c, ..., z, space, apostrophe, blank*

German, where umlauts are prevalent). Because the number of characters in the alphabet determines the number of nodes in the output layer, the output layer has different shapes for different languages.

- **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.

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.

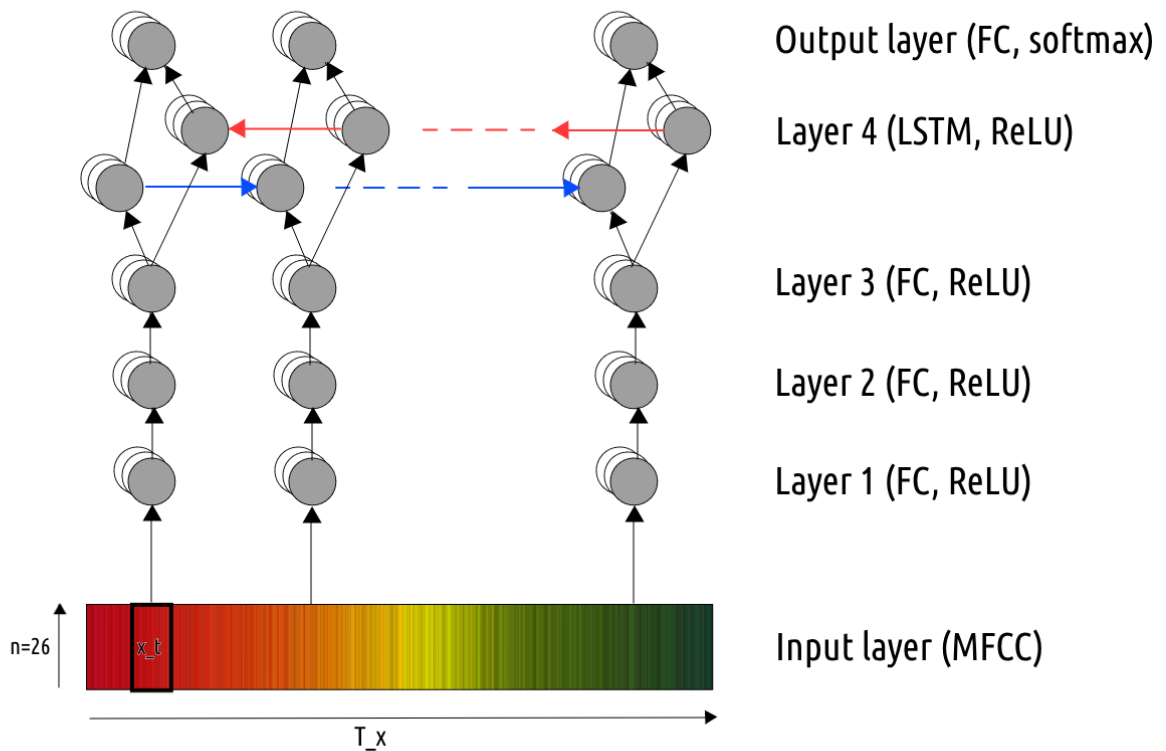


Figure 1 – Architecture of the simplified model. The cell type and the activation function is indicated in brackets for each layer (FC=Fully-Connected)

3 Integrating a Language Model

This chapter outlines the importance of LMs for ASR. It also describes how a LM was integrated into the simplified DS model as an attempt to improve the quality of the transcription.

3.1 Measuring and improving the performance of a STT engine

Although CTC is the cost that is optimized during training, the usual metrics to evaluate an STT system are WER and LER. These metrics correlate directly with the perceived quality of the system: Transcripts with a low WER and/or LER have a high similarity to the actual transcript and are therefore considered good.

The LER (sometimes also referred to as *Levenshtein Distance*) is defined as the mean normalized edit distance $ed(a, b)$ between two strings a and b . It operates on character level by counting the number of insertions (I), deletions (D) and substitutions (S) required to produce string a from string b . String a is the reference string, which in this project is the actual transcription of a speech segment (*ground truth* or *label*). String b is an inferred transcription produced by the simplified model (*prediction*).

The WER builds upon the LER and is therefore very similar. In contrast to LER however, WER operates on word level, i.e. it represents the number of words that need to be inserted, deleted or changed in an inferred transcription in order to arrive at the ground truth.

Both metrics can be normalized by dividing them by the length of the reference string i.e. the number of characters (LER) resp. the number of words (WER). If a single evaluation metric is required, the WER is often the better choice because it is more related to the way humans assess the quality of a STT engine: A transcription that might sound correct when read out loud, but is full of spelling mistakes, is not considered a good transcription. A LM can help inferring orthographically correct words from sequences of characters inferred by CTC and hence decrease the WER. Therefore, by using a LM the quality of transcriptions can improve considerably. Table 1 illustrates this with an example.

	transcript	LER	LER (norm.)	WER	WER (norm.)
GT	i put the vice president in charge of mission control	0	0.00	0	0.00
w/ LM	ii put he bice president in charge mission control	6	0.11	4	0.40
w/o LM	i put the vice president in charge mission control	3	0.06	1	0.10

Table 1 – Example for how a LM can help improve the quality of an inferred transcription by changing characters and words. Audio and ground truth (GT) were taken from the ReadyLingua corpus and the inference was made with the pre-trained DeepSpeech model.

3.2 Language Models in ASR

LMs model the probabilities token sequences. Because a sentence is a sequence of word-tokens, a LM can calculate its likelihood. Traditionally n -gram models (often simply referred to as *n-grams*) have been used for this task. n -grams are overlapping tuples of words whose probability can be approximated by training on massive text corpora. A special token `<unk>` is used for unknown tokens that do not appear in the training corpus. Because of combinatorial explosion and the dynamic nature of human language, the computational power and storage which are needed to train higher-order models explodes with increasing order n of the model. Most models are trained on an order of $n = 5$ or 6 .

Because the context of n -gram models is determined by its order they are somewhat restricted in that they do not take into account words outside the context to assess the probability of a sentence. 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 a good choice for many tasks (Jurafsky and Martin 2019). Because of their simplicity they are often faster to train and require significantly less training data than their neural counterparts.

3.3 A simple spell checker

The Mozilla implementation includes a 5-Gram LM, which is included as a part of the pre-trained model that can be downloaded from GitHub⁹. This LM was trained using *KenLM*. The LM is queried during training by 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. The LM is deeply baked in with the training process of the Mozilla implementation, which uses special hyperparameters to weight the score of the LM and the number of correctly inferred words.

According to Morais 2017 this tight integration is the culmination of various attempts to integrate a LM into the inference process. An early attempt used the LM as some sort of spell checker that was able to correct minor orthographic errors. Rather than including the LM-score for training, a spell-checker post-processes the inferences made by CTC *after* training. On one hand this deteriorates quality as has been shown by Mozilla, because no information from the LM is used during training. On the other post-processing the inferences is simpler and reduces complexity. It can also be implemented with the standard tools provided by Keras and does not need to be precompiled into C.

Post-processing the inferences with a spell checker was therefore the approach chosen for the simplified model. Because no additional weights for the LM have to be trained, the amount of training data is expected to be smaller than with those weights, which is the goal of the project in the first place.

