

Data Selection for the Training of Deep Neural Network in the Framework of Automatic Speech Recognition



Juan Karsten

Supervisor: Dominique Fohr and Irina Illina

Multispeech team, Inria Loria Nancy

Universite de Lorraine

A thesis submitted for the degree of

Master of Computer Science

2017

Acknowledgements

Firstly I would like to thank my supervisors: Dominique Fohr and Irina Illina for their amazing supports and knowledge while I was doing the internship. In addition, I highly appreciated the help of Multispeech team, especially Ismael Bada for helping me out tremendously in related to the script and Imran Sheikh for giving me overall knowledge of speech recognition.

I never forget to thank my family, my friends, and Feby Rara for giving me unlimited supports and love. Without them, this thesis will not be possible.

Abstract

Training a speech recognition system needs audio data and their corresponding exact transcriptions. However, manual captioning is expensive, labor intensive, and error-prone. Some sources, such as: TV broadcast, have closed captions by default. Closed captions are closed to the exact transcription, but not exactly the same. Building automatic speech recognition from the dataset may result in a bad model. Therefore, selecting data is important to train a highly performance model. In this work, we explore the lightly supervised approach which is an iterative process to select good data. We built a baseline model and proposed new approaches.

For the baseline model, we select the data based on phone matched error rate and average word duration. The process to build the model runs iteratively by training different acoustic models and selecting data generated by the acoustic models. The model is evaluated by recognizing a development set. The development set comprises of audio data and their corresponding manually annotated transcriptions. The comparison between the decoding results on the development set and the manually annotated transcriptions produces word error rate which is the metric to measure how good the model is. We stop the data selection iterative process when the word error rate stops decreasing.

We proposed three new models: averaging word matched error rate and average word duration of 3 speech recognizers, combining 3 recognizers by our proposed algorithm, and utilizing random training set in each iteration. The error rate decreases in the combination proposed model as well as utilizing random training set. However, the proposed models perform slightly poorer in the next iteration.

Keyword: automatic speech recognition, light supervised approach, error rate.

Contents

1	Introduction	1
1.1	Motivation	1
1.2	Contributions	4
1.3	Outlines of the thesis	4
2	Background	5
2.1	Automatic speech recognition	5
2.2	N-Gram language model	6
2.3	Acoustic model	7
2.3.1	Markov chain	7
2.3.1.1	Hidden markov model	8
2.3.2	Using HMM to acoustic modeling	8
2.4	Deep learning for acoustic model	10
2.4.0.1	Feedforward neural network	10
2.4.1	Time delay neural network(TDNN)	11
2.4.2	Hybrid acoustic model	13
3	Methodology	16
3.1	Lightly supervised data selection approach	16
3.2	State-of-the-art in data selection	18
3.3	Data	19
3.4	Baseline data selection method	19
3.4.1	Building the speech recognition systems	20
3.4.2	Metric	21
3.4.3	Data selection	23
3.5	New proposed data selection methods	25

4	Implementation	28
4.1	Data set	28
4.2	Experimental setup	29
4.2.1	Language model	30
4.2.1.1	Library: SRILM	30
4.2.1.2	Implementation	30
4.2.2	Acoustic model	32
4.2.2.1	Library: Kaldi	32
4.2.2.2	Acoustic model implementation	32
4.2.3	Recognition	33
4.2.3.1	Data Selection	34
4.2.3.2	Evaluation	34
4.3	Experiment resources	35
5	Experiment result	36
5.1	Baseline model	36
5.2	Proposed models	38
5.3	Execution time	40
5.4	Discussion	40
6	Conclusion and future works	43
6.1	Conclusion	43
6.2	Future works	44
6.2.1	Short term	44
6.2.2	Long term	45
	Bibliography	46

List of Figures

1.1	Automatic speech recognition [Vroegop, 2015]	1
2.1	A triphone HMM model consisting of three emitting states(beginning(S0), middle(S1), and end(S2) of a phone), a transition matrix A , and emission probability(observation likelihood) B [Silicon Intelligence, 2017] .	9
2.2	A perceptron [Nielsen, 2015]	10
2.3	Multi layer perceptron [Nielsen, 2015]	11
2.4	An example of a TDNN architecture. [Peddinti et al., 2015]	12
3.1	Training and decoding with the speech recognition system	20
3.2	An illustration of lattice [Gales and Young, 2007]	21
3.3	Data selection pipeline	23
3.4	Evaluation pipeline	25
4.1	The example	34
5.1	Average word duration of data selection iterations	38
5.2	Phone matched error rate of the iteration	39

Chapter 1

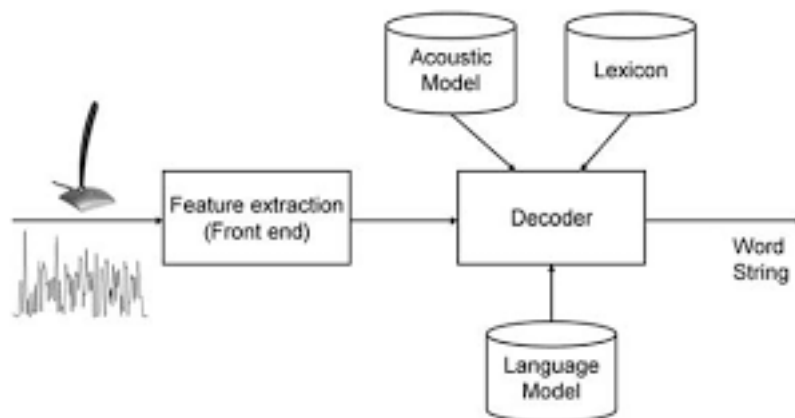
Introduction

1.1 Motivation

Automatic speech recognition(ASR) is one of the subfield of natural language processing(NLP) with many practical applications. Some examples of practical applications in speech recognition; are: automatic closed captioning for hearing-disabled persons [Patel and Rao, 2010], taking notes of conversations between doctors and patients [Klann and Szolovits, 2008], automatic closed captioning TV broadcasts [Woodland et al., 2015], and many more. This internship will focus on the application of automatic transcription on television broadcasts.

closed captions

Figure 1.1: Automatic speech recognition [Vroegop, 2015]



Picture 1.1 shows an illustration how the automatic speech recognition works. ASR receives acoustic speech inputs and outputs a word string. An automatic speech recognition is composed of three components: acoustic model, language model, and lexicon. Acoustic model is a statistical model which represents a relationship between

an acoustic input and its corresponding phonemes. Phoneme is a unit (or a group) of sound, for example: /k/, /l/, /sh/, etc. Lexicon is a dictionary which maps from words to their phonemes. Language model computes the probability of a word string. Hence, the lexicon is used by the acoustic model and the language model together. These three components work together to search the most likely word sentence (from all possible sentences) which represents what people said in the audio.

Despite of the rapid development of speech recognition, there are still many challenges in the field. One of the challenges is training a speech recognizer or adapting a speech recognizer to a new task which requires a massive amount of *transcribed* training data. The *transcribed* training data is audio data which has corresponding transcriptions what people said in the audio data. However, transcribing audio manually is labor intensive and also time consuming. There exist many unlimited supply of training data sources from internet, TV broadcasts, radio, as well as video streaming websites, like Youtube. Nevertheless, many of them only possess few corresponding transcriptions or no transcription at all. To build a speech recognition, we can use a training dataset from TV broadcasts which usually have *closed captions*. *closed captions* are close, but not exactly the same as what people said. Furthermore, *closed captions* are often badly aligned with the audio. These are some examples of good and bad *closed captions*, as well as noisy utterance.

1. An example of good *closed caption*. This segment/utterance has the same *closed caption* as the detailed transcription.
 - Real transcription: That is the very big question that leads to another big question. We're definitely thinking about putting an offer. That's great.
 - *Closed caption*: That is the very big question that leads to another big question. We're definitely thinking about putting an offer. That's great.
2. An example of bad *closed caption*. The *closed caption* is slightly different compared to the real transcription.
 - Real transcription: Russia started badly with the dropping at the hands of Spain. But, they got better and better. Spain looked unstoppable to start with but since then they have looked a little.
 - *Closed caption*: Russia started badly with at beating at the hands of Spain. Spain looked then they have looked a little

3. Noisy segment. Real transcription: Some segments in training audio may have background noise which makes it hard to be recognized by the ASR. Even, it is sometimes difficult to be recognized by humans.

[Clapping in the background] John Higgins developed very well. He is not five in front.[Still continues clapping in background]

This internship explored how to build a high performance speech recognition system by utilizing less supervision because we did not have correct transcription audio data. The main idea is to use an automatic speech recognizer to automatically transcribe audio data which can be leveraged as a training dataset. The steps of the main idea are as following: firstly, train on all data which are automatically *transcribed* by other ASR systems, recognize these data with the ASR which we trained, and then compare the decoding transcriptions with the *closed captions*. Finally, remove all segments which do not agree and train only the segments which agree.

The idea of using *untranscribed* data(or *unsupervised* training) has been proposed by Zavaliagkos and Colthurst in 1998 [Zavaliagkos and Colthurst, 1998], Kemp and Waibel in 1999 [Kemp and Waibel, 1999]. However, they only utilized small amount of data. Lamel et al were the first who proposed lightly supervised training data with a large amount of training data [Lamel et al., 2002]. In place of utilizing *untranscribed* data, they trained speech recognition systems by using audio data with *closed captions*. The process is called lightly supervised approach.

Lightly supervised approach is one of the technique to select "good" data. An acoustic model from another task(or another corpus) is utilized to recognize audio data. The decoding results are compared with the *closed captions* and removed if they disagree. These selected data are trained to generate a new acoustic model which is leveraged to recognize the same or different data. These decoding results are selected again for training a new acoustic model. This internship mainly focuses on the acoustic model component.

The first goal of this internship is to study and evaluate the state-of-the-art techniques to select the "good" training data. The second goal is to propose new solutions to select data. The outcome of this internship will be useful as a training dataset to build a high performance automatic speech recognition.

This internship is expected to contribute in solving data selection problem. Internet has infinitely many audio data, for example: podcast, internet radio, tv broadcast, as well as video streaming. But, mostly they do not have transcriptions at all. In

this internship, we conceived speech recognition systems by utilizing TV broadcasts and their corresponding *closed captions*.

1.2 Contributions

We make several contributions as following:

1. We give some backgrounds related to speech recognition which helps the readers to understand how the speech recognition works in the nutshell. This background helps to grasp the idea of data selection.
2. We did some literature reviews and proposed some data selection methods.
3. We implemented and evaluated some state-of-the-art data selection techniques.
4. We implemented and evaluated new approaches in data selection and assessed these approaches against the current state-of-the-art approach.

1.3 Outlines of the thesis

This report is organized as following. In chapter 2, we elaborate some backgrounds in the automatic speech recognition. The chapter gives necessary backgrounds to the readers how speech recognition works; moreover, the data selection method is further explained. Chapter 3 elaborates the data selection technique which was explored and experimented in this internship. Furthermore, the new data selection methods are also proposed. Chapter 4 tells the setup of experiment and evaluation measurement. Then, chapter 5 shows the result of experiments. Lastly, the report is concluded by chapter 6 which conveys conclusion and some suggestions for future works.

Chapter 2

Background

2.1 Automatic speech recognition

The goal of the automatic speech recognition(ASR) is to search the most likely sentence(\hat{W}) in a language L given acoustic observation O [Jurafsky and Martin, 2009].

Observation O can be seen as a sequence of sliced acoustic signal. Each successive index represents the consecutive slice of acoustic input.

$$O = o_1, o_2, o_3, \dots, o_t \quad (2.1)$$

A sentence is composed of a sequence of words.

$$W = w_1, w_2, w_3, \dots, w_n \quad (2.2)$$

Thus, an automatic speech recognition can be formalized as following:

$$\hat{W} = \operatorname{argmax}_{W \in L} P(W|O) \quad (2.3)$$

Applying the bayes rule to the equation 2.3 produces the equation bellow. Because observations O do not change for every possible sentence, we can assume $P(O)$ as a constant and ignore it.

$$\hat{W} = \operatorname{argmax}_{W \in L} \frac{P(O|W) \times P(W)}{P(O)} \quad (2.4)$$

$$= \operatorname{argmax}_{W \in L} P(O|W) \times P(W) \quad (2.5)$$

where \hat{W} is the most likely sentence, $P(W)$ is the prior probability of sentence W computed by the language model, and $P(O|W)$ is the observation likelihood calculated by the acoustic model. This internship mainly focuses on the acoustic model. We will discuss the language and acoustic model more in the next subsections.

