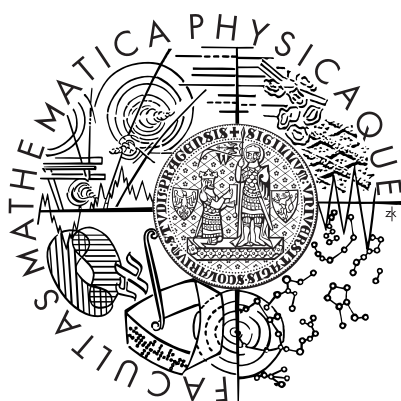


Charles University in Prague
Faculty of Mathematics and Physics

MASTER THESIS



Ondřej Plátek

Automatic speech recognition using Kaldi

Institute of Formal and Applied Linguistics
Supervisor: Ing. Mgr. Filip Jurčíček, Ph.D.
Study branch: Theoretical Computer Science

Prague 2013

I would like to thank my supervisor, Ing. Mgr. Filip Jurčíček, Ph.D., for his advice, guidance and for keeping me motivated. I would like to thank namely, Matěj Korvas for HTK scripts results and advice, Lukáš Žilka, David Marek and Ondřej Dušek for hacks in Vim, Bash and Perl, Marek Vašut for advice with shared library linking, Tomáš Martinec for C++ advice and Pavel Měnl for proofreading. I am also very grateful to the Kaldi team, which was very responsive and helpful. Expecially, Daniel Povey and Vassil Panayotov. Last but not least, I would like to thank my parents and Adéla Čiháková for all the help and support.

I declare that I wrote my master thesis independently and exclusively with the use of the cited sources. I agree with lending and publishing this thesis.

I understand that my work relates to the rights and obligations under the Act No. 121/2000 Coll., the Copyright Act, as amended, in particular the fact that the Charles University in Prague has the right to conclude a license agreement on the use of this work as a school work pursuant to Section 60 paragraph 1 of the Copyright Act.

Prague, August 1, 2013

Ondřej Plátek

Matematicko-fyzikální fakulta

Ústav formální a aplikované lingvistiky

Akademický rok: 2012/2013

ZADÁNÍ DIPLOMOVÉ PRÁCE

Jméno a příjmení: **Ondřej Plátek**studijní program: **Informatika**studijní obor: **teoretická informatika**

Děkan fakulty Vám podle zákona č. 111/1998 Sb. určuje tuto diplomovou práci:

Název práce: **Rozpoznávání řeči pomocí KALDI**

Zásady pro vypracování:

Jednou z důležitých komponent v dialogovém systému je modul rozpoznávání mluvené řeči. Tématem této práce bude seznámení se a využití open-source implementace výkonného rozpoznávače a systému trénování ASR Kaldi (<http://kaldi.sourceforge.net/>) pro dialogové systémy. Přestože KALDI již obsahuje ASR dekodéry, tak nejsou vhodné pro dialogové systémy z úvodu jejich malé optimalizace na rychlost a jejich velkého zpoždění v generování výsledku po ukončení promluvy. Proto hlavním cílem práce bude vyvinutí real-time rozpoznávače pro dialogové systémy minimalizující zpoždění a optimalizace na rychlost. Použité prostředky pro tuto optimalizaci mohou být například multi-vláknové dekodování nebo využití grafických karet pro obecné výpočty. Součástí této práce bude příprava akustického modelu a testování ve vyvíjeném dialogovém systému.

Seznam odborné literatury:

Psutka, J. and Müller, L. and Matoušek, J. and Radová, V. : Mluvíme s počítačem česky. p. 752, Academia, Prague, 2006.

Kaldi <http://kaldi.sourceforge.net/>Vedoucí diplomové práce: **Ing. Mgr. Jurčíček Filip, Ph.D.**

Navrhování oponenti:

Datum zadání diplomové práce: 8.10.2012

Termín odevzdání diplomové práce: dle harmonogramu příslušného akademického roku


.....
Vedoucí katedry


.....
Děkan

V Praze dne 14.11.2012

Univerzita Karlova v Praze
Matematicko-fyzikální fakulta
děkanát - studijní oddělení
121 16 Praha 2, Ke Karlovu 3
IČ: 00216208, DIČ: CZ00216208
tel.: 221 911 259, 221 911 111

Název práce: Rozpoznávání řeči pomocí Kaldi
Autor: Ondřej Plátek
Katedra: Ústav formální a aplikované lingvistiky
Vedoucí diplomové práce: Ing. Mgr. Filip Jurčíček, Ph.D.
E-mail vedoucího: jurcicek@ufal.mff.cuni.cz

Abstrakt: Tématem této práce je implementace výkonného rozpoznávače v open-source systému trénování ASR Kaldi (<http://kaldi.sourceforge.net/>) pro dialogové systémy. Kaldi již obsahuje ASR dekodéry, které však nejsou vhodné pro dialogové systémy. Hlavními důvody jsou jejich malá optimalizace na rychlost a jejich velké zpoždění v generování výsledku po ukončení promluvy. Cílem této práce je proto vyvinutí real-time rozpoznávače pro dialogové systémy optimalizovaného na rychlost a minimalizujícího zpoždění. Zrychlení může být realizováno například pomocí multi-vláknového dekodování nebo s využitím grafických karet pro obecné výpočty. Součástí práce je také příprava akustického modelu a testování ve vyvíjeném dialogovém systému "Vystadial".

Klíčová slova: ASR, rozpoznávání mluvené řeči, dekodér

Title: Automatic speech recognition using Kaldi
Author: Ondřej Plátek
Department: Ústav formální a aplikované lingvistiky
Supervisor: Ing. Mgr. Filip Jurčíček, Ph.D.
Supervisor's e-mail address: jurcicek@ufal.mff.cuni.cz

Abstract: The topic of this thesis is to implement efficient decoder for speech recognition training system ASR Kaldi (<http://kaldi.sourceforge.net/>). Kaldi is already deployed with decoders, but they are not convenient for dialog systems. The main goal of this thesis to develop a real-time decoder for a dialog system, which minimize latency and optimize speed. Methods used for speeding up the decoder are not limited to multi-threading decoding or usage of GPU cards for general computations. Part of this work is devoted to training an acoustic model and also testing it in the "Vystadial" dialog system.

Keywords: ASR, speech recognition, decoder

Contents

1	Introduction	3
1.1	The problem	3
1.2	The goals of the thesis	4
1.2.1	Training acoustic models	4
1.2.2	Development real-time speech recogniser	4
1.2.3	Integration into Alex SDS framework	4
2	Background	7
2.1	Automatic speech recognition	7
2.1.1	Speech parameterisation	8
2.1.2	Acoustic modelling	11
2.1.3	Language modelling	15
2.1.4	Speech decoding	15
2.1.5	Evaluating ASR quality	19
2.2	Hidden Markov Model Toolkit (HTK)	20
2.3	Julius decoding engine	21
2.4	Kaldi	22
2.4.1	Finite State Transducers	23
3	Acoustic model training	25
3.1	Vystadial acoustic data	25
3.2	Acoustic modelling scripts	26
3.3	Evaluation	28
3.3.1	Results	30
3.3.2	Kaldi and HTK comparison	31
4	Real time recogniser	33
4.1	OnlineLatgenRecogniser	33
4.1.1	<i>OnlineLatgenRecogniser</i> interface	34
4.1.2	<i>OnlLatticeFasterDecoder</i>	35
4.1.3	On-line feature pre-processing	36
4.1.4	Post-processing the lattice	37
4.2	PyOnlineLatgenRecogniser	38
4.3	Summary	39
5	Kaldi ASR in Alex SDS	41
5.1	Alex dialogue system architecture	41
5.2	Kaldi integration into SDS framework	43
5.2.1	<i>PyOnlineLatgenRecogniser</i> in Alex	43
5.2.2	Building in-domain decoding graph	44
5.3	Evaluation of <i>PyOnlineLatgenRecogniser</i> in Alex	45
6	Conclusion	49
A	Acronyms	51