The functionality of the spell checker can be summarized as follows (a more detailed and formal description can be found in the appendix):

1. Given an inference of space-separated word-tokens, process the words from left to right.
2. For each word check if it is contained in the vocabulary of the LM.
 - (a) If that is the case, continue with the next word.
 - (b) If not, create a list of variations with edit distance 1 and keep only those variations that appear in the vocabulary. Each of these variations is a possible continuation that can be scored by the LM.
3. If none of the variations appear in the vocabulary, create another list of variations with edit distance 2 to the original word. This list can be created recursively from the (complete) list of variations with edit distance 2. Again keep only those variations that appear in the vocabulary.
4. If none of the variations of the word with edit distance 2 are found in the vocabulary, use the original word as fallback. This can happen if the word is just gibberish or if the word is an actual valid word which does not appear in the training corpus for the LM and has therefore never been seen before. Note that in this step the word must not be substituted by the `<unk>` token because it may still be a valid word. Furthermore, replacing the word with the `<unk>` token can have a contrary effect on the alignment, because this token will most likely never appear in a valid transcript.

Above steps are repeated until the whole sentence is processed. For each word this yields a cascade of possible combinations. Each of these combinations can be scored by the LM as the sentence is being processed whereas only the n most likely prefixes are kept at each step (beam search). For this project, a beam width of 1.024 was used.

⁹<https://github.com/mozilla/DeepSpeech#getting-the-pre-trained-model>

3.3.1 Reducing the vocabulary size

The LM trained by Mozilla used texts from the *LibriSpeech* corpus¹⁰. Apart from lowercasing, the texts were not normalized or preprocessed. The resulting vocabulary is therefore very big and contains 973.673 unique words. Because no further preprocessing was done, it also contains some exotic words like "zzzz" and probably also misspelled words that happen to appear in the corpus. To train the LM, n -grams of order 4 and 5 were pruned with a threshold value of 1, meaning only 4- and 5-grams with a minimum count of 2 and higher are estimated¹¹. Because spelling errors are probably unique within the training corpus, 4- or 5-grams containing a misspelled word are unique too and are therefore pruned.

Above procedure might work well to estimate the likelihood of a sentence. For a simple spell checker however, such a big vocabulary might be counter-productive because it lowers the probability that an obviously wrong word is corrected because for some reason it found its way into the vocabulary. Vice versa a very large vocabulary raises the probability that a random-looking sequence of characters is "corrected" to a wrong word or that words are changed to said exotic or misspelled words. To prevent this, the original vocabulary was reduced into three vocabularies containing the 40.000, 80.000 and 160.000 most frequent words from the corpus each. These words make up 98.42%, 99.29% and 99.69% of the corpus.

To create the vocabularies, a list of unique words and their frequency was created from the corpus and sorted by frequency in descending order. Naturally, stop words like *the*, *and* or *of* are at the top of the list. The first 40.000, 80.000 resp. 160.000 words from this list were stored as the truncated vocabularies, the rest was discarded. Note that truncating the vocabulary only affects if and how words are exchanged by the spell checker during post-processing, not how a post-processed sentence's probability is estimated by the language model.

3.4 Further thoughts and considerations

The spell checker in this project uses the vocabulary with 80.000 words. This value was arbitrarily chosen and some unsystematic experiments were made to analyze the correctional capabilities of the spell checker. Because of time constraints and because it was unclear whether the spell checker would actually help improving the transcriptions, other vocabulary sizes were not evaluated. Further work may however try to find out an optimal vocabulary size for each language.

¹⁰<http://www.openslr.org/11>

¹¹see <https://github.com/mozilla/DeepSpeech/tree/master/data/lm>

4 Plotting a learning curve

This section describes how training progress was estimated by plotting a learning curve.

4.1 Previous corpora and their problems

The following two corpora were available for training from the IP8 project.

- **LibriSpeech (LS):** This corpus was created as an artifact of the IP8 project using raw data from OpenSLR. The raw is publicly available and can be downloaded¹². It consists of a number of audio files which were *partially* transcribed, i.e. there are parts in the audio for which the corresponding transcript is not exactly known (the audio contains *gaps*). The individual samples were obtained by exploiting metadata included in the download. The metadata includes a split into a training set (containing approximately 96% of the samples) and a validation resp. test set (each containing approximately 3% of the samples). The split was done manually into disjoint subset, i.e. ensuring each speaker was only included in one set. Additionally, other features like gender or accent were observed to achieve a similar distribution for each set. To leverage the efforts made by *OpenSLR*, this split was not changed.
- **ReadyLingua (RL):** This corpus was created from raw data provided by *ReadyLingua*. This data is proprietary and contains recordings in several languages which were manually aligned with their transcript. In contrast to the LS corpus, the raw data is fully aligned, i.e. there are no gaps in the audio. However, the metadata does not comprise a split into training-, validation- and test-set. Since the raw data contained recordings and transcripts in more than one language, separate splits were made for each language preserving a ratio of approximately 80/10/10% (relating to the total length of all recordings within each subset). Efforts were made to prevent samples from the same recording being assigned to different subsets. Other features were not observed, meaning the split into train/validation/test-set was done less carefully than in the LS corpus.

The model in the IP8 project was supposed to be trained on the LS corpus, because this corpus is much larger than the RL corpus. In the course of the project it became clear however that training on all samples from this corpus was not feasible within project time because training time would have taken more than two months. It also turned out that the LS corpus was probably less useful than initially assumed because the average sample length was much longer than the samples in the RL corpus. This made training even harder because convergence is much slower when training on long sequences. The RL corpus on the other hand consisted of shorter samples, but the total length of all samples was only a few hours compared to the 1000+ hours in the LS corpus.

4.2 The *CommonVoice* (CV) Corpus

Because of the aforementioned problems a new corpus was needed which combined the best of both worlds:

- it should contain a reasonable amount of speech samples to facilitate training an ASR model
- the average sample length should be short enough for the model to learn quickly.

The CV¹³ corpus is built and maintained and used actively by the Mozilla Foundation and exhibits both of these properties. This corpus is also used to train the Mozilla implementation of *DeepSpeech*. Datasets for various languages are being prepared and verified, each one containing speech samples of different contributors from all over the world. At the time of this writing, only the English dataset was available, but datasets for other languages will become publicly available at some time in the future. The English dataset comes pre-divided

¹²<http://www.openslr.org/12/>

¹³<https://voice.mozilla.org/en/data>

into training-, validation- and test-set of similar scale like the LS corpus. Each set consists of one audio file per sample and a CSV file containing the transcriptions for each sample.

For this project the English dataset was used to train the simplified model. Although still smaller than the LS corpus, the total length of all validated samples that can be used for training¹⁴ is much larger than the RL corpus while providing samples of similar length at the same time. Table 2 shows some statistics about the corpora described above.