2.2 N-Gram language model

A language model is a statistical model which computes the probability of word sequence. One of the most common language model is N-gram language model, which we leverage in this internship. The foundation mathematics of N-gram language model was proposed by Markov in 1913. Then, Shannon re-introduced N-gram to compute the likelihoods of English word sequences [Shannon, 2001]. N-gram language model has been applied in many applications, such as: speech recognition [Woodland et al., 2015], handwriting recognition [Poznanski and Wolf, 2016], as well as spelling correction [Kukich, 1992].

The goal of language model is to compute the probability ($P(w|h)$) of a word w given previous word history h . The way to compute this probability is by using a large corpus and counting the relative frequencies. For example: to count the word w_1 given history: w_2, w_3, \dots, w_n , the probability is as following:

$$P(w_1|w_2w_3\dots w_n) = \frac{C(w_1w_2w_3\dots w_n)}{C(w_2w_3\dots w_n)} \quad (2.6)$$

where $C(w_1w_2w_3\dots w_n)$ is the frequency of the sequence $w_1w_2w_3\dots w_n$.

Because we can create infinite number of possible sentences, it is more efficient to compute probability by applying the chain rule of probability.

$$\begin{aligned} P(w_1^n) &= P(w_1)P(w_2|w_1)P(w_3|w_1^2)\dots P(w_n|w_1^{n-1}) \\ &= \prod_{k=1}^n P(w_k|w_1^{k-1}) \end{aligned}$$

where w_1^{n-1} represents a sequence of words $w_1w_2\dots w_{n-1}$.

However, the probability of a word given preceding words is difficult to compute. Instead of computing the entire preceding history, N-gram only approximates the probability with last few words. For example: the bigram LM approximates the probability of a word given history by the probability of a word given the previous word only $P(w_n|w_{n-1})$. This approximation(or assumption) is called Markov assumption.

A major problem exists in N-Gram language model. Information inside a corpus might be limited. Some word sequences probably do not exist in the corpus, but they are perfectly acceptable word sequences. If we calculate the N-Gram likelihood for these sequences, we will get zero value. The solutions to the problem are *backoff* and *smoothing*. *Backoff* allows the N-gram LM to use the value of lower N-gram LM if we have zero probability in the higher N-gram. *Smoothing* is a technique to avoid zero

probability for a word sequence. If a word sequence does not exist in the corpus, the LM still produces a non-zero value even though it is small.

To evaluate the performance of a language model(LM), perplexity is used. The intuition behind perplexity is given two N-gram LMs, the better model is the one which predicts test set better than the other. Lower perplexity means that the model is better. Given a language model and a test corpus, perplexity is calculated as following [Bahl et al., 1983, Klakow and Peters, 2002]:

$$PP = \exp\left(-\frac{1}{N} \sum_{i=1}^N \log P(w_i|h_i)\right) \quad (2.7)$$

where N is the number of words in test corpus, w_i is the i -th word, and h_i is the history of word w_i .

2.3 Acoustic model

The acoustic model(AM) is commonly based on *Hidden Markov Model*(HMM). Thus, before talking about the acoustic model deeply, it is better to elaborate HMM first.

2.3.1 Markov chain

A weighted state finite automaton is a finite state automaton where the arc between nodes has probability how likely the path is taken. Markov chain is one kind of weighted state finite automata which the input sequence determines which state the automaton will go.

Markov chain can be defined formally as following:

- $Q = q_1 q_2 \dots q_N$ a set of states

- $A = \begin{bmatrix} a_{01} & a_{02} & a_{03} & \dots & a_{0n} \\ a_{11} & a_{12} & a_{13} & \dots & a_{1n} \\ \dots & & & & \\ a_{n1} & a_{n2} & a_{n3} & \dots & a_{nn} \end{bmatrix}$ is a transition probability matrix where a_{ij} represents the probability of moving from state i to state j . In addition to that, the sum of outgoing probability of every state must be summed to 1. $\sum_{j=1}^n a_{ij} = 1$.

- q_0 and q_F represent initial and final states.

2.3.1.1 Hidden markov model

Markov chain is only usable if the event is directly observable in the world. However, we are interested in some events which are not observable or hidden for some problem. For example: acoustic event is observable, but word event is not observable in speech recognition problem. Hidden markov model can fix this problem.

Hidden markov model is formalized as following:

- $Q = q_1 q_2 \dots q_n$ is a set of states

- $A = \begin{bmatrix} a_{01} a_{02} a_{03} \dots a_{0n} \\ a_{11} a_{12} a_{13} \dots a_{1n} \\ \dots \\ a_{n1} a_{n2} a_{n3} \dots a_{nn} \end{bmatrix}$ is a transition probability matrix where a_{ij} represents the probability of moving from state i to state j such that $\sum_{j=1}^n a_{ij} = 1$ for every i .

- $O = o_1 o_2 \dots o_T$ is a sequence of T observations.
- $B = b_i(o_t)$ is observation likelihoods(or emission probability) where state i emits probability given observation o_t .
- q_0, q_F are initial and final states.

Hidden markov model employs two assumptions:

1. *Markov assumption.* The state is only dependent to the previous state.

$$P(q_i | q_1 q_2 \dots q_{i-1}) = P(q_i | q_{i-1}) \quad (2.8)$$

2. *Output independence assumption.* Output observation o_i is only dependent of the state q_i which produces o_t and independent of other states and other observations.

$$P(o_i | o_1 o_2 \dots o_i \dots o_T q_1 q_2 \dots q_i \dots q_N) = P(o_i | q_i) \quad (2.9)$$

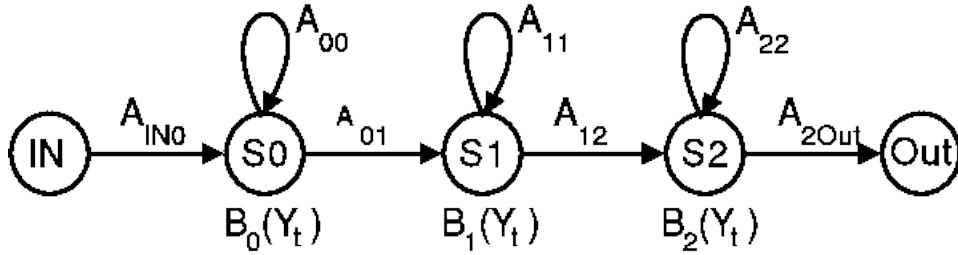
2.3.2 Using HMM to acoustic modeling

In section 2.3.1.1, HMM is explained. Now, it is time to apply HMM on acoustic model.

1. States

A state in acoustic model can represent a word, a phone, or a sub phone depending on the problem we are trying to solve. For the digit recognition problem, a state represents a sound of one number. A sequence of phones can represent as a word string. Hence, for more complicated tasks, like recognizing TV broadcast, a state represents a context-dependent sub phone(senone). A phone is usually influenced by the preceding and successive phone. Therefore, we use three HMM states(for each phone) which capture beginning, middle, and end of the phone.

Figure 2.1: A triphone HMM model consisting of three emitting states(beginning(S0), middle(S1), and end(S2) of a phone), a transition matrix A , and emission probability(observation likelihood) B [Silicon Intelligence, 2017]



2. Transition matrix

The transition matrix in speech recognition is the probability of moving from the current sub phone state to the successive sub phone state. Unlike HMM, Speech recognition HMM does not allow arbitrary transition because of the sequential nature of language and speech. A state can go to itself or to the next states, but it can not go back to the preceding state. This kind of left-to-right constrained HMM is called *Bakis network*.

3. Observation

The observation is a sequence of acoustic feature vector. The acoustic input is divided into small chunks and converted into a feature vector.

4. Observation likelihood

Observation likelihood is the probability of a feature vector which is generated by a sub phone state.

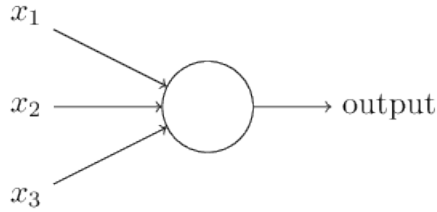
2.4 Deep learning for acoustic model

For our acoustic model, we used *Deep Neural Network Hidden Markov Model*(DNN-HMM). Beside DNN HMM, there exists *Gaussian Mixture Model Hidden Markov Model*(GMM-HMM). However, DNN HMM has been proven to outperform GMM DNN in practice [Peddinti et al., 2015]. We explain DNN HMM in this section

2.4.0.1 Feedforward neural network

Feedforward neural network(or sometimes called multilayer perceptron(MLP)) is one kind of neural network. MLP approximates any function $f'(\mathbf{x}) = \mathbf{y}$ which receives input x and outputs y . MLP can be represented as $f(\mathbf{x}; \theta) = \mathbf{y}$ and learn parameters θ which maximizes function approximations. Figure 2.2 represents a perceptron which receives three inputs and produces one output.

Figure 2.2: A perceptron [Nielsen, 2015]



Perceptron is a kind of artificial neuron developed in 1960 by Frank Rosenblatt [Rosenblatt, 1960]. A perceptron receives several inputs(x_1, x_2, \dots, x_n), does weighted sum $\sum_i x_i w_i$, and is thresholded by bias b to make a decision: zero or one.

$$output = \begin{cases} 0, & \text{if } \sum_i x_i w_i + b \leq 0 \\ 1, & \text{if } \sum_i x_i w_i + b > 0 \end{cases} \quad (2.10)$$

A perceptron is too simple to capture non linear pattern. Therefore, the perceptron is modified by applying an activation function into weighted sum. There are several common activation functions, such as:

1. Sigmoid

$$a(x) = \frac{1}{1 + e^{-x}} \quad (2.11)$$

2. Rectified linear unit(relu)

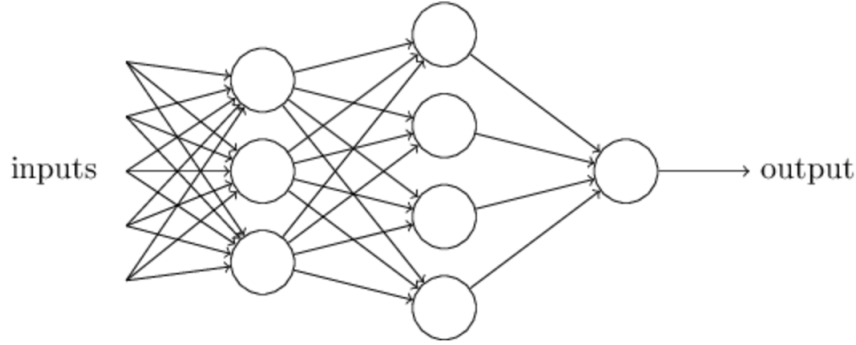
$$a(x) = \max(0, x) \quad (2.12)$$

3. Softmax

$$a(x) = \frac{e^x}{\sum_i e^{x_i}} \quad (2.13)$$

In addition to activation function, we need to make neuron to capture more complex features by stacking many neurons together. This network is called a *feed forward multi layer perceptron*. This neural network does not have connections between units which form a cycle. It is also called feedforward because the input values are propagated forward to output. The neural network receives inputs, does weighted sum, and applies activation function. The output of the units in one layer is propagated to the units in the next layer until reaching output layer.