B CD content **53**

The bibliography **54**

TODO: UNIT words with multiple spelling: real-time vs real time, online vs on-line, offline vs off-line

TODO: training: partial data evaluation

TODO: rewrite/deleteBatch interface vs oop on-line interface

TODO: cha:decoder: specify beam, max-active

TODO: add ml over whole thesis

TODO: dialog vs dialogue

*TODO: Fix "See Section XY". Tell what one can See in Section XY.
"Explain why one should look there".*

TODO: picture beam search

TODO: fix figure 3.2 (fig:parials) by introducing shapes not only colours

TODO: fix figure 3.3 (fig:parials_lm) bad axis description .. transform axis to %

1. Introduction

A spoken dialog is the most intuitive way of communication among people. The quality of a dialog largely depends on the quality of speech recognition because the reasoning and the answer is based on the recognised speech.

In this work, we build Automatic Speech Recognition for a dialog system called Alex. We see the added value of this thesis in:

- a new training scripts deployed with open acoustic data[17],
- but also in new C++ interface for speech recognition
- and its Python wrapper which is integrated into Alex dialogue system.

Training scripts, which use free publicly available data, evaluate the quality of trained Acoustic Models. We use Kaldi speech recognition toolkit[27] for acoustic modeling as well as for real-time recognition the textual representation from speech. The newly developed speech recogniser is deployed in the dialogue system Alex available at a public toll-free line 800 899 998.

1.1 The problem

The Automatic Speech Recognition (ASR) in a dialog system closely interacts with a Spoken Language Understanding unit. The Spoken Language Understanding (SLU) unit typically classifies the speech better if the speech recogniser outputs more than one hypothesis for one utterance. A word lattice effectively represents multiple hypothesis, so it is convenient for passing the hypothesis between ASR and SLU unit.

The Alex dialog system has used the HTK toolkit[46] and OpenJulius[18] lattice speech recogniser in order to train acoustic models respectively to decode lattices in real time. Unfortunately, our project members were experiencing crashes of OpenJulius during extracting lattices.

We were looking for another open source toolkit with a real-time speech recogniser because OpenJulius has a complicated source code and relatively slow development of both HTK and OpenJulius

We chose the Kaldi toolkit because its speech recognisers can produce high-quality lattices.[28] In addition, the Kaldi toolkit deploys modern training recipes, is actively maintained and is distributed under the permissive Apache 2.0 license¹. On the other hand, the speech recognisers in Kaldi do not support interface convenient for a dialog system which can process audio stream incrementally.

We developed the on-line interface suitable for a real-time use of a lattice recogniser, and we rewrote its feature preprocessing functionality to fit our interface. So far, the Kaldi developers focused on improving acoustic model training. In August 2012² Kaldi team published a demo version of an on-line one best hypothesis speech recogniser.

¹<http://www.apache.org/licenses/LICENSE-2.0>

²The changes were introduced by svn commit 1259.

1.2 The goals of the thesis

The goals of the thesis are presented in order as will be implemented:

1. Acoustic Models (AMs) will be trained to evaluate the new recogniser.
2. The new recogniser will be developed so its Python wrapper can be deployed into our dialogue system Alex.
3. Finally, we will integrate the recogniser into our Alex Spoken Dialog System (SDS) written in Python and evaluate its performance.

1.2.1 Training acoustic models

A Automatic Speech Recognition recogniser requires two pre trained components, an Acoustic Model and a Language Model. We focus on finding the best Acoustic Model for the Kaldi toolkit. The Language Model is changed dependently on targeted domain.

We aim at developing acoustic model training scripts using the Kaldi toolkit with such quality, that resulting AMs could be compared with the AMs trained with previously used HTK toolkit. The scripts will be developed for Czech and English transcribed acoustic data.

1.2.2 Development real-time speech recogniser

We should modify Kaldi speech recogniser in order to allow incremental speech recognition. The resulting incremental interface should be as simple as possible yet allow state-of-the art performance. In addition, we will implement such speech parametrisation and feature transformation preprocessing, so high-quality acoustic models can be used. Finally, we should compute the posterior probabilities of the word lattice representing multiple ASR hypotheses.

In addition, we may suggest potential speed improvements e.g. approximations, use of Graphics Processing Unit (GPU) or Deep Neural Networks (DNN) for speech processing[41].

1.2.3 Integration into Alex Spoken Dialog System framework

We should develop a thin wrapper which efficiently exposes the speech recognition interfaces into Python. We make sure that the lattices which are the output of the Kaldi recogniser can be also accessed from Python. The resulting recogniser should be integrated into Alex SDS and the decoding parameters are should be tuned to obtain best performance. The evaluation of speech recognition setup is important part of the integration.

Thesis outline

In Chapter 2 we introduce a fundamental theory of speech recognition for related areas to our work. In Sections 2.2 and 2.3 we describe alternatives

to Kaldi speech recognition toolkit. At the end of the chapter, we present OpenFST framework which allows the Kaldi library effectively implement many standard speech recognition operations. To obtain high-quality Acoustic Models, we develop training scripts for Czech and English data described in Chapter 3. In addition, we compare acoustic models trained by Kaldi and previously used HTK toolkit. Chapter 4 presents in detail the new Kaldi real-time recogniser and discuss its on-line properties. We distinguish the original work done by the Kaldi team and our improvements. Then in Chapter 5, we describe deployment of the real-time recogniser into dialogue system Alex, we suggest evaluation criteria and also evaluate the integrated recogniser accordingly. Finally, Chapter 6 summarises the thesis and concludes with future research directions.

2. Background

This chapter introduces basics of speech recognition related to this work. Section 2.1 introduces speech preprocessing, Acoustic Model (AM) and Language Model (LM) training, and explains important aspects of speech decoding. Next sections describe specific speech recognition software implementations. The Kaldi toolkit is described in Section 2.4, the HTK toolkit in Section 2.2 and the Julius decoder in Section 2.3.

The statistical methods for continuous speech recognition were established more than 30 years ago. The most popular statistical methods are based on acoustic modelling using Hidden Markov Models (HMMs) and n-grams LMs, which are also used in Kaldi, the toolkit of our choice. We introduce principles of speech recognition and present techniques which are used in Kaldi.