Corpus	Language	total audio length	train/dev/test	# samples	Ø sample length	Ø transcript length
LS	English	24days, 7 : 13 : 18	93.51/3.32/3.16%	166,510	12.60	183.84
RL	English	5 : 38 : 39	80.39/10.13/9.48%	6,334	3.20	51.81
RL	German	1 : 58 : 30	81.14/10.26/8.60%	2,397	2.89	45.55
CV	English	10days, 1 : 02 : 53	96.04/1.99/1.98%	201,252	4.31	48.07

Table 2 – Statistics about corpora that were available for training. The sample length is given in seconds, the transcript length as the number of characters.

4.3 Plotting the learning curve

The time needed to train an ASR model on all samples of the CV corpus is still too long for the available project time. We can however still get an estimate of the learning progress by plotting a *learning curve*. For this, exponentially increasing amounts of training data (1, 10, 100 and 1,000 minutes of transcribed audio) were used. Training was done making 30 full passes (*epochs*) over the training data. The training samples were sorted by the length of their audio signal¹⁵ and then zero-padded, yielding samples of the same length in each batch¹⁶. After each epoch, the progress was monitored by inferring the transcriptions for previously unseen samples from the validation set. Those samples were shuffled after each epoch. The CTC-loss for training and validation was plotted for each amount, yielding separate curves for the training- and the validation-loss. Comparing both curves allows for making statements about at what point the Neural Network starts to overfit.

Complementary to the CTC-loss, the mean values for the LER and WER metric over all samples in the validation set is calculated after each epoch, yielding the curves for the LER resp. WER. Observing these plots can give some insight about how well the network performs on unseen examples.

Both loss and metrics were compared along two dimensions:

- **The decoder dimension**, comparing the two distinct ways to decode a transcript from the probability distributions calculated by the model for each frame in the input signal
- **The LM dimension**, comparing inferences made with and without post-processing the decoded transcript with a spell-checker as described above

Both dimensions are described in more detail below.

4.3.1 Decoder dimension

In a nutshell, CTC aligns the T_y characters from a known transcription (*label*) with the T_x frames from the input audio signal during training. T_x is typically much larger than T_y and must not be shorter. The characters (*tokens*) in the label must come from an alphabet of size V , which for English are the 26 lowercased ASCII characters $a..z$, the space character and the apostrophe (because this character is very common in

¹⁴CSV file: cv-valid-train.csv

¹⁵sorting the samples is important because...

¹⁶note that the length of the samples could still vary between batches

abbreviations like e.g. "don't" or "isn't"). Additionally, CTC introduces a special token ϵ , called the *blank token*, which can be used to label unknown/silent frames or prevent collapsing (see below). Consequently, the number of characters in the alphabet used by the ASR in this project to recognize English is $V = 26 + 1 + 1 + 1 = 29$.

CTC is *alignment-free*, i.e. it does not require an alignment between the characters of a transcription and the frames of an audio signal. The only thing needed is the audio signal X itself plus its ground truth Y . Each token in the ground truth can be aligned with any number of frames in the input signal. Vice versa, repeated sequences of the same characters can be collapsed, whereas the ϵ token functions acts as a boundary within sequences of a token to prevent collapsing into one, when there should be two (such as in *g-g-o-o- ϵ -o-o-o-o-d-d-d*, which should collapse to *good* and not *god*).

For each frame input signal CTC calculates a probability distribution over the V characters in the alphabet. This yields a $V \times T_x$ probability matrix for the input signal. Because $T_x \gg T_y$, there is usually a vast amount of different valid alignments collapsing to the ground truth. The probability of each valid alignment can now simply be calculated by traversing the probability matrix from left to right and multiplying the probabilities of each character. Because calculating the probability of each valid alignment individually would be too slow and identical prefixes between valid alignments yield identical probabilities, a dynamic approach is usually chosen to calculate the probabilities whereas the intermediate probability for each prefix is saved once computed.

The most probable alignment is calculated by marginalizing (i.e. summing up) over the probabilities of the individual valid alignments. This calculation yields the CTC loss as a sum of products, which is differentiable and can therefore be optimized.

After training, a model using CTC will again output a $V \times T_x$ probability matrix for any previously unseen input. This matrix can be used to infer a transcription, a process also known as *decoding*. The CTC paper proposes two different decoding strategies that are applied before collapsing the characters Graves, Fernández, and Gomez 2006:

- **Best-Path (a.k.a. *greedy*) decoding:** This strategy only ever considers the most likely character at each time step. The transcription before collapsing will be a single path through the the probability matrix, whose probability will be the product of all elements along the path. This approach is easy to implement but does not take into account the fact that a single output can have many alignments, whose individual probability may be lower than the one found with this strategy.
- **Beam-Search decoding:** This strategy approximates the probability of the most probable transcription by following multiple paths simultaneously and only keeping the B most probable paths at each time step. The beam width is a hyperparameter that can be increased to get a more accurate transcription in exchange for higher computational cost.

Usually, Beam-Search decoding performs better than Best-Path decoding. For the sake of completeness, both decoding strategies were compared in this project. This will yield separate learning curves for the decoder dimension. For Beam-Search decoding, the Keras implementation was used, which proposes a default beam width of $B = 100$. This value was not changed.

4.3.2 LM dimension

Using a LM to post-process the inferred transcription with a rudimentary spell checker will not necessarily lead to more accurate transcription, especially if the edit distance between prediction and ground truth is large. Table 3 contains an example where the use of a spell checker is disadvantageous to the quality of a transcription.

In this example, $\langle oento \rangle$ was changed to $\langle onto \rangle$ because this was the most probable word with a maximum edit distance of 2 that was in the vocabulary. Similarly, $\langle appy \rangle$ was changed to $\langle app \rangle$. This

		LER
ground truth	i want to wish you a very happy thanksgiving	
prediction before spell-checking	oento wiceyouepery appy thangkive	0.4318
prediction after spell-checking	onto wiceyouepery app thangkive	0.4545

Table 3 – Example of a transcription whose LER was increased when using a spell checker

lead to a orthographically better sentence, but the LER is higher than without spell-checking.

It is generally expected that post-processing the inference as described above will lead to a lower WER, supposed the LER is already low enough, i.e. the prediction matches the ground truth already pretty well. If the LER value is too high, the spell checker might try too hard to find a word from the vocabulary. This might result in a changed sentence consisting of real words but whose similarity to the ground truth is lower than before the changes. Post-processing might then be counter-productive. Therefore, separate learning curves were plotted for inference with and without post-processing (the *LM dimension*)

4.4 Results and interpretation

As stated before, the most important metrics for this project are the LER and WER, which are measured on validation data and express the quality of the transcriptions produced by the network for unseen data. The loss however is still useful to make statements about at what point the network starts to overfit.

Figure 2 shows the learning curve for the CTC-loss. Obviously the training losses decrease steadily for all amounts of training data, converging to values between 30 (1.000 minutes) and 50 (1 minute). Because the network does not generalize well when being trained on only 1 or 10 minutes of audio data, its predictions are somewhat random. This may be an explanation for the jagged curve of the validation plot (dashed lines) for these amounts of training data. When training on 100 minutes, the plot for the validation loss is smoother, but starts to increase between epoch 10 and 15, meaning the network starts to overfit after that point. When training on 1.000 minutes the validation loss does not decrease after epoch 14 anymore and plateaus at a value of about 90, meaning that any training will not contribute to a more generalizable network.