Figure 2.3: Multi layer perceptron [Nielsen, 2015]



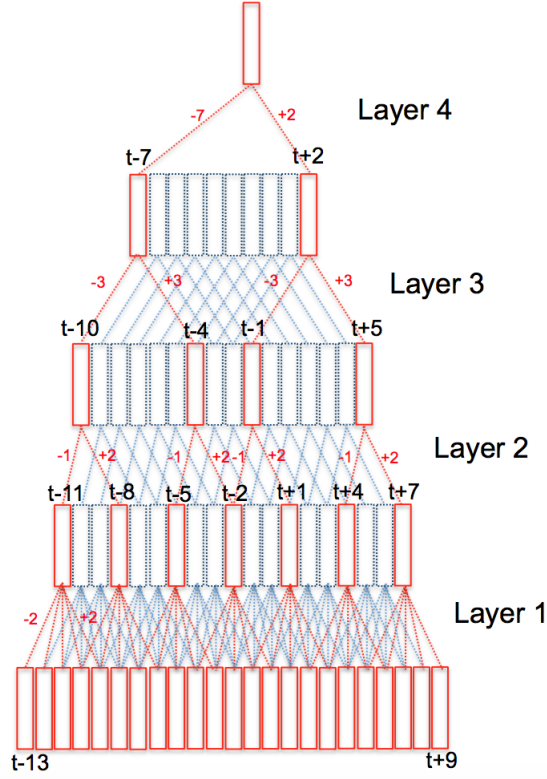
2.4.1 Time delay neural network(TDNN)

To model a temporal sequence neural network, like the speech recognition problem, the model must have the following properties: it must have multiple layers to learn complex non linear decision surface, represents relationships between events in time, invariant under translation in time, and does not require precise temporal alignment. From all of these properties, *Time Delay Neural Network*(TDNN) qualifies to model a temporal sequence neural network.

There exist two ways to capture long temporal contexts:

1. Feature representation which present the information of long temporal contexts. The examples of the feature representation are temporal patterns(TRAPs) [Hermansky and Sharma, 1999], multiscale spectro-temporal modulations representation [Mesgarani et al., 2004], deep scattering spectrum representation [Anden and Mallat, 2014], and modulation frequency feature representations [Thomas et al., 2009]. These feature representations can use the standard feed forward deep neural network to model the temporal sequence neural network.

Figure 2.4: An example of a TDNN architecture. [Peddinti et al., 2015]



2. Acoustic model which can learn long temporal contexts and is trained on short-term feature representations. Recurrent neural network(RNN) [Sak et al., 2014] and time delay neural network [Waibel et al., 1990] are two state-of-the-art acoustic models which are suitable to model this acoustic model.

Peddinti et al showed that TDNN performs better than RNN in term of parallelization during training [Peddinti et al., 2015]. Due to recurrent connections in RNN, RNN is not able to fully utilize parallelization in GPU(general processing unit) the same extent as TDNN. In addition, TDNN has been proven to have lower word error rate compared to RNN [Peddinti et al., 2015]. Thus, we will focus more on TDNN and trained acoustic models with TDNN in this internship.

Time delay neural network(TDNN) was first proposed by Waibel et al in 1990 [Waibel et al., 1990].TDNN is suitable to model long term temporal relationships within the range of N delays in short term speech features; for example: MFCC feature representations.

TDNN is basically a feedforward neural network. Learning procedure which we used is backpropagation. There are two passes procedure in the backpropagation. The first pass is forward pass where the neural network receives input, does weighted

sum, and applies activation function. The output values from the first layer is forward propagated until reaching the output layer. Then, square error is calculated between the expected value and the output of neural network. The second pass is backward pass where the derived square error value is propagated backward to update weights.

Typical TDNN computes all hidden activations in each layer. Figure 2.4 demonstrates an architecture of TDNN. TDNN receives acoustic inputs with some delays to the left and right contexts(before and after time t). Usually, the input context is asymmetric, i.e: more time context to the left than right. The main reasons are to reduce the latency of online decoding and improve word error rate. Each layer splice together contiguous temporal windows within a time range; for example: the third layer splices the frames within range $[t - 3, t + 3]$ in the figure 2.4. Furthermore, TDNN applies wider splice when going to higher layers in the networks as shown in the figure 2.4. The reason is we force the higher layer to learn wider temporal contexts.

Training a TDNN is still time consuming. The standard TDNN computes all hidden activations in each layer. Consequently, the training time of TDNN is approximately ten times compared to the standard DNN. Peddinti et al found that there exist large overlaps between input contexts at neighboring time steps; thus, they proposed subsampling [Peddinti et al., 2015]. The technique generally only splice no more than two temporal contexts with gap instead of contiguous temporal windows. Figure 2.4 illustrates the procedure of subsampling. In the standard TDNN, all hidden units(the red and blue squares represents hidden units) are calculated. In contrast with subsampling technique, only red hidden units are computed. In consequence, the training time of subsampled TDNN is approximately five times faster than the standard TDNN [Peddinti et al., 2015]. Moreover, subsampling is able to reduce the model size. Splicing contiguous temporal contexts will require a massive number of parameters. However, subsampling only splices less temporal contexts than the typical TDNN.

2.4.2 Hybrid acoustic model

An acoustic model approximates $P(o|y)$ which is the probability of observation o (short span of acoustic features) given HMM state label y . The HMM state label y is *senone*, a context-dependent sub phonetic state. The standard HMM approach to speech recognition system is gaussian mixture model(GMM) HMM. GMM-HMM is a HMM where $P(o|y)$ represents the observation likelihood of HMM state y and $P(o|y)$ is approximated by GMM acoustic model. A hybrid HMM resembles GMM-HMM, but

the hybrid HMM leverages a neural network to compute $P(o|y)$ instead of GMM [Maas et al., 2014].

As explained in 2.1, an acoustic model computes $P(o|y)$. However, a neural network can not directly compute $P(o|y)$. A neural network can be seen as a classifier of senone given acoustic input $P(y|o)$. By applying the bayes rule, $P(o|y)$ can be obtained as the following:

$$P(o|y) = \frac{P(y|o) \times P(o)}{P(y)} \quad (2.14)$$

$P(y)$ can be obtained from a distribution of senone in the training set. $P(o)$ is a probability of a particular acoustic input, which is difficult to model. Luckily, the observation o are fixed and can be seen as a constant. Hence, we can provide HMM observation likelihood with the unnormalized probability.

$$\frac{P(y|o)}{P(y)} \quad (2.15)$$

Our training data comprises of word level transcriptions. Each utterance(segment) only contains start and end times without time alignment for each word. To train a neural network acoustic model, each acoustic input frame must be labeled with its corresponding senone in each utterance. We can use a force alignment to produce a sequence of senone labels for each utterance. GMM acoustic model can generate the force alignment of the training data [Maas et al., 2014]. The aligned data(where each acoustic input frame is labeled with its corresponding senone) can be used to train a neural network acoustic model.

Applying the hybrid approach has been introduced over 20 years ago [Bourlard and Morgan, 1993, Renals et al., 1994]. Not until recently, the neural network HMM were not the standard approach. The standard approach for building speech recognition systems was GMM-HMM. However, the trained system(GMM-HMM) became complex when they tried to increase the representational capacity of GMM [Maas et al., 2014].

The modern DNN acoustic model is able to increase representational capacity better than GMM [Maas et al., 2014]. DNN acoustic model reduced substantial amount of word error rate compare to GMM [Dahl et al., 2011, Jaitly et al., 2012]. The HMM speech recognition system using DNN can be called DNN-HMM. There are several factors which contribute to the success of DNN-HMM just recently(compared to 20 years ago).

- The amount of training dataset has increased. DNN usually works well when there exists a massive volume of training data.
- Nowadays, training DNN acoustic models can be done on GPU which accelerates the execution time. A DNN has many hidden layers; hence, the training time is time consuming. The parallelism in GPU allows a DNN acoustic model to be trained more rapidly compared to the parallelism in CPU.
- We have more powerful neural network architectures recently with more number of parameters, more hidden layers(deeper network), and better initialization procedure.

Chapter 3

Methodology

3.1 Lightly supervised data selection approach

To generate a "good" speech recognition, one needs a massive number of training audio and their exact transcriptions. However, transcribing audio is labor intensive and time consuming too. There are unlimited supply of audio data in the internet, television, radio, and other sources. Nonetheless, they usually have a few or even no accurate transcription at all. Some television broadcasts have corresponding closed captions. Closed caption means that the transcription is close, but not exactly the same as what the audio says. Even, the closed captions are often badly time-aligned with the audio.

By utilizing these audio data with the corresponding closed captions, we hope to produce a high performance speech recognizer with less supervision. The main idea is to use the automatic speech recognizer to transcribe audio data which will automatically generate transcriptions, which later be used as training data [Lamel et al., 2002].

The idea of using untranscribed audio data has been explored before (Zavaliagos and Colthurst in 1998 [Zavaliagos and Colthurst, 1998] and Kemp & Waibel [Kemp and Waibel, 1999]), which is called as unsupervised acoustic model training. They utilized small amount of audio data without transcription to train acoustic models. The weakness of their experiments is not using massive volume of training data which is essential to train an acoustic model [Lamel et al., 2002].

Lightly supervised approach was proposed in 2002 by Lamel et al [Lamel et al., 2002]. In general, the lightly supervised approach operates as following:

1. Train an acoustic model on a small amount of manually annotated data.
2. Recognize(or automatically transcribe) a massive amount of data.

3. Align the automatic transcriptions with the closed captions. Some transcriptions and closed captions might disagree. We can remove or correct these segments.
4. Retrain a new acoustic model using the data which we selected in the previous step.
5. Reiterate from step 2.

These steps can be iterated several times as long as the error rate is decreasing. This method uses the idea of training acoustic models in less supervised manner because the training dataset(closed captions) is not the real detailed transcription.

Using closed caption as training data reduces effort in manual transcription. Detailed transcription process usually takes 20-40 times manual effort compare to live closed captioning. Closed caption is usually produced in real time and less costly compare to the detailed transcription. In addition, the manual transcription is possible to have error [Barras et al., 2001]. Additionally, some TV broadcasts, such as: CNN headline news, ABC world news tonight, BBC, already have closed captions which can be used as a training dataset. Nevertheless, some problems exist when using the data with closed captions as training dataset. Training using the closed captions faces several disadvantages compared to the real transcriptions: indication of non speech events, such as: coughing, speaker turn, acoustic conditions: background noise and music. Furthermore, some sentences might be paraphrased, deleted, and changed in word order.

Large text data are available in internet, such as: news, blog, and closed captions from TV broadcast. However, these data might be less related with the audio data. Some additional sources may be not contemporaneous with the data. Thus, it provides less supervision. We can interpolate these data with the closed captions to build a more general language model.

Despite the fact that closed captions are not accurate, it has been proven that closed captions provides enough supervision in building automatic speech recognition systems [Lamel et al., 2002]. Lamel found that an increase in training data improves accuracy even though it will increase the model size. Additionally, filter closed captions for the next iteration training set in the lightly supervised technique improves the accuracy slightly. The language model built from closed captions also effectively provides enough supervision in building the ASR. Therefore, lightly supervised technique is our main inspiration to do data selection in this internship.

3.2 State-of-the-art in data selection

Some researches have been conducted to investigate the lightly supervised approach to select data in the speech recognition problem [Lamel et al., 2002, Lecouteux, 2006, Chan and Woodland, , Mathias et al., , Stolcke et al., 2000].

Lamel et al proposed the lightly supervised approach which harnessed the closed caption as supervision to build a robust speech recognizer [Lamel et al., 2002]. In their experiment, they used a bootstrap acoustic model and biased language model to recognize the audio where the decoding result is compared with the closed caption. Then, the words which do not agree will be removed from training dataset. The new training dataset is utilized to train a new acoustic model. This filtering (removing not-agree words) is called *closed caption filtering*.

Chan and Woodland explored another way in filtering bad closed captions [Chan and Woodland,]. They leveraged the confidence measure metric to remove the bad segments. When decoding acoustic inputs, an ASR produces word hypothesis and their corresponding confidence measure. The confidence measure value from each word in each segment is averaged to produce one confidence value per segment. Finally, they determined a threshold confident value to filter out all potentially bad segments, where the confident value is lower than the threshold. Moreover, they also proposed to interpolate two different LMs created from two closed caption datasets and merge them into one single LM. This single LM is used as light supervision when decoding acoustic input together with the acoustic model and the lexicon.