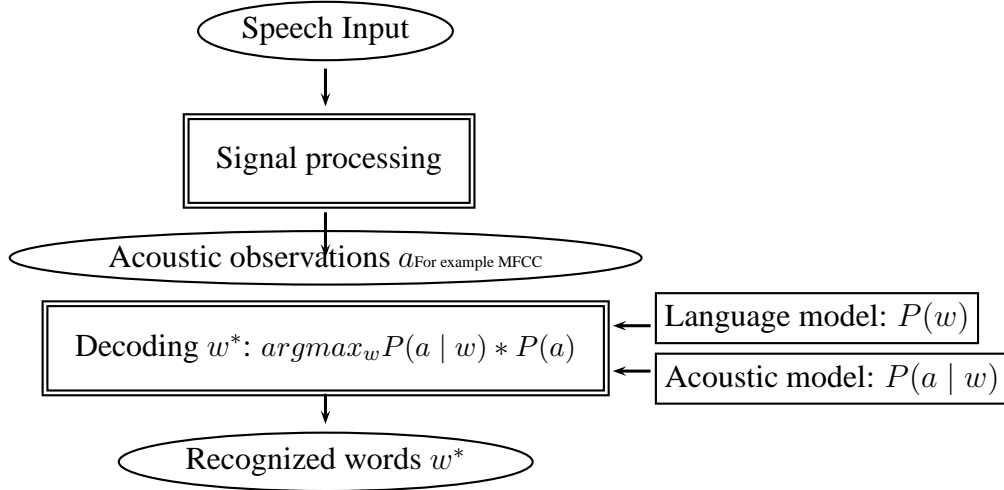


Figure 2.1: Architecture of statistical speech recognizer[23]

2.1 Automatic speech recognition

The goal of statistical ASR is to decode the most likely word sequence given speech. The term *decoding* finds its origin in HMM terminology. In speech recognition it is equivalent to *recognizing* the word sequence from the speech. Formally, we search for the most probable sequence of words w^* given the acoustic observations a as described in Equation 2.1. The best word sequence w^* does not depend on probability of the acoustic features $P(a)$ so it can be eliminated as shown on the second row of the equation.

$$\begin{aligned}
 w^* &= \operatorname{argmax}_w \{P(w | a)\} = \operatorname{argmax}_w \left\{ \frac{P(a | w) * P(w)}{P(a)} \right\} \\
 &= \operatorname{argmax}_w \{P(a | w) * P(w)\}
 \end{aligned} \tag{2.1}$$

The task of acoustic modelling is to estimate the parameters θ of a model so the probability $P(a \mid w; \theta)$ is as accurate as possible.¹ Similarly, the LM represents the probability $P(w)$.²

frames

The Figure 2.1 illustrate the process of decoding the most probable word hypotheses w^* for given speech utterance. First, the sampled audio signal is processed by speech parametrisation and feature transformations, so the decoding itself takes acoustic features a as input. The acoustic features a are computed on small overlapping windows of audio signal. The acoustic signal in one windows is known as a frame.

The decoding itself is performed time synchronously frame by frame using beam search. The beam search expands hypotheses from previous step by taking account the new frame features, and it computes probabilities of the expanded hypotheses using AM and LM. If the number of hypotheses exceeds the beam, the low probable hypotheses are pruned. After the decoding the last audio frame, all hypotheses represents whole utterance. The word labels w^* are extracted from the most probably hypothesis which survived the beam search.

Improving the accuracy of speech recognition engine depends mainly on improving AM and LM and also on parameters of the beam search such as threshold how many hypotheses are allowed at maximum.

2.1.1 Speech parameterisation

The goal of speech parameterization is to reduce the negative environmental influences on speech recognition. The speech varies in a number of aspects. Some of them are listed below:

- Differences among speakers pronunciation depends on gender, dialect, voice, etc.
- Environmental noises. In the dialogue system Alex where our ASR implementation is used the speech is typically recorded in a noisy street environment.
- The recorded channel. For example the telephone signal is reduced to frequency band between 300 to 3000Hz. The quality of mobile phone signal also influences the quality of the audio signal.

Different speech parametrisation may improve robustness of speech recognition for different recording conditions.

*acoustic
features*

Speech parametrisation extracts speech-distinctive acoustic features from raw waveform. The two most successful methods for speech parametrisation in last decades are Mel Frequency Cepstral Coefficients (MFCC)[8] and Perceptual Linear Prediction (PLP)[12]. Both MFCC and PLP transformations are applied on a sampled and quantized audio signal.³ For each window MFCC or PLP efficiently computes statistics with a reduced dimension. The

¹Acoustic modelling is described in Section 2.1.2.

²We describe language modelling in Section 2.1.3.

³In our experiments we use 16 kHz sampling frequency and 16 bit samples.

methods are very computationally effective and significantly improve the quality of recognised speech.

The toolkits used in our dialogue system, Kaldi and HTK toolkit, compute MFCC coefficients for given audio input in a similar way.⁴ Therefore, we choose MFCC as speech parametrisation technique for both toolkits, so we can compare them.

The MFCC statistics are computed for each frame. In Figure 2.2 there are 7 windows — frames with length of 25 ms and frame shift of 10 ms. The whole utterance lasts 85 ms.

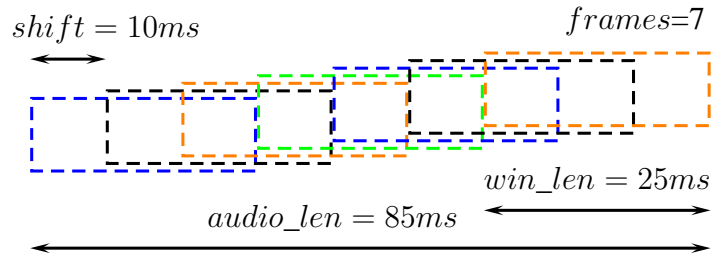


Figure 2.2: PLP or MFCC features are computed every 10 ms seconds in 25 ms windows. Audio length is $(frames - 1) * shift + win_len = 85ms$

Let us describe the MFCC computation for 25 ms window shifted by 10 ms and 16kHz audio sampling frequency. The $16000 * 0.025 = 400$ samples in one window are reduced to 13 static cepstral coefficients.

The MFCC static features are usually extended by time derivatives $\Delta + \Delta\Delta$ features [34]. As a result, MFCC $\Delta + \Delta\Delta$ extracts $13 + 13 + 13 = 39$ acoustic features for one frame. The original vector of 400 audio samples in one frame is reduced to vector of 39 MFCC $\Delta + \Delta\Delta$ acoustic features.

The MFCC features are computed by the following steps:

1. The audio samples are transformed into *frequency domain* by Discrete Fourier Transformation (DFT) in the window.
2. The frequency spectrum from the previous step is transformed onto the mel scale, a perceptual scale of frequencies, using triangular overlapping filters.
3. From the mel frequencies the logs of the powers are taken from each of the mel frequencies.
4. At the end the discrete cosine transform is applied on the list of mel log powers.
5. The MFCC coefficients are the amplitudes of the resulting spectrum.
6. The $\Delta + \Delta\Delta$ coefficients are computed from the current and previous static features. See Figure 2.3.

⁴The subtle differences are caused by implementation approaches, but does not effect the quality of MFCC coefficients in significant way.

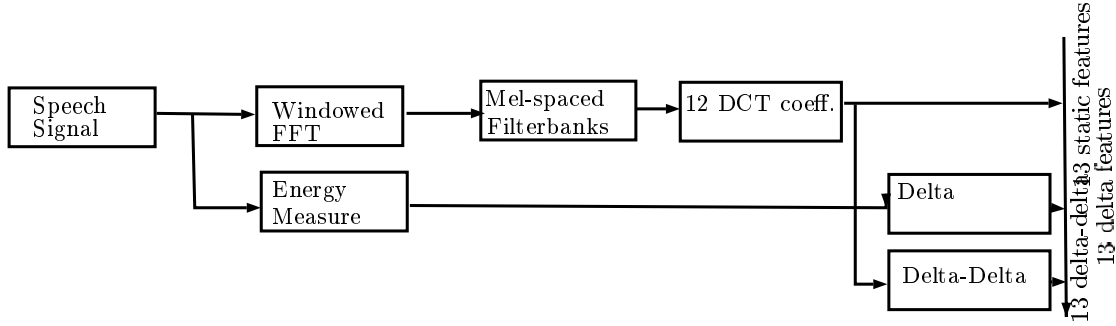


Figure 2.3: Typical setup with 39 features using MFCC.

Feature space transformations

Feature space transformations are usually applied in addition to MFCC or PLP parametrisation. The feature space transformations are also typically applied per frame, but they usually take into account context of several preceding (left context) and consecutive frames (right context).