Figure 3 shows how the average vakzes of the LER over all validation samples develops for the different amounts of training data. The plot on the left shows the results when using best-path decoding, the plot on the right for beam search decoding. The plots for all amounts of training data have been integrated in the same plot for the sake of a clearer representation. Both plots support the conclusions made for the CTC loss in that – except when training on 1.000 minutes of audio – the error rates do not decrease after epoch 15 anymore (in fact there is a slight increase). The plots for the LER when training on 1.000 minutes are almost identical for both decoding strategies. Surprisingly, the LER values continue to decline steadily, although only at a very slow rate, finishing with values of 0.54 (best-path decoding) resp. 0.52 (beam search decoding), meaning that the network got a bit more than half of the characters wrong, at the wrong position or entirely failed to predict them. The values when trying to correct the inferred transcriptions with a spell checker are slightly higher for both decoding strategies (0.55 and 0.53) meaning that post-processing did not help. This nourishes the assumption that a spell checker will probably only lower the WER and only do so if the LER is already low enough (see above).

Finally, figure 4 shows the development of the average WER values over all validation samples. Not surprisingly, the plots oscillate around a value of 1, meaning the network did not get any of the words right. Only when training on 1.000 minutes, the network was able to achieve a value below 1, but is still way higher than the 0.0065 achieved by the Mozilla implementation of *DeepSpeech*. It is noteworthy that the use of a spell checker marginally improves the results here.

In summary it can be said that the best results can be achieved when training on 1.000 minutes of audio,

which comes at no surprise. Training can be stopped after about 15 epochs however, because the results on validation data does not improve from there on. It can also be said that the network manages to get at least half of the characters right, yielding transcriptions from which it is sometimes hardly recognizable what the real transcript might be (especially for shorter sentences). It remains to be seen if this is enough for the following sequence alignment stage.

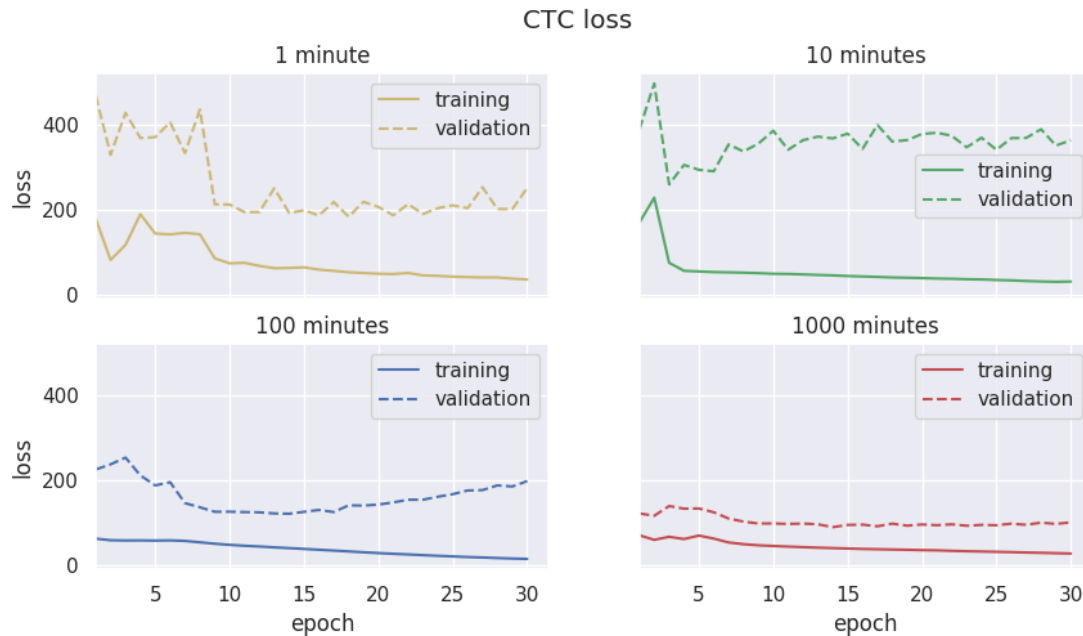


Figure 2 – Learning curve for the CTC-loss while training on 1/10/100/1000 minutes of transcribed audio from the CV corpus using the 5-gram LM provided by the Mozilla implementation of DeepSpeech

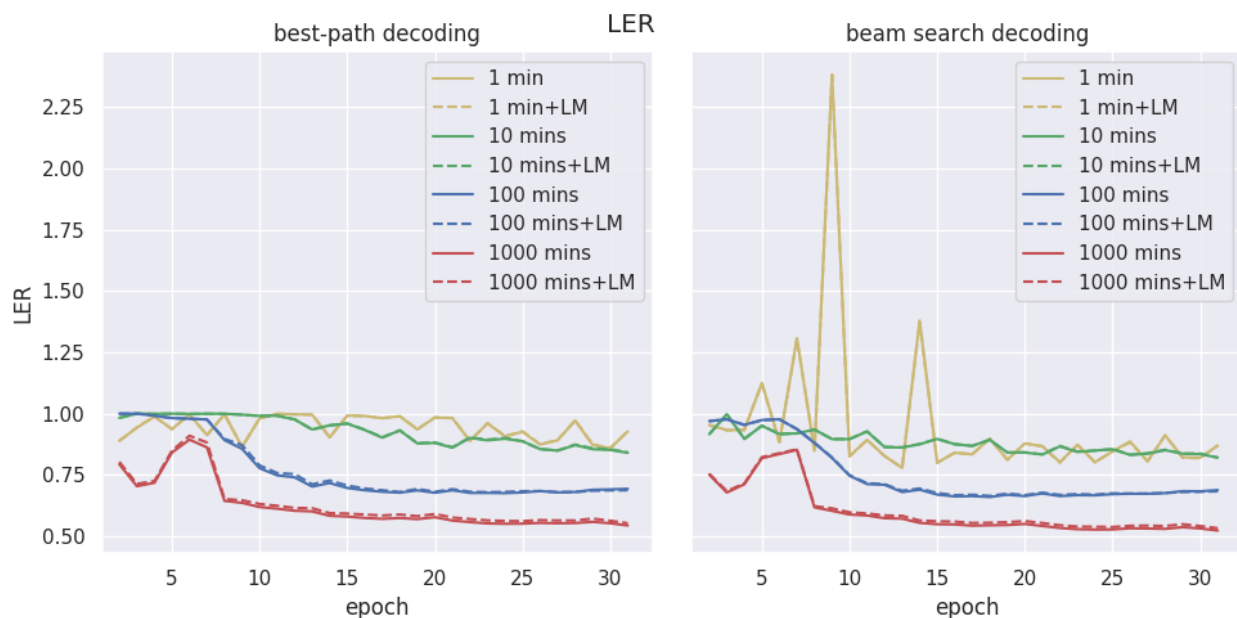


Figure 3 – Learning curve for the LER metric while training on 1/10/100/1000 minutes of transcribed audio from the CV corpus with and without spelling correction with a LM. For the lines where spelling was corrected, the 5-gram LM provided by the Mozilla implementation of DeepSpeech was used.