Matias et al applied lightly supervised approach on medical conversation data [Mathias et al.,]. There are unlimited supplies of medical speech; however, only informal medical reports are available. The medical reports are not the same as what the physicians said in the audio records. This problem matches perfectly with the lightly supervised approach. They explored frame level filtering in place of word level or segment level filtering. Basically, they forced align and compared the decoding results with the closed captions and mark words which are insertion, deletion, and substitution. All corresponding frames which are associated with those marked words are removed or filtered out from training set to produce a new cleaner training set.

The interesting data selection was proposed by Lanchantin et al from Cambridge university [Lanchantin et al., 2016]. They used the lightly supervised approach applied on the Multi Genre Broadcast (MGB) challenge. MGB challenge is a challenge to automatically transcribe TV broadcasts. The challenge organizers only provided

TV broadcast audio data and their corresponding closed captions. General TV broadcast data are usually recorded in highly diverse environment speech with background music, non speech event, and sound effect. In addition, some closed captions may be different compared to the real transcription due to deletion, insertion, substitution, and paraphrasing. These factors make the challenge interesting as well as complicated. Lanchantin et al proposed to tackle the problem inspired from the lightly supervised approach.

3.3 Data

Before talking about the methodology, it is better to talk about the dataset. We have the following dataset:

1. Training set has audio data with their corresponding closed captions. We call this set to be train.full which contains 200 hours of training data. In addition to the full training set, we randomly select a small subset of data which is called train.small. This subset of data is a random selection of 100 hour data from train.full(We call the 100 hour data as train.100). This training set is used for training and building speech recognition systems.
2. Test set has audio data and their corresponding manual transcriptions. These manual transcriptions are the same as what people said in the audio. Because we have the exact transcriptions, we can use this test set to evaluate our data selection method.
3. Additional large text corpus which is utilized for building language model. This additional corpus is different than the closed captions in training set.

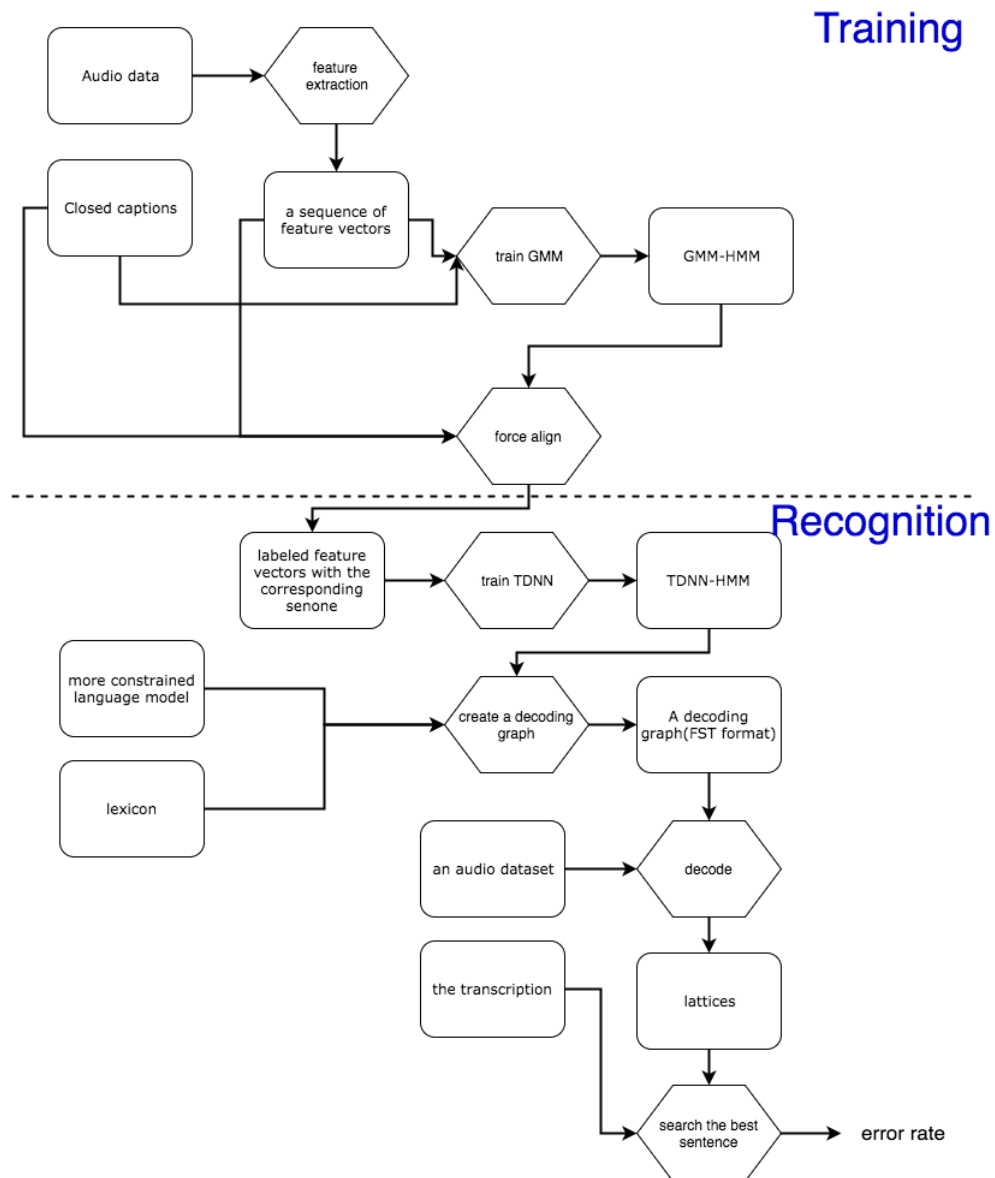
3.4 Baseline data selection method

Inspired by the lightly supervised approach, we developed our data selection method. Firstly, we dive deeply into how to build the acoustic model as well as the language model and combine them into a speech recognition system for data selection. Then, this section continues with the general explanation of the data selection. This baseline method is a benchmark for our proposed model to determine whether the proposed model thrives against the baseline.

3.4.1 Building the speech recognition systems

The previous section explains the data selection process. Nevertheless, it does not elaborate the way to train the acoustic model and to decode. This section elaborates how the speech recognizer works together as illustrated in the figure 3.1.

Figure 3.1: Training and decoding with the speech recognition system



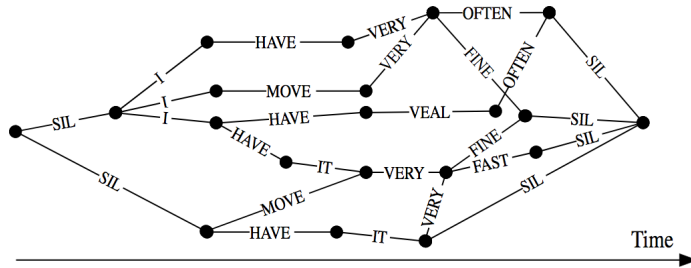
Before training an acoustic model, it is necessary to turn the audio data into a sequence of acoustic feature. The audio will be divided into small chunks and each chunk will be converted into a feature vector.

The hybrid HMM acoustic model is the approach to build the acoustic model as explained in the section 2.4.2.

Firstly, train a GMM-HMM, force align, and label each short-span acoustic inputs. After that, TDNN acoustic model is trained. The TDNN acoustic model, the less constrained language model, and the lexicon are combined together to make a decoding graph. The decoding graph is based on *Weighted Finite State Transducer*(WFST). WFST is a finite automata whose state transition is labeled with input, output, and a weight which encodes a probability to go to the next state. The transition maps the input symbol sequence to an output string. WFST is able to encode HMM, lexicon, and n-gram language model [Mohri et al., 2008]. Thus, WFST is a perfect representation which maps the acoustic input to a word string.

When decoding an audio dataset, the decoder does not return the most probable sentence given the acoustic inputs. Instead, it returns a directed graph containing the N-best set of hypotheses. This directed graph is called *lattice*. The nodes in *lattice* represent points in time, while the arcs represent a lot of information, such as: the start and end times, the acoustic model probabilities, the language model probabilities, the sequence of phones hypothesis, and phone duration.

Figure 3.2: An illustration of lattice [Gales and Young, 2007]



The *lattice* can be used later for decoding and rescoring with different language model or acoustic model [Murveit et al., 1993]. From the generated lattice, we can search the best sequence with different and more sophisticated configuration (but slower to do) or more complex language model. Searching the best sequence can be done simply by tracing the path in the lattice which results in the biggest probability.

3.4.2 Metric

Before talking about the data selection method, it is better to elaborate metrics which we used for selecting the data. Here are the metrics:

1. Average word duration(AWD) and average phone duration(APD). This metric is to detect if there exist errors in aligning the start and end time of a segment or recognizing audio data. A duration of a word can not exceed an upper limit threshold and the duration can not be lower than a bottom limit threshold.

$$AWD = \frac{\text{utterance duration in second}}{\text{number of words in the utterance}} \quad (3.1)$$

$$APD = \frac{\text{utterance duration in second}}{\text{number of phones in the utterance}} \quad (3.2)$$

2. Phone error rate and word error rate. This metric is used to measure the performance of a speech recognition system.

$$\text{error_rate} = 100 \times \frac{\text{substitution} + \text{deletion} + \text{insertion}}{\text{substitution} + \text{deletion} + \text{correct words}} \quad (3.3)$$

We usually obtained error rate by comparing between the real transcriptions and the decoding transcriptions produced by the speech recognition systems. Phone error rate is computed by the comparison in phone level; in contrast, word error rate is obtained by the comparison in the word level.

Our training set only has closed captions(not the real transcription). To avoid the confusion when comparing the closed captions and the decoding result, we name them as phone matched error rate(PMER) and word matched error rate(WMER).

Here is the example how to compute AWD, PMER, and WMER. We compare the closed caption and its corresponding decoding transcription of a segment.

- ref: there aren't that many parts in THE story
- hyp: there aren't that many parts in *** story

The ref line is the closed caption, while the hyp line is the decoding transcription. The start and end time is 521.2 and 523.92 second respectively. Because there are 7 words in the decoding result, the AWD is $\frac{523.92-521.2}{7} = 0.39$. There is one deletion, zero substitution, and zero insertion; thus, the WMER is $100 \times \frac{0+1+0}{0+1+7} = 12.5$.

3.4.3 Data selection

The data selection will be divided into two process pipelines which are intertwined together. The first pipeline is to select data and the second one is to evaluate how good the data selection is. The evaluation pipeline is also important to know when to stop the iteration of the data selection pipeline.