The linear or affine transformations are expressed by matrix multiplications Ax respectively Ax^+ . The matrix A represents the transformation. The x is the input vector and Ax are the transformed features. The affine transformations uses extended vector $(x^+)^T = (x_1, \dots, x_n, 1)$ and matrix $A : (n+1) * (n+1)$.

There is a large variety of available transformations. Dependently on acoustic data one should choose the most appropriate one. Some transformations are estimated discriminatively, some use generative models. Some are speaker dependent, some speaker independent.

We list some of Kaldi transformations in order to illustrate rich choice of feature transformations in Kaldi toolkit.

- Heteroscedastic Linear Discriminant Analysis (HLDA)[9].
- Linear Discriminant Analysis (LDA)[11] is typically used with MLLT for speaker independent training.
- Maximum Likelihood Linear Transform (MLLT) also known as Semi-Tied Covariance (STC)[11]
- Exponential Transform (ET)[32]. It uses small number of speaker specific parameters for adaptation on speaker.
- Cepstral Mean and Variance Normalisation (CMVN)[21]. Typically normalise the cepstrum mean and variance per speaker.

In our acoustic modelling scripts, see Chapter 3, we use two non-speaker adaptive feature transformations, which can be computed with very small context. The first transformation, $\Delta + \Delta\Delta$ for MFCC coefficients, was already introduced. The second transformation, LDA and MLLT, is described briefly below.

Linear Discriminant Analysis and MLLT feature transformation

The LDA+MLLT is an alternative setup to $\Delta + \Delta\Delta$ transformation in our training scripts. We use it also on top of MFCC coefficients. Using several spliced MFCC vectors the LDA+MLLT searches for the best dynamic transformation.

The combination of LDA and MLLT applies the feature transformation in two steps: LDA reduces the feature dimension and MLLT applies linear simple transformation[11]. Whereas, the HLDA estimates dimension reduction and space transformation in one step.[9] The combination LDA and MLLT performs very similar feature transformation to HLDA and gains significant improvements over $\Delta + \Delta\Delta$ transformation similarly as HLDA[9][11].

2.1.2 Acoustic modelling

Acoustic modelling is arguably the heart of speech recognition. The AM estimates the probability $P(a|w;\theta)$ of generating acoustic features a for given words w and thus directly affects speech recognition quality as seen in Equation 2.1.

Acoustic modelling has only partial information available for training AM parameters θ because the corresponding textual transcription is time-unaligned. The hidden information of the word (time) alignment in a utterance makes acoustic model training more challenging. Modern speech recognition toolkits use Hidden Markov Model for modelling uncertainty between acoustic features and the corresponding transcription.

Choice of training units

The most successful acoustic modelling methods do not estimate the $P(a|w)$ directly, but estimate probability $P(a|f_1f_2f_3f_4)$ of generating acoustic features a for phones $w = f_1f_2f_3f_4$ which forms the pronunciation of the word w . Moreover, the triphones are used even more successfully for estimating probability of acoustic features give word pronunciation.

Phone is the smallest contrastive unit of speech. Let us see few examples of words and their phonetic transcriptions according CMU dictionary[42].

phone

- *youngest* & Y AH1 NG G AH0 S T
- *youngman* & Y AH1 NG M AE2 N
- *earned* & ER1 N D
- *ear* with two transcribed pronunciations IY1 R and IH1 R

The CMU dictionary distinguishes among several variations for each vowel e.g. AH1 and AH0. It also stores two possible pronunciations for the word *ear*.

The acoustic features for a phone significantly depend on its context. The previous and the following phone strongly influence the sound of the middle phone.

The triphone is a sequence of three phones and captures the context of single phone. As a result, acoustic properties of the triphones vary much less according to the context than phones. Let us note that certain combinations of prefixes have the same effect on the central phone, e.g. *q* and *k* has the same effect on *i*. In order to reduce the number of triphones for acoustic modelling, these triphones are clustered together.

triphone

Hidden Markov Models (HMMs)

The HMMs is a very powerful statistical method for characterizing observed data samples of a discrete-time series with an unknown state. [13]. In case of speech recognition the hidden states typically represent monophones or triphones and we observe samples of the acoustic features.

Hidden Markov Models have two type of parameters *transition probabilities among states* and *probabilistic distribution for generating observation in given state*. These parameters need to be estimated in AM training.⁵

transition probability

The transition probability is a probability of changing state q to state u . Each transition is represented as arc $e = qu$ between the states q and u , see 2.4. The probability is typically represented as the weight w_e of arc e .

Importantly, an HMM use self loop arc $e = uu$ for all states to model acoustic features which are generated several times from the same state u . As a result, an HMM is able to model variable length of phones.

The Markov model emits an observation during traversal over its arcs. The Hidden Markov Model emits the observation stochastically based on the probabilistic distribution related to the visited state. In speech recognition, a multivariate Gaussian distribution is typically used to model observation probabilities of HMM states. The Gaussian distribution models probability of emitting acoustic features in given state. The parameters of the Gaussian distribution are estimated for each state individually. However, the states are usually clustered during AM training and the states within a cluster share same parameters to the Gaussian distribution.

Training HMM

The Kaldi uses Viterbi training and the HTK toolkit uses Expectation Maximization algorithms to train HMM Acoustic Model. The toolkits models the observation probabilities using multivariate Gaussian distribution with dimension of the acoustic features a .

Typically, the transition probabilities are initialised with values uniformly distributed. The observation probabilities are usually initialized by multivariate Gaussian distribution with μ and Σ set to global mean and global covariance matrix estimated on all training acoustic data.

Let us describe how the Expectation Maximization (EM) algorithm operates for one pair of training data consisting of acoustic features a and corresponding text speech transcription t . We create HMM t' , where each state represents one monophone.⁶ The monophones are extracted from transcription t using pronunciation dictionary. In Figure 2.4 the utterance *how do you do* was expanded to monophone HMM model. Given the HMM model for transcription t and acoustic features a the parameters of the model are estimated. It should be obvious that only states representing phones in transcription can be trained by training pair (a, t) . Consequently, one needs lot of training data to robustly estimate parameters of all states.

The EM algorithm iterates following steps in order to update parameters of transition and observation probabilities:

- The observation probabilities are computed using HMM t' .
- **E-step**: Based on the observation probabilities the observation are align to states of HMM t' .
- **M-step**: Based on the alignment of observation to states the t' parameters are re-estimated.

The **E-step** finds a distribution for the alignment between HMM t' and tran-

⁵Both kind of parameters are denoted together as θ in Equation 2.1.

⁶We describe the identical training procedure for simplicity on monophones. The state-of-the-art AMs use triphones.

scription t using Maximum Likelihood Estimation (MLE)[11] and observation probabilities. MLE takes into account all possible alignments and its probabilities to compute the resulting distribution. The Baum-Welsh equations can be derived from the fact, that the MLE criterion is also used for finding the most probable distribution in **M-step**. [13]

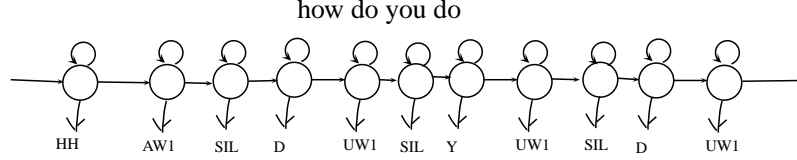


Figure 2.4: Markov monophone model for four words. Such an HMM is constructed for monophone Viterbi training and reference transcriptions *how do you do*. The parameters of the HMM model are updated according Equation 2.9, 2.10 and 2.11.