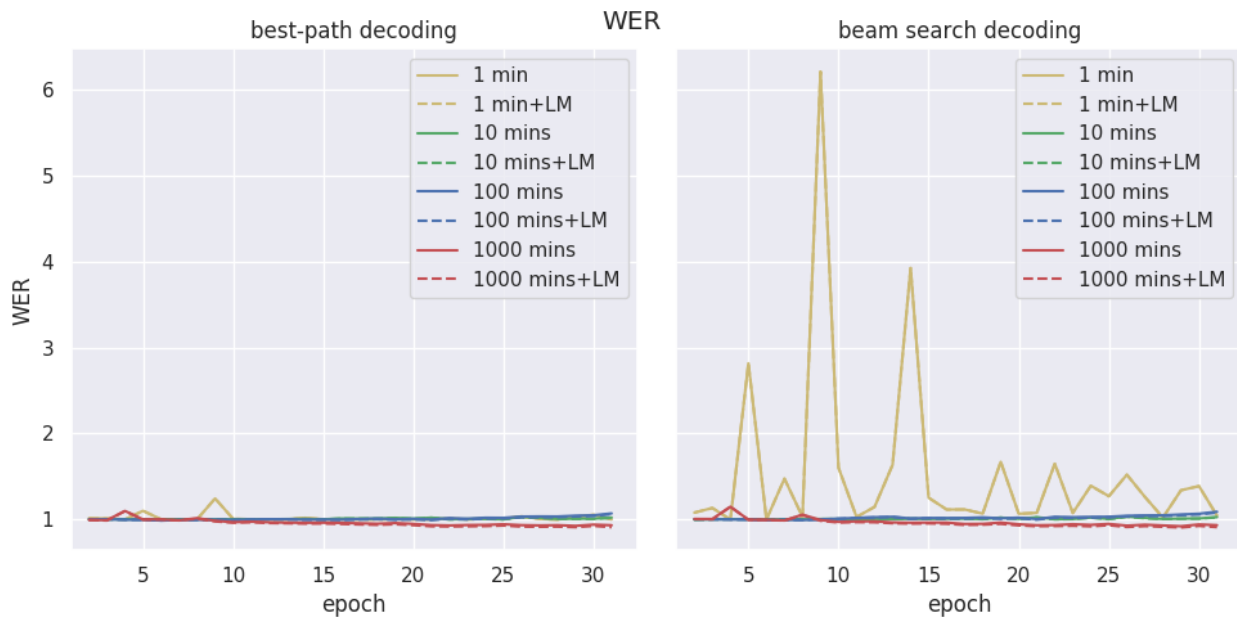


Figure 4 – Learning curve for the WER metric while training on 1/10/100/1000 minutes of transcribed audio from the CV corpus with and without spelling correction with a LM. For the lines where spelling was corrected, the 5-gram LM provided by the Mozilla implementation of DeepSpeech was used.

4.5 Regularization

As an attempt to prevent overfitting (or at least postpone it to later epochs), the network has been regularized by inserting dropouts after each layer. The rate of each dropout has been set to 0.1, meaning a random 10% of the unit weights in each layer will be zeroed out.

4.6 Final thoughts and considerations

Above results were achieved with a spell checker using a vocabulary of 80.000 words and the 5-gram LM from Mozilla. This did not help very much, but it might be possible that a different vocabulary size will produce better results. It is also possible that a different Optimizer, different dropout rates or integrating the LM score into the cost function (like Mozilla did) will produce better results. Finally, it might be fruitful to train on smaller batches as it has been observed that larger batches degrade the quality of a model by Keskar et al. 2016.

All these ideas produces many more combinations to try out, but preparing and running them is very time consuming. Because the LER of about 0.5 (1.000 minutes, no spell checker) looks promising, I decided to leave it at this for the moment and see how far I get.

5 Measuring the performance of the pipeline

Above results reflect the performance of the ASR model alone. To get some insight about the quality of alignments produced by the whole pipeline, a simple web application was implemented that highlights the aligned parts of the transcript as the audio file is being played. This is very useful for an informal review, because the subjective quality of the alignments can be examined interactively. However, this method is not very systematic and infeasible for larger amounts of test data. To get a clearer sense of how well the pipeline performs, steps were taken to run large numbers of previously unseen samples through the pipeline and measure the quality of the final product (the alignments). This section describes how this was done.

5.1 The quality of alignments

Assessing the quality of alignments is not trivial because there is often no reference alignment to compare to. Even if there is one, assessing the quality of an alignment is somewhat subjective because a different alignment does not necessarily need to be worse or better. Quantifying the quality of a result is difficult for an alignment pipeline because there is a massive number of theoretically possible alignments for each audio/text combination. We can however derive a few objective criteria that make up a good alignment:

1. The aligned partial transcripts should not overlap each other
2. The alignments should neither start nor end within word boundaries
3. The aligned partial transcripts should cover the as much of the original transcript as possible
4. The aligned partial transcripts should be at the correct position (i.e. they should cover the actually spoken text)

The first criterion is enforced by changing the type of algorithm used for sequence alignment from a local to a global alignment algorithm. During the IP8 project, the *Smith-Waterman* algorithm was used in the LSA stage, which finds a local optimum for each transcript individually. This was replaced by the *Needle-Wunsch* algorithm, which finds an optimal alignment for all partial transcripts at once.

The second criterion is ensured by adjusting some of the alignments produced by *Needle-Wunsch* so that they fall exactly on word boundaries.

The remaining two criteria can be quantified with the following metrics (note the correlation¹⁷):

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

Table 4 – Metrics to evaluate the quality of alignments

Because the first metric measures how much of the target transcript is covered by the alignments, it is somewhat similar to the Recall R . The second metric measures how well the produced results match up with the underlying parts of the transcript and is therefore similar to Precision P . Both metrics can therefore be reduced to the F-score, a single number usually used to evaluate classification results:

$$F = 2 \cdot \frac{P \cdot R}{P + R}$$

¹⁷positive correlation: higher is better, negative correlation: lower is better

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

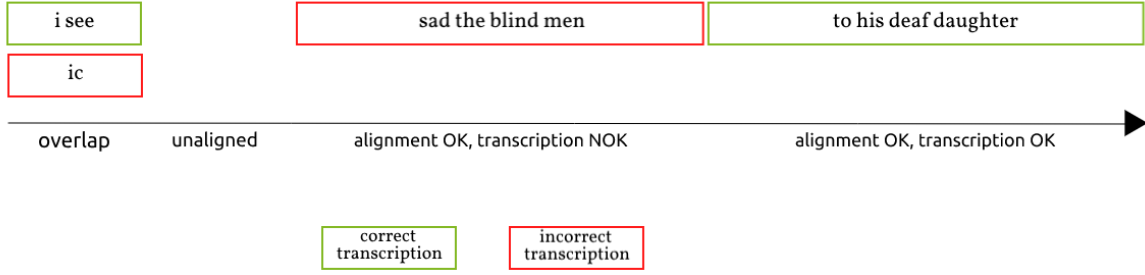


Figure 5 – Example for a partially correct alignment of parts of a sentence

Figure 5 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 $t_1 = \text{i see}$, $t_2 = \text{i c}$, $t_3 = \text{sad the blind men}$, $t_4 = \text{to his blind daughter}$ The metrics can be calculated as follows:

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

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

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

5.2 Test and results

The pipeline was evaluated using the model with the lowest average validation-LER. For English, this was the model with dropouts. To prevent overfitting, training was stopped after x epochs (early stopping). The pre-trained DS model was used as a reference model for comparison.