Figure 3.3: Data selection pipeline

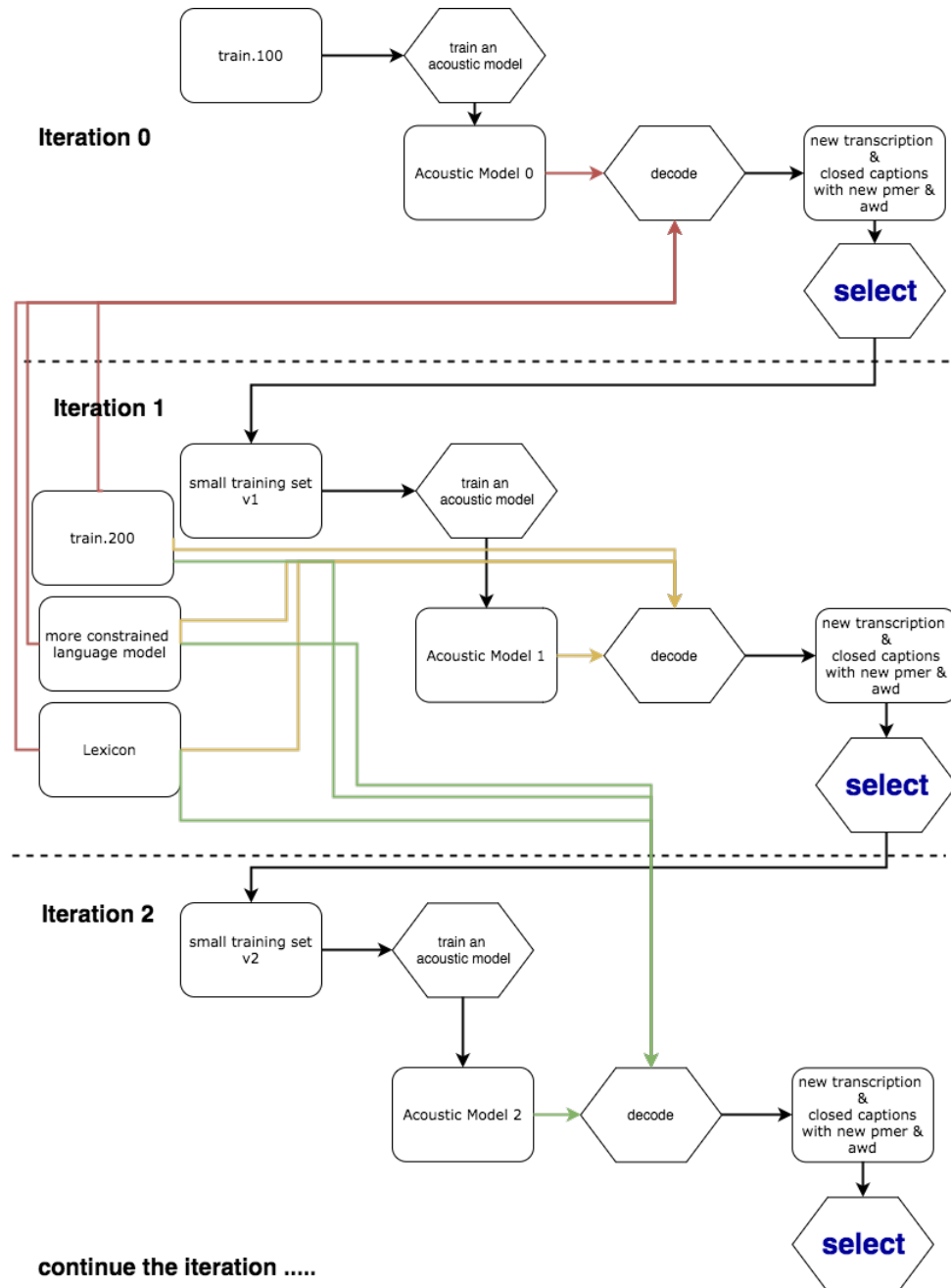


Figure 3.3 demonstrates the overall data selection pipeline. Here is the explanation how it works:

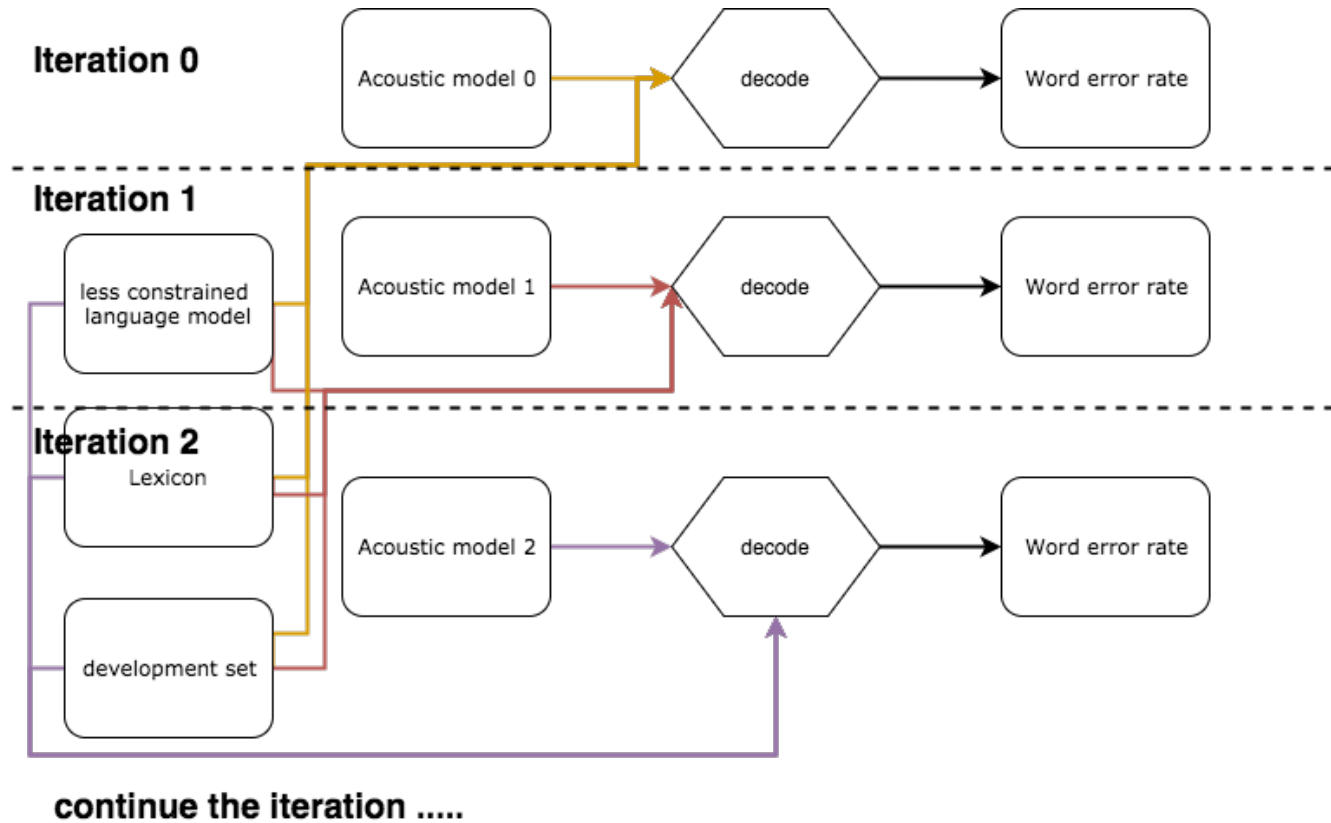
1. Build two language models: a more constrained language model and a less constrained language model. The more constrained is built from the closed captions of the full training set. In contrast, the less constrained language model is built from interpolation of the closed caption language model and the additional large text corpus language model. The more and less constrained language models are utilized in the data selection and evaluation pipeline respectively.
2. Randomly select a subset of the training set. The small training set is used to train an initial acoustic model(acoustic model 0).
3. Train an initial acoustic model(acoustic model 0) by utilizing the audio and the closed captions in the small training set.
4. Decode the full training set with the acoustic model 0, the more constrained language model, and the lexicon. This will produce the new decoding results and new values of PMER and AWD. The new values of PMER and WMER are obtained from comparing the closed captions and the decoding results.
5. Select the closed captions from the training set based on PMER and AWD. This will be the new training set for training the next acoustic model. The selection works as following:
 - (a) Select the segments which have AWD within the range $0.166 \leq AWD \leq 0.65$. This range follows the range which was used in [Lanchantin et al., 2016]. AWD represents the length of word duration. If the AWD is high, it means each word may have long duration. If AWD is too tiny or large, it means that the corresponding segment is incorrect or wrongly align. Hence, it must be removed.
 - (b) Sort the segments based on the new value of PMER ascendingly.
 - (c) Choose the top segments to make a new small training set. Top segments are segments which have the lowest PMER.
6. Continue the step 3-5 until the data selection does not improve anymore.

The last step from the data selection pipeline says to do iteration over and over again until the technique does not improve. The technique improves when selecting

the closed captions which are closed to the real transcription. However, until some points, it may select the data which can not improve the acoustic model or even the same data in the next iteration. Therefore, the evaluation pipeline is necessary to evaluate how good the data selection is and stop the iteration. Figure 3.4 illustrates the evaluation pipeline. Here is the explanation how it works:

1. In each iteration, a new acoustic model is trained.
2. By utilizing the acoustic model, the less constrained LM, and the lexicon, the development set is decoded. By comparing the decoding result and the real transcription in the development set, one obtains the WER. If the WER decreases in the next iteration compared to the previous iteration, the data selection improves; otherwise, it stops improving and we can stop the iteration.

Figure 3.4: Evaluation pipeline



3.5 New proposed data selection methods

The general idea of the proposed methods is to combine three different automatic speech recognition(ASR) systems by varying the language model but with the same

acoustic model.

The language models are built with different constraints(the least, medium, and the most constraint language models). The proposed methods work almost the same as the baseline method: training acoustic model with the subset of training data(audio data and their closed captions), recognizing and selecting from the full training data, and repeating the steps to do more data selection. However, the difference is when recognizing and selecting data. Three different ASRs recognize the full training set, combine, and select the decoding results as following:

1. Average PMER and AWD of three ASRs.

We built three speech recognition systems and each recognition system recognized the same training set; consequently, we had three transcriptions which have the same amount of segments. After that, average these transcriptions into one transcriptions with the following algorithm.

Listing 3.1 the proposed averaging algorithm

```

1: For each segment:
2:     Average the value of PMER and AWD of three corresponding
3:     segments(from three decoding transcriptions).
4:     Choose the closed captions using the averaged PMER and AWD
5:     as the new training set.
```

After averaging the value and producing a new training set with new value of PMER and AWD, we can select the data with the same method as the baseline method.

2. Combine PMER and AWD of three ASRs with a proposed combination algorithm

Here is the pseudocode how the algorithm works.

Listing 3.2 the proposed combination algorithm

```

1: For each segment:
2:     Keep only segments which satisfy  $0.166 < AWD < 0.65$ 
3:     and  $0.03 < APD < 0.25$  range
4:     If one ASR segment has zero PMER,
5:         Select the segment and corresponding closed
6:         captions
7:     Else
8:         If two segments have the same phone sequence
9:             Select the segment and the corresponding
10:            decoding result transcription
11:         Else
```

12: Sort by PMER. Choose top segments with
13: the lowest PMER. The closed captions will
14: be chosen.

3. Utilize a randomly different training set for each iteration.

In our baseline model, we always recognize the same full training set and select the best subset from the training set. Instead of using the same training set, we randomly choose a new 200 hours training set. The intuition behind this is randomness may increase the performance of the model.

Chapter 4

Implementation

4.1 Data set

For this internship, we used TV broadcasts downloaded from the internet.

- The audio files are taken from TV broadcasts. Total duration of audio is around 200 hours.
- TV broadcast transcriptions(*closed captions*) of the audio files are provided. The *closed captions* are close to what the audio says, but not exactly the same as the speakers of the audio said. We call this training set as train.200, containing audio data and the corresponding *closed captions*.
- 640 million words of TV subtitles are provided. In addition to that, the 640M subtitles are filtered to avoid overlap with the *closed captions* from the TV broadcasts.
- The development set contains around 8 hours of audio data and their corresponding transcriptions. The transcriptions in the development set are closed to what people said in the audio.

We prepared two data sets: training set and development set. The training dataset has a total duration of around 200 hours of audio. To evaluate our model, we have a development set. Some metadata are provided in our data, such as: speaker identity of the transcript, genre of the show, date and time of the show, as well as television channel where the show was aired. In addition, each segment has start and end time when the segment was shown in the TV broadcast. All files in full training set have *closed captions*. In contrast, all files in development set have *closed captions* as well as manually annotated human transcription. The transcript is manually and carefully

annotated by human. Thus, the transcript is believed to be the closest transcription spoken by speakers in audio files.

Table 4.1 shows the statistics of the training set and the development set. The second column represents the number of TV shows in each genre; the third column represents the total duration of each genre. Lastly, the forth column(the last column) shows the total duration of speech when the speakers spoke in the audio.

Table 4.1: Dataset duration

No	Data set	#shows	Duration(h)	Aligned speech (h)
1	training set	274	193	149
2	development set	12	8	6

After obtaining about 200 hours of data, a 100 hours subset was created randomly from 200 training set. The 200 hours and 100 hours training set are named as train.200 and train.100 respectively.