Maximum Likelihood Estimation method

The MLE is a general approach to setting statistical model parameters. It searches for best parameters θ^* in order to maximize the likelihood function f for Independent and Identically Distributed (IID) training data illustrated in Equation 2.4. For IID training data holds Equation 2.2 describing data joint probability.

$$f(x_1, x_2, x_3, \dots, x_n | \theta) = f(x_1 | \theta) * f(x_2 | \theta) * \dots * f(x_n | \theta) \quad (2.2)$$

The likelihood function can be derived from Equation 2.2 assuming training data fixed and parameter θ free as described in Equation 2.3.

$$\mathcal{L}(\theta | x_1, \dots, x_n) = \sum_{i=1}^n \log(f(x_i | \theta)) \quad (2.3)$$

$$\theta^* = \operatorname{argmax}_{\theta} \mathcal{L}(\theta | x_1, \dots, x_n) \quad (2.4)$$

Viterbi training of acoustic models

On the other hand, the Kaldi toolkit applies the Viterbi criterion in assigning the acoustic observation to HMM states. The Viterbi training approximates EM algorithm by choosing single best alignment and maximizing the posterior probability for the chosen alignment. Latest work suggest that Viterbi training is just as effective for continuous speech recognition as Baum-Welch algorithm [36]. Moreover, Viterbi training needs much less computational resources.

We detail the Viterbi training since it is used in the Kaldi toolkit for acoustic model training and also a very similar algorithm is used for Viterbi decoding.

Given set of training observations $O^r, 1 \leq r \leq R$ and HMM state sequence $1 < j < N$ the observation sequence is aligned to the state sequence via Viterbi alignment.[5] The best alignment T results from maximising Equation 2.5 for $1 < i < N$.

$$\phi_N(T) = \max_i [\phi_i(T) a_{iN}] \quad (2.5)$$

The $\phi_i(o_t)$ from Equation 2.5 is computed recursively according Equation 2.6

$$\phi_i(o_t) = b_j(o_t) \max \left\{ \begin{array}{l} \phi_j(t-1) a_{jj} \\ \phi_{j-1}(t-1) a_{jj} - 1 \end{array} \right. \quad (2.6)$$

The initial conditions are $\phi_1(1) = 1$ and $\phi_j(1) = a_{1j} b_j(o_1)$, for $1 < j < N$. In our case the likelihoods are modeled as mixture Gaussian densities, so the output probability $b_j(o_t)$ is defined as in Equation 2.7.

$$b_j(o_t) = \sum_{m=1}^{M_j} c_{jm} \mathcal{N}(o_t; \mu_{jm}, \Sigma_{jm}) \quad (2.7)$$

The M_j represents number of mixture components in state j , c_{jm} is the weight of m^{th} component and $\mathcal{N}(o_t; \mu_{jm}, \Sigma_{jm})$ is multivariate Gaussian with mean vector μ and covariance Σ .

Firstly, model parameters are updated based on the single-best alignment of individual observation to states and Gaussian components within states. Secondly, transition probabilities are estimated from the relative frequencies, Equation 2.13 where A_{ij} denotes the number of transitions from state i to state j .

$$\hat{a}_{ij} = \frac{A_{ij}}{\sum_{k=2}^N A_{ik}} \quad (2.8)$$

The indicator function $\psi_{jm}^r(t)$ is used for updating means and covariance matrix from statistics. It returns one if o_t^r is associated with mixture component m of state j and is zero otherwise. The mean vector and covariance matrix is updated according Equations 2.9 and 2.10.

$$\hat{\mu}_{jm} = \frac{\sum_{r=1}^R \sum_{t=1}^{T_r} \psi_{jm}^r(t) o_t^r}{\sum_{r=1}^R \sum_{t=1}^{T_r} \psi_{jm}^r(t)} \quad (2.9)$$

$$\hat{\Sigma}_{jm} = \frac{\sum_{r=1}^R \sum_{t=1}^{T_r} \psi_{jm}^r(t) (o_t^r - \hat{\mu}_{jm})(o_t^r - \hat{\mu}_{jm})'}{\sum_{r=1}^R \sum_{t=1}^{T_r} \psi_{jm}^r(t)} \quad (2.10)$$

Finally, the mixture weights are computed based on the number of observations allocated to each component.⁷

$$c_{jm} = \frac{\sum_{r=1}^R \sum_{t=1}^{T_r} \psi_{jm}^r(t)}{\sum_{r=1}^R \sum_{t=1}^{T_r} \sum_{l=1}^M \psi_{jl}^r(t)} \quad (2.11)$$

generative training To conclude, AM are trained using MLE or Viterbi training, which approximates the theoretically optimal MLE Baum-Welsh training; however, in practice Viterbi training performs as well as MLE modelling. Baum-Welsh or Viterbi training aim at modelling likelihood of spoken utterance and perform so called generative training. However, discriminative methods, which re-estimates generative AMs, perform better.

Discriminative training

discriminative training The discriminative training uses its objective function and likelihood of generative models to discriminate – boost differences between high probable and low probable hypotheses. The discriminative training is typically initialised by acoustic generative model from Baum-Welsh or Viterbi training. Then, the likelihood from

⁷The Viterbi equations has the same notation as in [5].

the generative model is boosted according an objective function and the AM is re-estimated. Following objective functions and accordingly named discriminative training methods are used in our training scripts:

- Maximum Mutual Information[7]
- Boosted Maximum Mutual Information[29]
- Minimum Phone Error[26]

For details how the methods are initialized and its usage in Kaldi see Chapter 3.

2.1.3 Language modelling

A Language Model effectively reduces and more importantly prioritise the AM hypothesis. A probability of acoustic features given words transcription $P(a|w)$ estimated by AM is combined with the probability of the words transcription $P(w)$ estimated by LM for given domain in order to compute posterior probability of transcription $P(w|a) = \frac{P(a|w)*P(w)}{P(a)}$.

The statistical LM assigns a given word sequence its probability according Equation 2.12. The most used, n-gram LMs compute the probability of k word sequence W according Equation 2.12.[?] The Markov assumption approximates the probability by assuming that only the most recent $n - 1$ words are relevant when predicting next word. We call the number n an order of LM.

LM order

$$P(W) = P(w_k, w_{k-1}, w_{k-2}, \dots, w_1) \approx \prod_{i=1}^k P(w_i | w_{i-n+1}^{i-1}) \quad (2.12)$$

The probabilities $P(w_i | w_{i-n+1}^{i-1})$ for each word w_i are estimated using relative frequencies of the n-grams, (n-1)-grams, (n-2)-grams, ... on training data. The MLE is used for estimating estimating relative frequencies r according Equation 2.13.

$$r(w_i | w_{i-n+1}^{i-1}) = \frac{f(w_{i-n+1}^i)}{f(w_{i-n+1}^{i-1})} \quad (2.13)$$

The Equation 2.13 is intuitive but many valid and even reasonable utterances are missing or too few. Consequently, the numerator might be zero and the relative frequency may be undefined. This is known as sparse data problem. Smoothing techniques are often used to estimate the higher n-gram relative frequencies based on the lower frequencies.[?]. In principle, the predictive accuracy of the language model can be improved by increasing the order of the n-gram. However, doing so further exacerbates the sparse data problem.[?]

sparse data

*LM
smoothing*

The LM estimates the probability by counting the relative frequencies on text corpus which is typically chosen according the targeted ASR domain. For example, in training scripts which are described in Chapter 3 we train the LM only on text transcriptions from the training data using Witten-Bell smoothing.[43]

2.1.4 Speech decoding

The speech HMM decoders find the most probable word sequences by searching phone sequences which corresponds to the words. The phones are typically represented as triphones in AM.