P , R and F were calculated for English and German by running each sample from the test set of the respective corpus through the pipeline.

5.3 Pipeline performance using an ASR model for a different language

Due to an implementation error the samples from the German test set were run through the pipeline using an ASR for English at some point in the testing phase. The error was corrected, but surprisingly enough, the transcripts produced were still accurate enough for alignment. Apparently, the Needle-Wunsch algorithm used in the global alignment stage only needs very little resemblance of the partial transcripts to the real transcripts in order to produce alignments that are usable, although maybe a bit worse.

6 Forced Alignment for other languages

So far, only audio and transcripts in English were considered. A fully automated solution however should be able to align text and audio in any other language. Because of linguistic characteristics like sound patterns and morphology the results might vary a lot between languages when tested underwise identical circumstances. To get some intuition about the influence of language and whether above conclusions are transferable to other languages, the pipeline was evaluated on the German samples received from *ReadyLingua*.

6.1 Inferring German transcripts

Enabling the pipeline to handle German samples means training a German ASR as its core element. This required minimal modifications to the network architecture, because German transcripts use a different alphabet. As mentioned before, the apostrophe is far less common in German than in English and was therefore dropped. On the other hand, umlauts are very common in German and were added to the 26 ASCII characters. Since the alphabet represents all possible labels, the output layer consisted of 31 units (one for each character in the alphabet, the three umlauts, space plus a blank token) instead of the 29 units used for English.

Training an ASR model for German and plotting a learning curve also required amounts of training data on a similar scale like the CV corpus used for English. Since at the time of this writing, the CV was still a work in progress, datasets for languages other than English were not available. High-Quality ASR corpora are generally hard to find, especially considering the number of samples needed to train a RNN. There are corpora for ASR in German, but those are either not freely available or their quality is unknown. An extensive list of German corpora for various purposes can be found at the Bavarian Archive for Speech Signals (BAS)¹⁸. Some of the corpora on this list are free for scientific usage. However, not all of these corpora are targeted at ASR and their quality is often unknown.

6.2 Data augmentation

Integrating new raw data means preprocessing the audio (e.g. resampling) and the text (e.g. normalization, tokenization) to make sure it exhibits the same properties as the other corpora and the data conforms to the expected format. This step is usually very time consuming, often using up most of the project time. Because no ASR corpus for German was readily available, training was done on the data received from *ReadyLingua* as a start and see how far we get. The alignment between audio and transcript in this corpus was done manually and is therefore very accurate. Audio and text were already preprocessed in the IP8 project and the metadata was processed and stored as a corpus. The individual training samples could therefore be transformed to the expected format with comparably little effort. Also, the samples exhibited similar properties (average audio and transcript length) like the CV corpus (refer to table 2). However, the total length of the samples in the training set was only about one and a half hours, which was much less than the 1000+ minutes of the CV corpus and certainly not enough for the 1.000 minutes needed to plot a learning curve like to the one made for English.

An easy way to get more training data is augmenting existing data by synthesizing new data from it. This is particularly easy for audio data, which can be distorted in order to get new samples corresponding to the same transcript. The following distortions were applied in isolation to each sample in the training set:

- **Shifting:** The frames in the input signal were zero-padded with samples corresponding to a random value between 0.5 and 1.5 seconds, shifting the signal to the right, i.e. starting the signal later. This resulted in one additional synthesized sample for each original sample. Shifting to the left was not done to prevent cropping parts of the speech signal.
- **Echo:** The presence of echo can be generated with the Pysndfx library¹⁹ using random values for delay

¹⁸<https://www.phonetik.uni-muenchen.de/Bas/BasKorporaeng.html>

¹⁹<https://github.com/carlthome/python-audio-effects>

and damping. This resulted in one additional sample.

- **Pitch:** The pitch of the signal was increased or decreased. Increasing and decreasing was done using two different random factors, resulting in two additional samples. This can be seen as a rudimentary way to simulate a female from a male speaker or vice versa.
- **Speed:** Faster or slower speaking rates can be simulated by "stretching" or "compressing" the signal while preserving the pitch. Similar to the change in pitch, two different random factors were used to change the tempo. This resulted in two additional samples.
- **Volume:** The loudness of the speaker was artificially reduced or increased by a random value within the range of $[-15.. -5]$ resp. $[5..15]$ db. This resulted in two additional samples.

With above methods eight synthesized samples can be created for each original sample from the corpus. It turned out however that this was still not enough to plot a learning curve. To augment the data to the 1.000 minutes needed, additional samples were created using a random combination of the distortions. The random parameters differed from the ones used before to prevent overfitting to the distortion. Table 5 shows the corpus statistics before and after data augmentation.

	total audio length	# samples	Ø sample length (seconds)
before augmentation	1 : 36 : 09	1,700	2.89
after augmentation	16 : 40 : 00	18,955	3.16

Table 5 – Comparison of RL corpus before and after data augmentation (training set only)

6.3 Creating a Language Model for German

Since the ASR stage in the pipeline uses a spell-checker querying a LM to post-process the results a 5-gram model similar to the one created by Mozilla needed to be trained first.

6.4 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. The probability of a sentence can be computed from the probabilities of each word (1-grams) with given all its preceding words in the sentence using conditional probability. Getting statistically relevant high numbers for each combination of words requires huge text corpora. However, language is dynamic and new sentences can be created all the time so that no corpus would be big enough. To handle this, n -grams approximate the probability of a combination of words by only considering the history of the last n words (n denoting the order). However, above problem is still valid for n -grams of any order: Because of combinatorial explosion n -grams suffer from sparsity with increasing order.

6.4.1 Perplexity, discount and smoothing

To evaluate an n -gram LM a metric called *perplexity* is usually used, which is the normalized inverse probability on a test set. The perplexity can be interpreted as the grade to which the LM is "confused" by a certain n -gram. A high perplexity therefore corresponds to a low probability. Since the perplexity carries the probability of a certain n -gram in the denominator, the perplexity for OOV- n -grams cannot be calculated (division by zero). To handle this 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 for n -grams that do not appear in the training corpus. Smoothing therefore shaves off a bit of the probability mass from the known n -grams and moves it to the unknown n -grams. The factor with which the probability of a known n -gram is reduced is called *discount*.