Table 4.2: Statistics of train.200 per genre

Genre	#shows	Duration(h)	Aligned speech(h)
advice	40	26	22
children	51	19	14
comedy	19	8	6
competition	30	21	16
documentary	29	22	14
drama	20	14	10
events	24	38	29
news	61	42	37

Table 4.2 shows that the statistics of our training data. Our data has 8 genres. The second column(from left) shows the number of shows corresponding to their genre; the third columns shows the total audio duration of each genre; the forth column shows the total speech duration. The data with the news genre and events genre have the most duration and the second most duration respectively. In contrast, the comedy genre data have the least duration.

4.2 Experimental setup

In this section, we explained how we did the experiments. We elaborated how to build language models, the acoustic model, and how to decode the data set for the data selection and the model evaluation.

4.2.1 Language model

4.2.1.1 Library: SRILM

SRILM is a toolkit written in C++, for developing statistical language models for speech recognition and other natural language applications [Stolcke, 2002]. SRILM allows creating a language model from a text corpus, interpolating several language models, pruning a language model, and estimating the perplexity of a language model using test corpus.

The code bellow is an example how to make a 4-gram language model with input from *CORPUS* and interpolated with Kneser-ney interpolation [Kneser and Ney, 1993]. The syntax is pretty straightforward and easy to understand. The resulted language model will be written in arpa language model format.

```
ngram-count -order 4 -kndiscount -interpolate -sort -text CORPUS -lm LM
```

4.2.1.2 Implementation

This internship experimented with three different kinds of language model for speech recognition. All language models are based on four gram language model.

1. LM.200

This language model, named as LM.200, was produced from the *closed captions* of 200 hours training set(train.200). LM.200 was used for the baseline speech recognizer to select the "good" subset of data. Table 4.3 displays the statistics of the language model. The closed caption dataset has 33,914 words. Section 3.4.3 mentions the more constrained language model for the data selection pipeline. In the internship, we utilized LM.200 as the more constrained language model.

Table 4.3: Statistics of LM.200

No.	File	Information
1	word.200	33,914 words
2	LM.200	uni-gram=33916 bi-gram=402299 tri-gram=111868 4-gram=64801

2. LM.7weeks+subtitles.limited.1e-9

A language model was generated from the *closed captions* of 1600 hours transcription(7weeks TV broadcast *closed captions*) and interpolated with a language model from the big TV subtitles with ratio 0.9/0.1. The ratio means

that estimation of language model of *closed captions* is multiplied by 0.9, the estimation of big subtitle language model is multiplied by 0.1, and and they are summed together. The language model was limited by top 160,000 frequent word list and pruned by 10^{-9} resulting in LM.7weeks+subtitles.limited.1e-9. Section 3.4.3 mentions the less constrained language model for the evaluation pipeline. In the internship, we utilized LM.7weeks+subtitles.limited.1e-9 as the less constrained language model.

No.	File	Information
1	word.160k	160,000 words
2	LM.7weeks	ngram 1=91836 ngram 2=1817881 ngram 3=980925 ngram 4=847736
3	LM.subtitles	ngram 1=756644 ngram 2=27144330 ngram 3=37162518 ngram 4=57912280
4	LM.7weeks+subtitles	ngram 1=768522 ngram 2=27441072 ngram 3=37312673 ngram 4=58183256
5	LM.7weeks+subtitles.limited	ngram 1=160002 ngram 2=24926818 ngram 3=32102763 ngram 4=44253079
5	LM.7weeks+subtitles.limited.1e-9	ngram 1= 160002 ngram 2= 5251197 ngram 3= 3876171 ngram 4= 2453067

3. LM.genres

Eight language models were generated from each genre(LM.documentary, LM.news, LM.events, LM.drama, LM.competition, LM.comedy, LM.children, and LM.advice). Then, each language model was interpolated with the big subtitle LM with ratio 0.9/0.1, limited by the top 160,000 word list, and pruned by 10^{-9} threshold. The language model was used in the new proposed data selection only. To build language models, we utilized SRILM.

No.	File	Information
1	LM.documentary	ngram 1=37542 ngram 2=437631 ngram 3=135632 ngram 4=90381
2	LM.news	ngram 1=44588 ngram 2=661625 ngram 3=305473 ngram 4=255125
3	LM.events	ngram 1=23717 ngram 2=306938 ngram 3=125934 ngram 4=95203
4	LM.drama	ngram 1=21345 ngram 2=202078 ngram 3=60920 ngram 4=40073
5	LM.competition	ngram 1=33269 ngram 2=385257 ngram 3=124584 ngram 4=89642
6	LM.comedy	ngram 1=20180 ngram 2=165073 ngram 3=39353 ngram 4=23877
7	LM.children	ngram 1=27029 ngram 2=287141 ngram 3=84063 ngram 4=57851
8	LM.advice	ngram 1=30545 ngram 2=397766 ngram 3=154416 ngram 4=108794

4.2.2 Acoustic model

4.2.2.1 Library: Kaldi

Kaldi is a toolkit, written in C++, which comprises several tools to experiment and build a speech recognition system [Povey et al., 2011]. Kaldi supports several operations to build a speech recognition system, such as: extracting feature vectors, acoustic modeling, training acoustic models, and creating decoding graph as well as decoding audio inputs into lattices. Kaldi provides several features, such as:

- It supports finite state transducer(FST), which is based on OpenFST library. The FST is used to build a decoding graph with the inputs of acoustic model, language model, and lexicon.
- It is extensible and provides complete recipes. Kaldi gives many recipes to build acoustic model, such as: gaussian mixture model, deep neural network, etc.

4.2.2.2 Acoustic model implementation

We utilized TED-LIUM’s recipe(nnet3 in Kaldi) to train the TDNN acoustic model. The feature vectors are MFCC, consisting of 40 feature values. The TDNN architecture has 6 hidden layers, where each hidden unit utilizes rectified linear unit(RELU) activation function. The TDNN outputs around 9000 output nodes with softmax

activation function. Moreover, the TDNN architecture is implemented by leveraging the subsampling technique. Table shows the parameters and subsampling splice configuration we used in the architecture.

Table 4.4: Parameters of theTDNN acoustic model

Parameter	Method
initial learning rate	0.0015
final learning rate	0.00015
mini batch size	512
epoch	3
subsampling splice	
layer 1	$t - 2, t - 1, t, t + 1, t + 2$
layer 2	$t - 1, t + 2$
layer 3	$t - 3, t + 3$
layer 4	$t - 7, t + 2$
layer 5	$t - 3, t + 3$
layer 6	t

4.2.3 Recognition

After training the acoustic models and building the language models, we recognize a data set. The recognition is useful for data selection(computing the new word error rate as well as phone error rate) and evaluate how good our data selection is. In overall, the recognition(or decoding) process works as following:

1. Create a decoding graph from the acoustic model, the language model, and the lexicon. In Kaldi’s implementation, it transforms the acoustic model, the language model, and the lexicon into a weighted finite state transducer(WFST).
2. Decode the dataset using the WFST. In Kaldi, the decoding process will result in a lattice. Lattice is a decoding graph which represents some variants of the recognition. Providing only one best recognition is not practical; thus lattice is the way to represent the recognition/decoding result.
3. Compute the best path of words or phones in the lattice and calculate the WER as well as PER. Kaldi produces two ctm files(word and phone ctm files) which are the best decoding results which the algorithm can obtain. The ctm file consists of words(or phone) in the transcription with its start and end times relative to the audio.

Sc-lite is a tool for evaluating and scoring the output of speech recognition systems. By utilizing sc-lite, the ctm file(the hypothesized word/phone from the

ASR) is compared with the transcriptions(the *closed captions* or the detailed transcription) and the error rate is calculated.

4.2.3.1 Data Selection

To select the good data for the next iteration, we need to re-recognize the training set and calculate the new value of matched word phone matched error rate, and average word duration. Here are the steps for computing the value of PMER and WMER:

1. Recognize the training set with the trained acoustic model, the language model, and the lexicon. This will result new decoding transcriptions.
2. By utilizing scilite, we can compare the *closed captions* and decoding transcriptions. After comparing, scilite automatically generates the numbers of correct words, insertions, deletions, and substitutions which is written in a file with extension ".prf". Figure 4.1 illustrates a segment with its corresponding *closed caption* and decoding transcription. REF represents the closed caption, while HYP represents the decoding transcription. After knowing the number of correct words, insertions, deletions, and substitutions, the values of PMER, WMER, and AWD can be computed.

Figure 4.1: The example

```

Speaker sentences  0:  spk-0001  utt# 1 of 99225
id: (spk-0001-001)
Labels: <o,f0,male>
File: 20080401_004000_bbcone_miracle_on_the_estate
Channel: a
Scores: (#C #S #D #I) 21 0 0 0
Ref times: t1= 20.13 t2= 26.85
REF: in      england about   five   hundred years ago   whole towns and villages got together once a year to
HYP: in      england about   five   hundred years ago   whole towns and villages got together once a year to

```

3. After computing the new value of PMER, WMER, and AWD for each segment, we select segments with AWD in the range of $0.165 \leq AWD \leq 0.66$ and then sort the selected data based on PMER. We only choose top 100 hours data which has the lowest PMER.

4.2.3.2 Evaluation

To evaluate how good our data selection is, we need to evaluate the automatic speech recognizer. In every data selection iteration, we train an acoustic model, build a speech recognizer (with the acoustic model, the language model(LM.7weeks+subtitles.limited.1e-9), as well as the lexicon), and recognize the development set. After decoding and

computing WER of development set, we compare the WER with the previous iteration. If the new WER is lower, we can continue the iteration; otherwise, stop the iteration.

4.3 Experiment resources

To run our experiment, we utilized grid 5000 server cluster. Grid 5000 is a computing cluster with a focus on parallel and distributed computing [Team, 2017]. It has several key features:

- It provides a large amount of resources, over 1000 nodes, and massive volume of hard disk.
- The experiment is highly configurable where we can specify number of nodes and the criteria of the nodes which we will use.
- We can do advanced monitoring of memory usage, cpu usage, and the log of our experiment which later be used to analyze the experiment.

The following is the description of resources which we leveraged for this internship based on the tasks:

1. Train acoustic models: 1 computer with GPU and 8 computers without GPU. At least one computer with GPU is compulsory because we can harness GPU to parallelization the computation of DNN training.
2. Create decoding graph only needs one computer because the Kaldi recipe does not need parallelization for this process.
3. Decoding the ASR leverages 20 computers. Usually one computer can recognize one hour of audio within one hour. Therefore, the more computer we use, the faster the decoding process will be.
4. The error rate computation and data selection process only utilizes one computer.

Chapter 5

Experiment result

5.1 Baseline model

As presented in the previous chapters, the baseline data selection works as following:

1. train a GMM-HMM and then TDNN acoustic model(AM) on 100 hour randomly selected training set(with closed caption)
2. recognize the training set using this acoustic model
3. compute phone matched error rate and average word duration
4. reject segments(utterances) which does not satisfy the AWD range($0.165 \leq AWD \leq 0.66$)
5. sort and select new data based on the lowest phone matched error rate. Retrain an acoustic model using the selected data. Repeat steps 1-5.

To evaluate how good the data selection is, we recognize the development set and compute word error rate. If the word error rate decreases, the iteration can continue; otherwise, stop the iteration.

Figure 5.1 shows the average word duration of data selection iteration 0, 1, and 2. The X-axis is the value of average word duration, while the Y-axis represents the cumulative duration of audio in percentage. All utterances are sorted based on AWD to draw this figure. Two vertical line represents boundaries to select utterances. The utterances which are inside these two lines(the value of AWD is in the range $0.165 \leq AWD \leq 0.66$) will be selected. The result shows that less than 15% data are rejected. Selection by using AWD is important to reject utterances which are badly aligned.

Figure 5.2 demonstrates the phone matched error rate of all iterations. The X-axis and Y-axis represent the value of PMER of each utterance and the duration of audio in percentage respectively. All utterances are sorted based on PMER to generate this figure. The value of PMER decreases slightly over several iterations where more and more utterances have less PMER in the next iteration. Nevertheless, the iteration stops in iteration three because the data selection converges(no further improvement in PMER).