Using combination of AM and LM probabilities as described in Equation 2.1 does not produce the most accurate speech transcriptions. Typically a Language

Model Weight (LMW) w_{lm} is used to improve speech recognition accuracy. It is tuned on development set and balances the impact of the two models. Using the LMW the words best sequence is found according Equation 2.14.

$$w^* = \operatorname{argmax}_w \{P(w \mid a)\} = \operatorname{argmax}_w \{P(a \mid w) * P(w)^{w_{lm}}\} \quad (2.14)$$

The ASR is a pattern recognition task as well as a search problem. In speech recognition, making a search decision is also referred to as decoding.[13]

For a word recognition the AM limits the possible phone sequences only to words in lexicon — the words in training data. The word recognition is nowadays the most successful form of ASR. The HMM sequences represent phone sequence which forms words only as illustrated on Figures 2.5 and 2.6. The words are connected via HMM model which represent inter word silence.

For isolated word recognition the HMMs are evaluated for each word possibility. Using the forward algorithm for each HMM h_w , we are able to compute the probability of every word w given the acoustic observations. The isolated word recognition becomes a simple recognition problem, where we select the most probable HMM h^* from a finite set of word HMMs.

Note that the HMM training is identical for continuous ASR and isolated word recognition, but the decoding is more complicated for continuous ASR where we aim to decode word sequences.

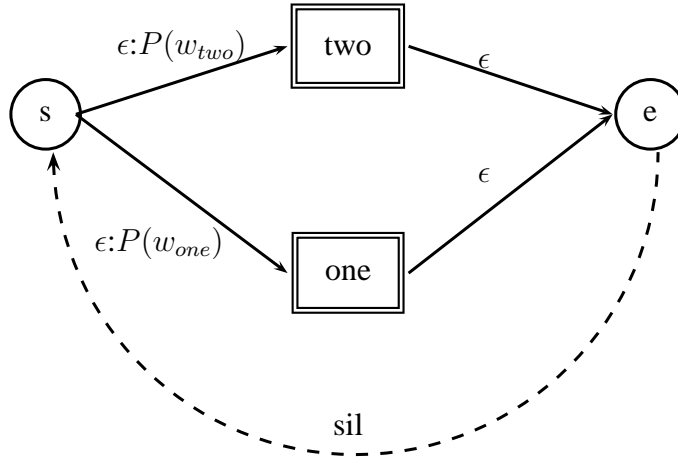


Figure 2.5: Diagram of how LM is combined with HMMs.

Let us introduce simple example of continuous word ASR. Imagine a LM of order 1 modelling only two words - *one* and *two* each uniformly distributed⁸. We want to decode any possible sequence of words *one*, *two*. The ϵ transition at the end of words HMMs to final state e allow us introduce HMM silence model which connect the final state e and start state s . Consequently, the words are chained using silence HMM model as illustrated in Figure 2.5. Expanded monophone HMM for words *one* and *two* is shown in Figure 2.6. Note that LM weight $P(w)$ can be stored on the in the ϵ transition at the beginning.

Even the simple HMM network in Figure 2.6 can become large search problem partially because the search space of words grows exponentially in number of words in utterance and partially because the word boundaries are unknown. Each

⁸If a LM of order 1 assigns to every word equal probability, we say it has order 0

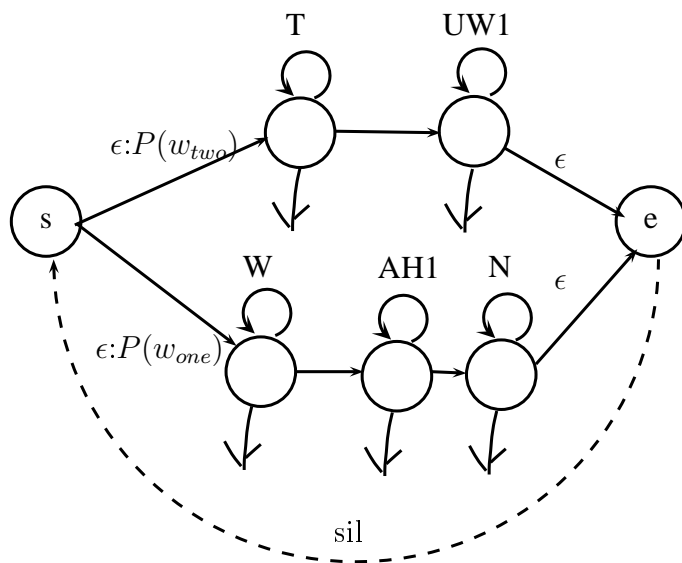


Figure 2.6: Expanded HMMs for words *one* and *two*. The arrows at HMM states illustrate that every observation of acoustic features can be generated according to the statistical distribution. Note that if a speaker says *two* a HMM model with well trained parameters should output higher probability for HMM representing *two*. For longer speaker sequences e.g. *two one one two ...* the HMMs are connected over the ϵ transitions, and a search is used for selecting the most probable sequence.

word can begin in any moment and last with decreasing probability ad infinity, so the search space explode. In addition, higher order LMs increase the decoding complexity even more.

One can see that the search space of speech recognition problem is enormous and still has to be solved in very short time for real-time applications.

Natural choice for one-best hypothesis is Viterbi beam search[13]. It uses a dynamic programming to compute the new best partial hypotheses for new audio data based on partial hypotheses from previous step. *TODO: picture beam search*

The Viterbi algorithm is breadth-first search algorithm and a beam is used to limit number of nodes, which are expanded from current set of nodes to next iteration. We list few alternatives how to set up the beam for speech decoding.

- *Fixed beam* guarantees maximum size of memory footprint and fast decoding.
- *Relative one-best hypothesis comparison* effectively discards most of the improbable hypothesis if the one-best hypothesis is significantly better than alternatives and keeps lot of alternatives if one-best hypothesis is weak. The relative one-best hypothesis comparison naturally broaden the beam in uncertain region, but does not guarantee no hard limits e.g., maximum number of nodes expanded.
- *Combination* of methods applies the strictest criteria on beam in each iteration.

Numerical stability

The hypotheses which are represented as the paths of states are typically rather long in the search graph and lot of hypothesis are assigned with tiny probabilities. In

order to keep the numeric stability, the probabilities are expressed in a logarithmic arithmetic.

In order to use the shortest distance measure to find the most probable path we use formula 2.15 derived from equation 2.14. The $C(a | w)$ and $C(w)$ are costs with range between zero and one, where cost of one corresponds to zero probability $C(1) \cong P(0)$ and cost of infinity corresponds to one probability $C(\infty) \cong P(1)$.

$$\begin{aligned} w^* &= \operatorname{argmin}_w \left\{ \log \left(\frac{1}{C(a | w) * C(w)^{w_{lm}}} \right) \right\} \\ &= \operatorname{argmin}_w \left\{ -\log(C(a | w) * C(w)^{w_{lm}}) \right\} \\ &= \operatorname{argmin}_w \left\{ -\log(C(a | w)) - w_{lm} * \log(C(w)) \right\} \end{aligned} \quad (2.15)$$

Decoding formats

The one-best hypothesis outputs only single sequence of words despite the fact that other sequences of words are often almost as probable as the best hypothesis. Formats which are able to represent alternative hypothesis provide better results for further processing than one-best hypothesis because the alternatives may cover almost all probable hypotheses. We present n-best list and lattice formats, which both are able to represent alternative hypothesis.

N-best list is an extension to the one-best hypothesis format. In n-best list is included apart from the most probable word sequence, also the second, third, ..., n-th most probable hypothesis.