6.4.2 Kneser-Ney Smoothing

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 because it discounts frequent n -grams too much (i.e. gives too much probability to unknown n -grams). A better way of smoothing is achieved using *Kneser-Ney Smoothing*. For unseen n -grams, *Kneser-Ney Smoothing* estimates the probability of a particular word w being the continuation of a context based on the number of context it has appeared in the training corpus. For any previously unseen n -gram, a word that appears in only few contexts (e.g. the word *Kong*, which only follows the words *King* or *Hong* in most corpora) will yield a lower probability than a word that has appeared in many contexts, even if the word itself may be very frequent. The intuition behind this is that such a word is assumed less likely to be the novel continuation for any new n -gram than a word that has already proved to be the continuation of many n -grams.

6.5 Creating a raw text corpus

To train a n -gram model for German, a raw text corpus of German Wikipedia articles was used as corpus. Like the English n -gram from Mozilla KenLM (Heafield 2011) was used to estimate the probabilities. The articles were pre-processed to meet the requirements of *KenLM*. It was normalized as follows

- remove Wiki markup
- remove punctuation
- make everything lowercase
- **Unidecoding**: translate accentuated characters (like è, é, ê, etc.) and special characters (like the German ß) to their most similar ASCII-equivalent (e resp. ss). This process helps accounting for ambiguous spelling variants of the same word and misspelled words. It also reduces the number of unique words by reducing different versions to a normalized variant. A special case are umlauts. Although also not part of the ASCII code set, they were kept as-is because they are very common in German.
- **Tokenization**: 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).
- **Numeric tokens**: Word tokens that are purely numeric (such as year numbers) are replaced with the special token <num>. Although such tokens occur frequently in the Wikipedia articles, they are unwanted in the corpus because they represent values and do not carry any semantic meaning. Because there is a infinite number of possible numeric tokens, they were all collapsed to the same normalized token.

The corpus was saved as text file containing one normalized sentence per line. The special tokens <s> and </s> are used to mark beginnings and endings of sentences as well as the <unk> token which is traditionally used to represent OOV words. They are however not part of the corpus because they are added automatically by *KenLM*.

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 Recognition²⁰)

```

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

Like for the English spell checker, three vocabularies containing the 40.000, 80.000 and 120.000 most frequent words from the corpus was created. The words from these vocabularies make up 87.75%, 90.86% resp. 93.36% of the total number of words in the corpus. It is expected that the optimal number of words in the vocabulary is higher for German than for English. This is due to the fact that different flexions of the same word are very common in German due to grammatical conjugations (different forms for the same verb) and declinations (different cases for the same noun). Therefore German tends to apply a wider range of words and the size of vocabulary had to be increased. Handling the different flexions would require lemmatization and/or stemming the corpus in order to reduce them to a common base form. This has not been done for simplicity and time constraints. It is also doubtful whether this would actually help improving the quality of inferred transcripts, since humans do not speak in lemmata or stems.

6.6 Training the LM

The final corpus contained data from 2,221,101 Wikipedia articles (42,229,452 sentences, 712,167,726 words, 8,341,157 unique words). This corpus was used to train a 5-gram LM using *KenLM*. *KenLM* uses *Kneser-Ney Smoothing* and some optimization techniques called *quantization* and *pointer compression*.

6.6.1 Data structures

n -grams can be represented with a prefix-tree structure (called *Trie*)²¹, which allows for pruning. n -grams of order 2 and higher can be pruned by setting a threshold value for each order. n -grams whose frequency is below the threshold will be discarded. *KenLM* does not support unigram pruning.

6.6.2 Quantization

To save memory, the amount of bits used to store the non-negative log-probabilities can be reduced with the parameter q to as little as $q = 2$ bits at the expense of accuracy. This reduction yields $2^q - 1$ possible bins. The value of each bin is calculated by equally distributing the probabilities over these bins and computing the average. Note that the quantization is done separately for each order and unigram probabilities are not quantized.

6.6.3 Pointer Compression

To use memory even more efficiently, the pointers which are used to store n -grams and their probabilities can be compressed. Such pointers are used to represent e.g. word IDs (for q -grams) and are stored as sorted integer-arrays. Additionally, These integers can be compressed using a lossless technique from Raj

²⁰<https://de.wikipedia.org/wiki/Spracherkennung>

²¹note 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

and Whittaker 2003 by removing leading bits from the pointers and store them implicitly into a table of offsets. The parameter a controls the maximum number of bits to remove. There is a time-space trade-off meaning that a higher value of a will lead to a smaller memory footprint at the cost of a slower training time.

6.6.4 Building the model

The a 5-gram LM was trained on the German Wikipedia corpus using the Trie data structure and the same parameters like the model downloaded from *DeepSpeech* ($q = 8$ and $a = 255$). Like the *DeepSpeech* model 4- and 5-grams were pruned by setting a minimum frequency of 1.

6.7 Evaluating the LM

Literature suggests two methods to evaluate a LM: Extrinsic and intrinsic evaluation.

6.7.1 Extrinsic and intrinsic evaluation

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, models can only be compared if they use the same vocabulary and always encode characteristics of the training corpus (Jurafsky and Martin 2019). Since the sentences in a corpus of legal documents use different structures and word distributions than a corpus of children's books, two models trained on these corpora will not be comparable. Evaluating the created German Wikipedia corpus intrinsically would therefore require training a reference model on the same corpus, which can become very time consuming.

6.7.2 Evaluation of KenLM

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). The highly optimized performance was a big advantage especially for this project because it enabled testing the model on a local machine. The probabilistic performance of *KenLM* has been evaluated by training a 5-gram model on a 126 billion token corpus (393 million 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.

6.7.3 Evaluation of the German Wikipedia LM

Because of time constraints and because *KenLM* has already been extensively evaluated on English I resigned from evaluating my German LM intrinsically, even though the corpus used for training is not as big as the one used in Heafield et al. 2013. *KenLM* is to date widely recognized as the best performing LM available, which is also 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, two different experiments were made:

- **Experiment 1:** The probability calculated for valid German sentences was compared against variants of the same sentences with the words in randomized order.

- **Experiment 2:** The LM was used together with its vocabulary to build a simple word predictor.

Both experiments are explained in more depth below.

6.7.4 Evaluation 1: Comparing scores of randomized sentences

The first experiment tests the validity of the probabilities (*scores*) calculated by the LM. For this, an arbitrary choice of 5 valid sentences in German was used. To ensure the sentences could not have been seen during training, the following 5 sentences were taken from a newspaper printed after the creation of the Wikipedia dump:

```

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

```

Listing 2 – Representation in corpus

All sentences have been normalized the same way sentences were preprocessed for training. For each of them the score was calculated. Then the words were shuffled and the score was calculated again. A good LM should calculate a (much) higher probability for the original sentence, because the shuffled sentence is most likely just gibberish. Table 6 shows the results of the comparison. It is evident that the probabilities for the shuffled sentences are much lower than for the sentences where the words appear in the correct order. The probabilities calculated by the LM are therefore deemed valid.

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 6 – Comparison of log10-probabilities calculated for news sentences and a permutation of their words

6.7.5 Experiment 2: Word predictor

The second experiment tests whether the trained LM is able to continue a sentence given its beginning. For this each word from the vocabulary is appended and the score of the resulting stumps is calculated. The most likely continuation can be estimated by sorting the resulting list in descending order (the probabilities are $\log_1 0$ -based, i.e. negative) and taking the first element. This behavior can be applied iteratively to construct a sentence from a stump. For this experiment a sentence was started with the stump « *Ein 2007 erschienenenes* ». Afterwards a word from the five most probable continuations was appended. The extended stump was then again fed into the LM. This process was repeated until some kind of sentence ending was encountered. Each extended stump was preprocessed the same way the sentences were preprocessed for training (lowercasing, replacing numbers with <num>, etc.). Figure 6 shows the path taken through the predictions. Note that the predictions for the second and third word of the stump after typing the first word are shown in grey for illustrative purposes, although they were not considered for continuation.

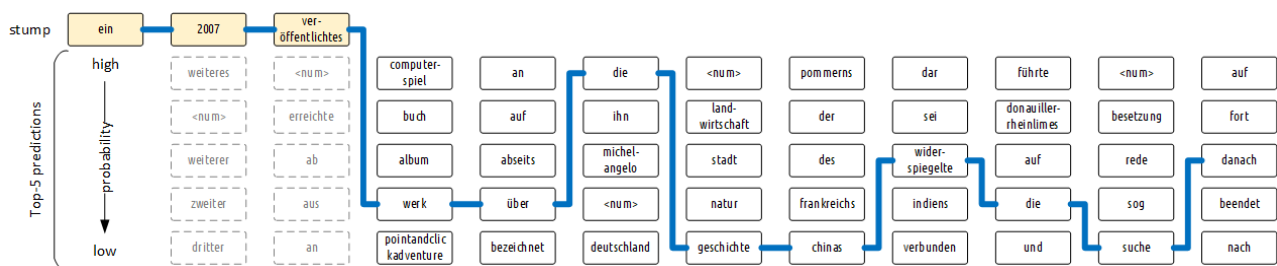


Figure 6 – Word predictions of the trained 5-gram model for continuations of the stump « *Ein 2007 erschienenenes ...* ». The blue path represents a grammatically valid German sentence.

Although prediction was slow we can observe that the words suggested by the LM are generally grammatically correct continuations and often make sense, although the probability for some of the predicted words (like *Michelangelo*) is sometimes unexplicably high. Nevertheless it was possible to create a valid German sentence from the stump using only the suggested words. The LM even seems to have captured some notion about grammatical concepts like German cases (e.g. that « *die Geschichte Chinas* » is more likely than « *die Geschichte China* »). On the other hand we can observe that the meaningfulness of the suggestions decreases with the progress because some long-distance relationships between words are lost for small values of n .

6.8 Results

The following figures show the learning curve for training on the German with synthetisation.

7 Conclusion

Was man auch noch machen könnte:

- Zusätzliche Daten/Korpora: Bavarian ARchive for Speech Signals (<http://www.bas.uni-muenchen.de/forschung/Bas/BasK>)
- Hunspell Checker: <http://hunspell.github.io/> –; Python modul: <https://github.com/blatinier/pyhunspell> –;
- Dictionaries gibt's hier <https://github.com/woorm/dictionaries> - Transfer Learning (müsste man aber zuerst ein geeignetes Modell finden) und Layer freeze –; Layer freeze

List of Figures

1	Architecture of the simplified model. The cell type and the activation function is indicated in brackets for each layer (FC=Fully-Connected)	6
2	Learning curve for the CTC-loss while training on 1/10/100/1000 minutes of transcribed audio from the CV corpus using the 5-gram LM provided by the Mozilla implementation of <i>DeepSpeech</i>	14
3	Learning curve for the LER metric while training on 1/10/100/1000 minutes of transcribed audio from the CV corpus with and without spelling correction with a LM. For the lines where spelling was corrected, the 5-gram LM provided by the Mozilla implementation of <i>DeepSpeech</i> was used.	14
4	Learning curve for the WER metric while training on 1/10/100/1000 minutes of transcribed audio from the CV corpus with and without spelling correction with a LM. For the lines where spelling was corrected, the 5-gram LM provided by the Mozilla implementation of <i>DeepSpeech</i> was used.	15
5	Example for a partially correct alignment of parts of a sentence	17
6	Word predictions of the trained 5-gram model for continuations of the stump « <i>Ein 2007 erschienenenes</i> ... ». The blue path represents a grammatically valid German sentence.	25

List of Tables

1	Example for how a LM can help improve the quality of an inferred transcription by changing characters and words. Audio and ground truth (GT) were taken from the <i>ReadyLingua</i> corpus and the inference was made with the pre-trained <i>DeepSpeech</i> model.	7
2	Statistics about corpora that were available for training. The sample length is given in seconds, the transcript length as the number of characters.	11
3	Example of a transcription whose LER was increased when using a spell checker	13
4	Metrics to evaluate the quality of alignments	16
5	Comparison of RL corpus before and after data augmentation (training set only)	20
6	Comparison of log10-probabilities calculated for news sentences and a permutation of their words	24

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://khefield.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://khefield.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/>.
- Keskar, Nitish Shirish et al. (2016). "On Large-Batch Training for Deep Learning: Generalization Gap and Sharp Minima". In: *CoRR* abs/1609.04836. arXiv: 1609.04836. URL: <http://arxiv.org/abs/1609.04836>.
- Kunze, Julius et al. (2017). "Transfer Learning for Speech Recognition on a Budget". In: *CoRR* abs/1706.00290. arXiv: 1706.00290. URL: <http://arxiv.org/abs/1706.00290>.
- 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).
- Raj, B. and E. W. D. Whittaker (Apr. 2003). "Lossless compression of language model structure and word identifiers". In: *2003 IEEE International Conference on Acoustics, Speech, and Signal Processing, 2003. Proceedings. (ICASSP '03)*. Vol. 1, pp. I–I. DOI: 10.1109/ICASSP.2003.1198799.

Acronyms used in this document

ASR	Automatic Speech Recognition
BAS	Bavarian Archive for Speech Signals
CTC	Connectionist Temporal Classification
CV	CommonVoice
DS	Deep Speech
E2E	end-to-end
FA	Forced Alignment
FHNW	University of Applied Sciences
GCS	Google Cloud Speech
GPU	Graphics Processing Unit
GRU	Gated Recurrent Unit
GSA	Global Sequence Alignment
LER	Label Error Rate
LM	Language Model
LS	LibriSpeech
LSTM	Long Short Term Memory
LSA	Local Sequence Alignment
MFCC	Mel-Frequency Cepstral Coefficients
NN	Neural Network
RL	ReadyLingua
RNN	Recurrent Neural Network
SGD	Stochastic Gradient Descent
STT	Speech-To-Text
OOV	Out Of Vocabulary
SRILM	the SRI Language Modelling Toolkit
SW	Smith Waterman
VAD	Voice Activity Detection
WER	Word Error Rate

The simple spell checker in detail

- 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.

$$P = \{(w_j, w_{j+1}) \mid w_j \in C_j \wedge w_{j+1} \in C_{j+1}\},$$

$$C_j \in \{C_i^0, C_i^1, C_i^2, C_i^{>2}\}$$

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

8 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, November 22, 2018

Daniel Tiefenauer