Table 5.1: WER and PER of acoustic models

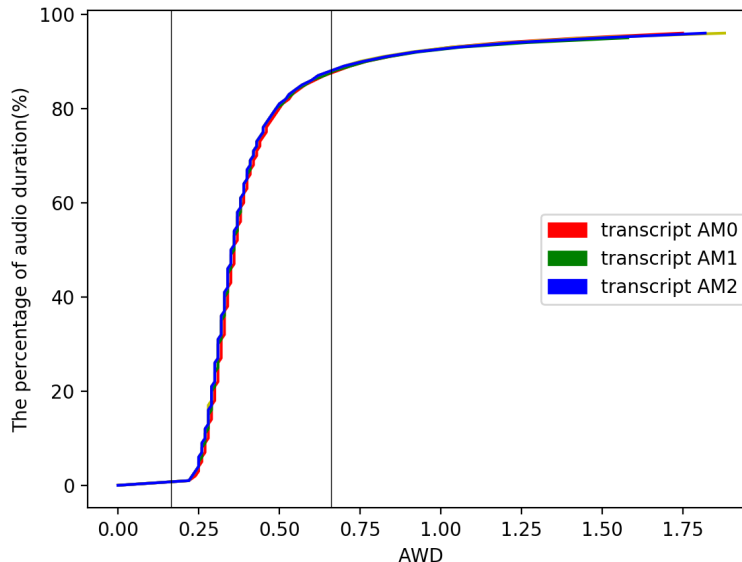
No.	Model	WER(%)	PER(%)	PMER thresh(%)
1	baseline 0	43.9	35.8	36.0
2	baseline decoding	42.3	34.3	45.98
3	baseline WMER	41.2	34.1	
4	baseline 1	41.2	33.8	35.71
5	baseline 2	40.5	33.3	35.44
6	baseline 3	40.5	33.3	

Table 5.1 shows the result of the recognition experiments. The second column denotes the model name; the third and forth represents the value of WER and PER respectively of the corresponding model. The last column shows the PMER threshold. When we finish building a speech recognizer, we recognize the train.200 training set and select 100 hours data for the next iteration. The selection is thresholded by PMER threshold, shown in the last column.

- Baseline 0(row 1) is the initial data selection iteration(iteration 0) where the training set is randomly 100 hours selection from 200 hours of training set. We evaluated baseline 0 which resulted in 43.9% WER and 35.8% PER. Then, we re-recognize train.200 and selected top 100 hours data with PMER threshold 36%.
- Baseline decoding(row 2) utilizes the audio data and the decoding results(in place of the closed captions) of the previous iteration(baseline 0) as the new training set.

Baseline WMER(row 3) and baseline 1(row 4) use the audio data and the closed captions(with new PMER and AWD) as the new training set. However, baseline WMER sorts utterances based on WMER instead of PMER. In contrast, baseline 1 utilized the same method as described in 3.4.3. Consequently, using closed captions and sorting by PMER are the best; thus, the next iteration only experiments with

Figure 5.1: Average word duration of data selection iterations



closed captions as training set and sorting by PMER. Baseline 2 and 3 are the second and third(last) iteration respectively. They have the same value of WER; hence, the iteration stops because the data selection has converged.

5.2 Proposed models

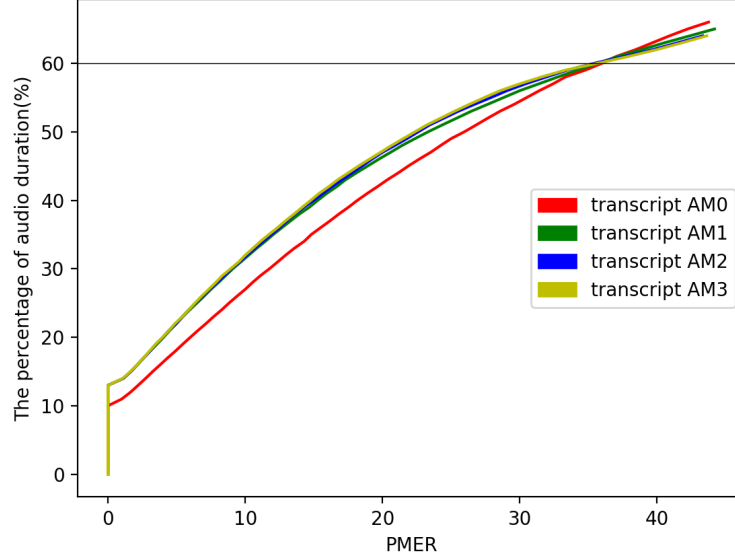
In addition to the baseline model, we proposed three new models which we hoped it would increase the performance. Table 5.2 demonstrates the result of our technique. Baseline 1 and 2 are the same model as shown in Table 5.1, while other rows show the results of the proposed models. We put baseline 1 and 2 to compare with the proposed models.

As explained in Section 3.5, we proposed three new models:

1. Average PMER and AWD of three ASRs.

By using the baseline acoustic model in baseline 1, three different ASRs were created by varying three different language models(as mentioned in 4.2.1.2). Each ASR from the three ASRs recognized the full training set and calculated new PMER and AWD. Each corresponding segment in three ASRs averages their PMER and AWD. The averaged training set is selected with the data selection method and trained to generate a new acoustic model. The new acoustic model with LM.7weeks+sub and the lexicon is evaluated as shown in model proposed

Figure 5.2: Phone matched error rate of the iteration



average(row 3). Because the model does not improve the performance as good as baseline 2, we do not continue to the next iteration.

2. Combine three ASRs.

We combined the recognition of three ASRs from Baseline 1 by using the algorithm explained in Listing 3.2, and then, we train a new acoustic model from the data. We evaluated the performance of the new model as shown in row 4. The value of WER of row 4(proposed combination 1) and WER of row 2(baseline 2) are the same; but, the PER decreases by 0.1. Therefore, we continue the iteration to see if it can further reduce the error rate.

Then, three ASRs of proposed combination 1 were combined and trained a new acoustic model based on the combination algorithm(Listing 3.2). We evaluated the new model and showed the result in row 5. Unfortunately, the performance of the model becomes poorer. The reason behind this is the combination algorithm tries to fix the bad closed captions by replacing with the decoding transcriptions if two ASRs agree(has exactly the same decoding results).

3. Utilize a randomly different random training set.

We select a different full training set, and then, we used the baseline model 2 to recognize the new training set. The closed captions of the new training set and their corresponding PMER and AWD are used for training a new acoustic

model. The new model is evaluated and written in row 6. It shows better performance by reducing WER as well as PER by 0.1. Again, we do another iteration, select a new training set randomly, and evaluate the new model(row 7). Unfortunately, the performance becomes worse compared to baseline 2. There is a luck factor in this technique. If we choose a very good training set in the next iteration, we may obtain a better performance model. Otherwise, the error rate can get worse as shown in row 7.

Table 5.2: WER and PER of baseline versus proposed model

No.	model name	WER(%)	PER(%)
1	baseline 1	41.2	33.8
2	baseline 2	40.5	33.3
3	proposed average	41	33.7
4	proposed combination 1	40.5	33.2
5	proposed combination 2	41.6	33.7
6	different data set 1	40.4	33.2
7	different data set 2	40.7	33.6

5.3 Execution time

Table 5.3 shows the execution time of each data selection iteration. The most demanding process is to train an acoustic model. The second most demanding process is to recognize the full training set and select a good subset of data from it. This is time consuming due to a massive volume of training data(around 200 hours of TV broadcasts). In overall, the execution time of each iteration roughly takes 43 hours.

Table 5.3: WER and PER of baseline versus proposed model

Iteration	Train AM	Recognize & data selection	Calculate WER	Decoding graph with LM.200	Decoding graph with LM.7weeks+subs	Total
0	21h	16h	1h	6m	5h	43h
1	24h	14h	1h	5m	3h	43h

5.4 Discussion

This chapter provides the results of experiments of baseline and proposed approaches. In this section, we give a summary of our experiments and offer possible explanations for the obtained results.

In baseline model, we showed a chart of PMER(see 5.2) produced by four iteration baseline models. The chart shows that the next iteration obtains more segments with zero PMER and lower PMER. The value of PMER decreases slightly in each iteration until it converges in the last iteration(iteration 2 and 3). The reason behind this is the data selection method has converged in Baseline 2. Furthermore, we compared the selected segments for training Baseline 2 and 3. They have exactly the same subset of selected segments/utterances(with slightly reduced PMER values in Baseline 3).

In addition to comparing PMER in the cumulative chart, we also evaluate each baseline model by computing word error rate of the development set. We did four iterations(baseline 0, 1, 2, and 3). The WER of the baseline model keeps decreasing in each iteration. In other word, we has shown that the data selection technique(inspired from [Lanchantin et al., 2016]) works well in our data. Model baseline 2 and 3 have the same WER which denotes that the data selection has converged. Moreover, we experimented by modifying the baseline technique.

The first modification is we tried to use the decoding result as the training set for the next iteration. The result is not as good as using our original technique(using closed captions as the training set). The intuition of the bad performance is the PER of baseline 0 is 35.8%; i.e. around 1 from 3 phones are incorrectly recognized. Therefore, the decoding result may be more terrible than the real transcriptions or even more terrible than the closed captions. The second modification is we tried to sort the segments by WMER in place of PMER. The result is worse than the baseline method. This is because we train an acoustic model which predicts sub phones(senone) from the acoustic inputs. Therefore, sorting by PMER is more likely to produce a better acoustic model(and of course better ASR because we use the same language model and lexicon) than sorting by WMER.

Beside the baseline model, we experimented with three different proposed models. The idea behind the proposed models is by combining three different ASRs(the same acoustic model but different language models). The first proposed model is by averaging PMER and AWD and do data selection the same as the baseline method. The result is slightly worse than the baseline. The second proposed model is by combining the recognition of three ASRs by the algorithm proposed in Listing 3.2. We obtained slightly better PER and exactly the same value of WER compared to Baseline 2. Thus, we continued the iteration which resulted in the bad performance model. The reasoning is our combination algorithm tried to fix the closed captions by replacing with the decoding results if two ASRs agree. We hypothesize that our three ASRs are not really different because we only utilize different language models(built with

somewhat the same text corpus) and the same acoustic model. Therefore, if one ASR makes a mistake, the other probably makes the same mistake. The last proposed model is to recognize randomly different training set for each iteration and do data selection from the training set. The first iteration shows better model than Baseline 2. However, the second iteration performs slightly worse than Baseline 2. We believe there is a luck factor in selecting the random training set. If the randomly selected set is well chosen, we obtain a better model; otherwise, the model will perform poorer.

In conclusion, we have experimented with the baseline model and we also proposed some methods. However, our proposed method does not work as well as we expected. The main reason behind it is we do not use and combine highly different automatic speech recognitions. Instead, we used the same acoustic model and three different language models, which are built from somewhat the same corpus. Therefore, it is necessary to experiment and investigate more the combination of totally different acoustic model(trained on different architecture, different acoustic input, etc) and totally different language model(different text corpus, different parameters, etc). Due to the limited time constraint of the internship(the execution time of one iteration is long, see 5.3), we do not have opportunity to experiment and investigate further.

Chapter 6

Conclusion and future works

6.1 Conclusion

In this internship, we explored baseline models and new proposed models of data selection for building automatic speech recognition. The models are inspired by lightly supervised technique. To build an automatic speech recognition system, one needs audio data and their corresponding exact transcriptions as the training set. However, we only have audio data and closed captions as our training data. Closed captions are closed to the real transcription, but not exactly the same. Therefore, we need to apply data selection to choose the subset of good data which are used to build the automatic speech recognition. This process is called lightly supervised technique.