0.5 hi how are you
0.2 hi where are you
0.1 bey how are you

Figure 2.7: Example of 3-best list output with posterior probability for each path. N-best list in Kaldi can be easily extracted from lattices. Corresponding example lattice is in Figure 2.8.

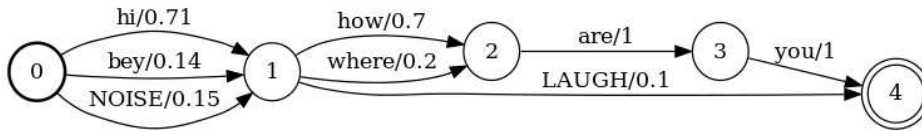


Figure 2.8: Word posterior lattice. Common parts of hypotheses are effectively represented. All outgoing arcs for each node sum to 1.0.

lattice

A lattice is a convenient type of ASR output. It effectively represents the alternative hypotheses by sharing their common parts. Example of word lattice in Figure 2.8 shows word posterior lattice. Each hypothesis is represented as sequence of arcs from starts to final node. The words and their weights are associated with the arcs. The posterior probability of a hypothesis is computed as a product of posterior probabilities of each word in the hypothesis.

It is useful to capture how the quality of each hypothesis contrasts to its alternatives or even provide an absolute quality measure. Typically likelihood and posterior probability is associated with each word sequence to express its quality.

The likelihood measure can be used only for relative comparison, whereas posterior probability is an normalised absolute measure. For some applications the likelihood measure is sufficient, other applications for example dialogue systems prefer posterior probabilities. Note that posterior probabilities for n-best lists typically do not sum to one and may need to be renormalised, because n-best list omits some hypothesis which are used to compute the posterior probability in lattices. See Figure 2.7 for such example for 3-best list.

2.1.5 Evaluating Automatic Speech Recognition quality

The accuracy of a speech recognizer is typically measured using Word Error Rate. The Word Error Rate (WER) measure is computed on one-best ASR hypotheses and their human transcriptions. The WER is computed as a minimum edit distance on words between the ASR output and reference transcription. Following edit operations are used *substitution*, *deletion*, *insertion* to compute the minimum edit distance as illustrated in 2.16. The effective implementation for computing WER uses dynamic programming and is not computationally intensive because ASR hypotheses are typically quite short.

WER

$$WER = 100 * \frac{\min_dist(decoded_{AM,LM}(a), t, edit_operation = \{Subs, Del, Ins\})}{\# \text{ words in } t} \quad (2.16)$$

Note that WER is an error function so the ideal value is zero because for $WER = 0$ the one-best hypothesis $decoded(a)$ and the reference transcription t are identical. The WER value of 100 show that every single word is different between $decoded(a)$ and reference t if the number of words in ASR output and reference are equal. Despite the fact that WER resembles percentage format, it can be bigger than 100. See the third example in Figure 2.9.

reference

```

decoded(a) = 'hi hi hi hi'
t='hi hi ha ha'
WER = 100 * ( 2 / 4) = 50

decoded(a) = 'how do you do'
t='how do you do''
WER = 100 * ( 0 / 4) = 0

decoded(a) = 'hi hi hi hi'
t='hello'
WER = 100 * ( 4 / 1) = 400

```

Figure 2.9: WER captures the ASR one-best hypotheses accuracy.

Note that the data used for evaluation should not be used in AM training because we are evaluating the ability to decode unknown speech. We should also measure the ASR quality on speech from a speaker who does not appear in speech training data because we usually want to decode speech of an unheard speaker.

Alternative measures

The Sentence Error Rate measures how many decoded utterances $decoded(a)$

SER

match exactly its reference t for all pairs (a, t) in test set T .

$$SER = \frac{\sum_{\{(a,t) \in T; decoded(a)=t\}} 1}{|T|} \quad (2.17)$$

*oracle
WER*

If the n-best list or lattice is used the one-best hypothesis is extracted to compute WER or Sentence Error Rate (SER). On the other hand, we are using n-best lists or lattice because the one-best hypothesis might be wrong and the alternative hypothesis may be closer to the reference. In order to evaluate quality of alternative hypotheses one may use oracle WER which reports the WER of the best hypotheses in n-best list or in lattice. The lattices with rich alternatives gain much lower oracle WER than short n-best lists or even one-best hypotheses. The rich alternatives contain additional information and for example a dialogue system Spoken Language Understanding component may exploit the alternatives.

Measuring speed

In this thesis we are especially concerned about the speed of speech decoding because the implemented decoder is used in a real-time Spoken Dialog System.

*Real Time
Factor*

A very natural measure of a speech decoding speed is Real Time Factor, which expresses how much the recognizer decodes slower than the user speaks. We measure the Real Time Factor (RTF) for each recording as described in Equation 2.18.

$$RTF = \frac{time(decode(a))}{length(a)} \quad (2.18)$$

For real-time decoding in a dialogue system we need smaller than one $RTF < 1.0$. In other words, the decoding of an utterance should take less time than a user needed for pronouncing the utterance. With $RTF < 1.0$ the hypothesis is decoded immediately after the user finishes the speech.

The decoding is performed while user is speaking, but extracting the ASR hypothesis output is triggered at the very end of the speech. The users have to wait at least the time when the ASR hypotheses is extracted.

latency

In real-time SDS the critical measure is a delay how long the user has to wait for its answer. The latency measures the time between the end of the user speech and the time when a decoder returns the hypothesis, which is the most important speed measure for ASR component in SDS. Note if $RTF < 1.0$ then the latency corresponds to time of ASR hypotheses extraction.

2.2 HTK

The HTK toolkit is a set of command line tools, sample scripts and library for training and decoding HMM focused on speech recognition. With the toolkit are distributed two decoders *HVite* and *HDecode*, which are not designed for real-time applications.

Functionality of the core library can be accessed through command line executables. The command line programs are typically combined in training scripts to train acoustic and language models. In Figure 2.10 the acoustic models are labelled as "HMMs" and the language models in HTK are represented in "Networks". The trained models are used in one of HTK decoders e.g. *HVite* for decoding transcriptions, which can be evaluated using *HResults*.

The HTK library use Baum-Welch algorithm to train acoustic models. The *HVite* decoder uses token passing algorithm and Viterbi criterion.[47] Only unigram

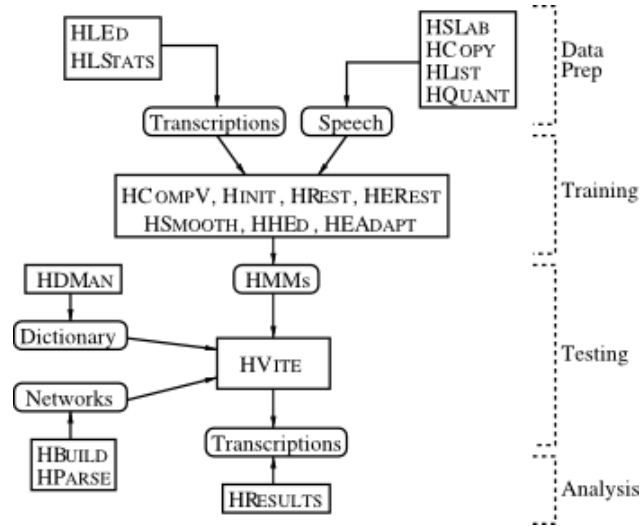


Figure 2.10: Figure 2.2 from HTK Book 3.4[47]

and bigram LMs can be used with *HVite*. The termHDecode decoder can handle bigram or trigram language models.