Our baseline model starts from training an acoustic model from a random subset of the full training set. Then, the automatic speech recognition, consisting of the acoustic model, more constrained language model, and the lexicon, recognizes the full training set to produce new value of phone matched error rate(PMER) and average word duration(AWD). The new value of PMER and AWD are utilized to select data. This process is repeated over several iterations until the process does not improve. The process is evaluated by recognizing the development set because the development set has the exact transcriptions to measure word error rate. If the word error rate decreases in the next iteration, the iteration continues; otherwise, stop the iteration. We did the experiments for several iterations until the forth iteration(the word error rate of the third and the forth is the same). In addition, we experimented by modifying the baseline system. We selected the segments based on WMER and chose the decoding results as the training set. However, our baseline model(without modification) works better compared to the modified baseline models. Beside the baseline model, we proposed three new models. The idea behind this is to combine

three different speech recognizers by varying the language models and using the same acoustic model.

The first proposed model is averaging the PMER and AWD of the transcriptions of three recognizers. Our experiment shows that the average model does not perform better than the baseline. The second proposed model is combining three speech recognizers by our proposed algorithm. This process improves the model slightly; however, the second iteration of the model performs worse than the first iteration. This is because fixing the closed captions(replacing with the decoding results) makes the transcriptions to be worse than the closed captions. The last model is randomly choosing the training set for each iteration. Our last model appears to improve the WER slightly; nonetheless, it becomes worse in the next iteration because the random selection of training data involves luck. If the training set is well selected, the WER decreases; otherwise, the model performs poorer. To conclude, the combination system seems to perform slightly better than the baseline. However, we need to make stricter rules in our combination system to prevent worsening the transcriptions in the next iteration. Moreover, it is necessary to have quite different speech recognizers to maximize the improvement of the process.

6.2 Future works

Our work shows some promising directions in improving the lightly supervised approach. This section disseminates some suggestions for short term as well as long term future works.

6.2.1 Short term

One of the problem in our proposed models is when the combination algorithm fixes the closed captions by replacing with the decoding results when two speech recognizer agrees. This process tends to make the transcriptions to be worse compared to the closed captions and increases the error rate of the model. One of the solution is by making the rule to be stricter. We hope to be able to modify the combination algorithm(replacing the closed captions when three speech recognizers agree) and run more experiments to see if it decreases the error rate.

We also wish to experiment with more distinct three speech recognizers. Furthermore, varying the acoustic model(in term of architecture and corpus) and language model(corpus and type of language model) probably will increase the performance of our proposed models.

6.2.2 Long term

Training an automatic speech recognition requires a large amount of training data. However, manual transcription is expensive in term of time and manpower. There are unlimited supply of audio data in the internet, TV, and radio. Mostly, they have no corresponding exact transcriptions or even no transcription at all. One of the solution is to experiment with the TV broadcasts with closed captions(not the same as the detailed transcriptions); likewise, what we did in this work.

However, closed captioning is still time-consuming and costly. Some data, such as TV broadcasts, are usually closed-captioned, but other vast sources of data, such as: audio data from internet, are not closed captioned. These data are usually in massive volume and richer in variety and content. We need to research the way to harness these data to build a highly performing automatic speech recognition.

Our work here is only the early step to build a less supervised automatic speech recognition. We still utilized closed captions as a light supervision to build the automatic speech recognition. We hope that in the future we can research and harness the data without transcriptions at all for training a highly performing unsupervised automatic speech recognition.

Bibliography

- [Anden and Mallat, 2014] Anden, J. and Mallat, S. (2014). Deep scattering spectrum. *IEEE Transactions on Signal Processing*, 62(16):4114–4128.
- [Bahl et al., 1983] Bahl, L. R., Jelinek, F., and Mercer, R. L. (1983). A maximum likelihood approach to continuous speech recognition. *IEEE Trans. Pattern Anal. Mach. Intell.*, 5(2):179–190.
- [Barras et al., 2001] Barras, C., Geoffrois, E., Wu, Z., and Liberman, M. (2001). Transcriber: Development and use of a tool for assisting speech corpora production. *Speech Communication*, 33(1-2):5–22.
- [Bourlard and Morgan, 1993] Bourlard, H. A. and Morgan, N. (1993). *Connectionist Speech Recognition: A Hybrid Approach*. Kluwer Academic Publishers, Norwell, MA, USA.
- [Chan and Woodland,] Chan, H. and Woodland, P. Improving broadcast news transcription by lightly supervised discriminative training. In *2004 IEEE International Conference on Acoustics, Speech, and Signal Processing*. IEEE.
- [Dahl et al., 2011] Dahl, G. E., Yu, D., Deng, L., and Acero, A. (2011). Large vocabulary continuous speech recognition with context-dependent DBN-HMMS. In *2011 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE.
- [Gales and Young, 2007] Gales, M. and Young, S. (2007). The application of hidden markov models in speech recognition. *Found. Trends Signal Process.*, 1(3):195–304.
- [Hermansky and Sharma, 1999] Hermansky, H. and Sharma, S. (1999). Temporal patterns (traps) in asr of noisy speech. In *in Proc. ICASSP*, pages 289–292.
- [Jaitly et al., 2012] Jaitly, N., Nguyen, P., Senior, A., and Vanhoucke, V. (2012). Application of pretrained deep neural networks to large vocabulary speech recognition. In *Proceedings of Interspeech 2012*.

- [Jurafsky and Martin, 2009] Jurafsky, D. and Martin, J. H. (2009). *Speech and Language Processing (2Nd Edition)*. Prentice-Hall, Inc., Upper Saddle River, NJ, USA.
- [Kemp and Waibel, 1999] Kemp, T. and Waibel, A. (1999). Unsupervised training of a speech recognizer: Recent experiments. In *in Proc. EUROSPEECH*, pages 2725–2728.
- [Klakow and Peters, 2002] Klakow, D. and Peters, J. (2002). Testing the correlation of word error rate and perplexity. *Speech Commun.*, 38(1):19–28.
- [Klann and Szolovits, 2008] Klann, J. and Szolovits, P. (2008). An intelligent listening framework for capturing encounter notes from a doctor-patient dialog. In *2008 IEEE International Conference on Bioinformatics and Biomedicine Workshops*. IEEE.
- [Kneser and Ney, 1993] Kneser, R. and Ney, H. (1993). Improved clustering techniques for class-based statistical language modelling. In *Third European Conference on Speech Communication and Technology, EUROSPEECH*.
- [Kukich, 1992] Kukich, K. (1992). Technique for automatically correcting words in text. *ACM Computing Surveys*, 24(4):377–439.
- [Lamel et al., 2002] Lamel, L., Gauvain, J.-L., and Adda, G. (2002). Lightly supervised and unsupervised acoustic model training. 16:115–129.
- [Lanchantin et al., 2016] Lanchantin, P., Gales, M. J., Karanasou, P., Liu, X., Qian, Y., Wang, L., Woodland, P., and Zhang, C. (2016). Selection of multi-genre broadcast data for the training of automatic speech recognition systems. In *Interspeech 2016*. ISCA.
- [Lecouteux, 2006] Lecouteux, B. (2006). Imperfect transcript driven speech recognition. In *In: Proceedings of InterSpeech?06. Pitsburg*.
- [Maas et al., 2014] Maas, A. L., Qi, P., Xie, Z., Hannun, A. Y., Lengerich, C. T., Jurafsky, D., and Ng, A. Y. (2014). Building dnn acoustic models for large vocabulary speech recognition.
- [Mathias et al.,] Mathias, L., Yegnanarayanan, G., and Fritsch, J. Discriminative training of acoustic models applied to domains with unreliable transcripts. In *Proceedings. (ICASSP 05). IEEE International Conference on Acoustics, Speech, and Signal Processing, 2005*. IEEE.

- [Mesgarani et al., 2004] Mesgarani, N., Shamma, S., and et al. (2004). Speech discrimination based on multiscale spectro-temporal modulations. In *IN PROC. ICASSP*, pages 601–604.
- [Mohri et al., 2008] Mohri, M., Pereira, F., and Riley, M. (2008). Speech recognition with weighted finite-state transducers. In *Springer Handbook of Speech Processing*, pages 559–584. Springer Berlin Heidelberg.
- [Murveit et al., 1993] Murveit, H., Butzberger, J., Digalakis, V., and Weintraub, M. (1993). Large-vocabulary dictation using sri’s decipher speech recognition system: progressive search techniques. In *IEEE International Conference on Acoustics Speech and Signal Processing*. IEEE.
- [Nielsen, 2015] Nielsen, M. A. (2015). Neural networks and deep learning. <http://neuralnetworksanddeeplearning.com/>. Accessed: 2017-08-20.
- [Patel and Rao, 2010] Patel, I. and Rao, Y. S. (2010). Realtime speech processing for automated cue-generation. In *2010 INTERNATIONAL CONFERENCE ON COMMUNICATION CONTROL AND COMPUTING TECHNOLOGIES*. IEEE.
- [Peddinti et al., 2015] Peddinti, V., Povey, D., and Khudanpur, S. (2015). A time delay neural network architecture for efficient modeling of long temporal contexts. In *INTERSPEECH*.
- [Povey et al., 2011] Povey, D., Ghoshal, A., Boulianne, G., Burget, L., Glembek, O., Goel, N., Hannemann, M., Motlicek, P., Qian, Y., Schwarz, P., Silovsky, J., Stemmer, G., and Vesely, K. (2011). The kaldi speech recognition toolkit. In *IEEE 2011 Workshop on Automatic Speech Recognition and Understanding*. IEEE Signal Processing Society. IEEE Catalog No.: CFP11SRW-USB.
- [Poznanski and Wolf, 2016] Poznanski, A. and Wolf, L. (2016). CNN-n-gram for HandwritingWord recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE.
- [Renals et al., 1994] Renals, S., Morgan, N., Boulard, H., Cohen, M., and Franco, H. (1994). Connectionist probability estimators in HMM speech recognition. volume 2, pages 161–174. Institute of Electrical and Electronics Engineers (IEEE).
- [Rosenblatt, 1960] Rosenblatt, F. (1960). Perceptron simulation experiments. *Proceedings of the IRE*, 48(3):301–309.

- [Sak et al., 2014] Sak, H., Senior, A., and Beaufays, F. (2014). Long short-term memory based recurrent neural network architectures for large vocabulary speech recognition.
- [Shannon, 2001] Shannon, C. E. (2001). A mathematical theory of communication. *SIGMOBILE Mob. Comput. Commun. Rev.*, 5(1):3–55.
- [Silicon Intelligence, 2017] Silicon Intelligence (2017). Acoustic model. [Online; accessed July 25, 2017].
- [Stolcke, 2002] Stolcke, A. (2002). Srilm - an extensible language modeling toolkit. pages 901–904.
- [Stolcke et al., 2000] Stolcke, A., Wang, W., Vergyri, D., Ramana, V., Gadde, R., and Zheng, J. (2000). An efficient repair procedure for quick transcriptions. In *in Proc. ICSLP*, pages 1961–1964.
- [Team, 2017] Team, G. (cited August 2017). Grid5000:home. <https://www.grid5000.fr/mediawiki/index.php/Grid5000:Home>.
- [Thomas et al., 2009] Thomas, S., Ganapathy, S., and Hermansky, H. (2009). Phoneme recognition using spectral envelope and modulation frequency features. In *2009 IEEE International Conference on Acoustics, Speech and Signal Processing*. IEEE.
- [Vroegop, 2015] Vroegop, K. (2015). Javaone ? day 3 ? listen and watch. <https://technology.first8.nl/javaone-day-3-listen-and-watch/>. Accessed: 2017-08-20.
- [Waibel et al., 1990] Waibel, A., Hanazawa, T., Hinton, G., Shikano, K., and Lang, K. J. (1990). Readings in speech recognition. chapter Phoneme Recognition Using Time-delay Neural Networks, pages 393–404. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA.
- [Woodland et al., 2015] Woodland, P. C., Liu, X., Qian, Y., Zhang, C., Gales, M. J. F., Karanasou, P., Lanchantin, P., and Wang, L. (2015). Cambridge university transcription systems for the multi-genre broadcast challenge. In *2015 IEEE Workshop on Automatic Speech Recognition and Understanding (ASRU)*. IEEE.
- [Zavaliagkos and Colthurst, 1998] Zavaliagkos, G. and Colthurst, T. (1998). Utilizing untranscribed training data to improve performance.