Let us stress that we use high quality Bash and Perl scripts for training HTK AM from Vertanen improved by Matěj Korvas.[40][17]

The HTK toolkit is licensed under a special license⁹. The *HDecode* has very similar license condition but can be only used for research purposes.¹⁰

2.3 Julius decoding engine

Julius is a large vocabulary continuous speech decoder which can use AM in HTK format for decoding.[18] Julius is BSD licensed¹¹ and performs almost real-time decoding.

Julius is a two pass decoder. In the first pass, the decoding is performed using time synchronous beam search. The second pass re-ranks and further prunes the extracted hypothesis from the pass one. Bigram LM is used for the first pass and more complex trigram LM is used for re-ranking.

Before the implementation of this thesis was finished the Alex SDS team had been interested in Julius because its ability of real-time decoding and confusion network¹² output format.

The Alex team abandoned the Julius decoder for software issues e.g., crashes of the decoder. The crashes appeared during extracting confusion networks from Julius. In addition, the crashes were hard to detect because Julius used to run in a separate process.

⁹<http://htk.eng.cam.ac.uk/docs/license.shtml>

¹⁰You need to register even to see the license:
http://htk.eng.cam.ac.uk/prot-docs/hdecode_register.shtml

¹¹<http://www.lininfo.org/bsdlicense.html>

¹²A confusion network is approximation of a lattice described in Section 2.1.4.

2.4 Kaldi

Kaldi is a speech recognition toolkit consisting of a library, command line programs and scripts for acoustic modelling. Kaldi deploys several decoders for evaluation Kaldi AMs. Kaldi uses Viterbi training for estimating AMs. Only in special cases of speaker adaptive discriminative training the extended Baum-Welch algorithm is also used[27].

The architecture of the Kaldi toolkit could be separated to Kaldi library and training scripts. The scripts access the functionality of Kaldi library through command line programs. The C++ Kaldi library is based on the *OpenFST*[3] library and it uses optimized libraries for linear algebra such as BLAS and LAPACK. Related functionality is usually grouped in one namespace in C++ code, which corresponds to one directory on file system. The examples of the namespaces or directories can be seen in Figure 2.11

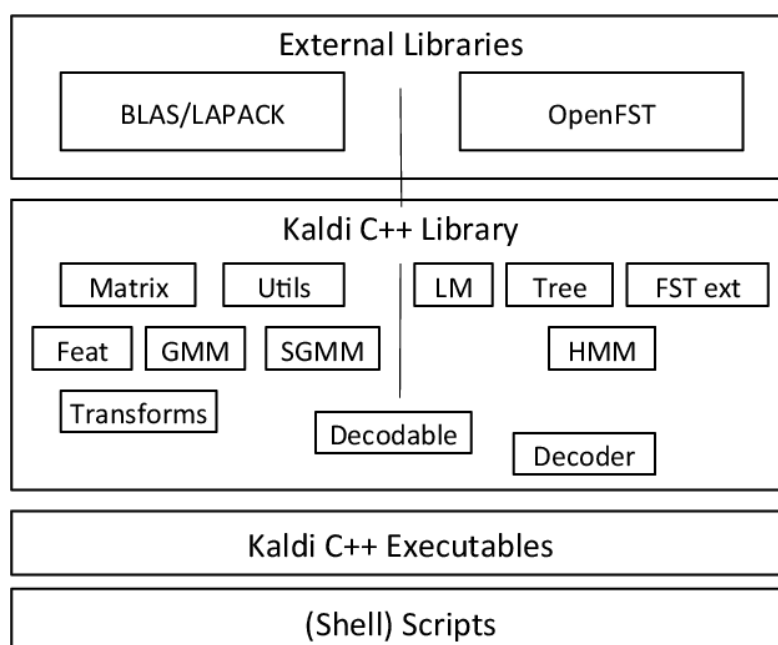


Figure 2.11: Kaldi toolkit architecture[27]

Kaldi uses executables which load its input from files and typically store results again to files. Alternatively, the output of one Kaldi program can be feed into next command using system pipes. There are usually many alternatives for every speech recognition tasks as seen in list of executables below:

1. Speech parametrisation
 - *apply-mfcc*
 - *compute-mfcc-feats*
 - *compute-plp-feats*
 - ...
2. Feature transformation
 - *apply-cmvn*

- *compute-cmvn-stats*
- *acc-lda*
- *fmpe-apply-transform*
- ...

3. Decoders

- *gmm-latgen-faster*
- *gmm-latgen-faster-parallel*
- *gmm-latgen-biglm-faster*
- ...

4. Evaluation and utilities

- *compute-wer*
- *show-alignments*
- ...

In addition, Kaldi provides very useful standardized scripts which wrap Kaldi executables or add new functionality. The scripts are located in *utils* and *steps* directories and are used in many training scripts recipes for different corpus data. In this thesis we created a new training recipe using the Kaldi infrastructure and Czech and English training corpus [17]. The recipe, the data and acoustic modelling scripts are described in Chapter 3.

2.4.1 Finite State Transducers

The Finite State Transducer framework and its implementation OpenFST determines the shape of the Kaldi data structures. Kaldi uses Finite State Transducer (FST) as underlying representation for LM, partially for AM, lexicon and also for representing transformation between text, pronunciation and triphones.

The FST framework provides well studied graph operations[19] which can be effectively used for acoustic modelling. Using the FST framework the speech decoding task is expressed as a beam search in a graph, which is well studied problem.

The FST graphs used for AM model training and speech decoding can be constructed as sequence of standardized OpenFST operations.[19]. Decoding is performed on so called *decoding graph HCLG* which is constructed from simple FST graphs as illustrated in Equation 2.19.

$$HCLG = H \circ C \circ L \circ G \quad (2.19)$$

. The symbol \circ represents an associative binary operation of composition on FSTs. We briefly explain the functionality of the transducers from Equation 2.19:

1. G is an acceptor that encodes the grammar or language model.
2. L represents the lexicon. Its input symbols are phones. Its output symbols are words.
3. C represents the relationship between context-dependent phones on input and phones on output.

4. H contains the HMM definitions, that take as input id number of Probability Density Functions (PDFs) and return context-dependent phones.

Following one liner illustrates how Kaldi decoding graph is created using standard FST operations¹³. [19]

$$HCLG = asl(min(rds(det(H'omin(det(Comin(det(LoG)))))))) \quad (2.20)$$

Semiring

Most of the operations operate on paths in the decoding graph. Path is a sequence of edges which have weights and an input and an output labels. Based on the weight type and weight path operations we distinguish several semirings.

Formally, a *semiring* $(\mathcal{K}, \oplus, \otimes, \bar{0}, \bar{1})$ is an algebraic structure on set \mathcal{K} with operations \oplus and \otimes . The binary operations multiplication \oplus and addition \otimes have identity element $\bar{0}$ respectively $\bar{1}$. The (\mathcal{K}, \oplus) forms commutative monoid and (\mathcal{K}, \otimes) forms just a monoid. The multiplication is left and right distributive over addition. Moreover, multiplication by $\bar{0}$ annihilates any member of \mathcal{K} to *zero*. Table 2.4.1 shows useful semirings in OpenFST.

Name	\mathcal{K}	\oplus	\otimes	$\bar{0}$	$\bar{1}$
Real	$[0, \infty)$	+	*	0	1
Log	$(-\infty, \infty)$	$-\log(e^{-x} + e^{-y})$	+	∞	0
Tropical	$(-\infty, \infty)$	min	+	∞	0

Table 2.1: Semirings used in speech recognition.[35]

¹³Kaldi tutorial on building *HCLG*: http://kaldi.sourceforge.net/graph_recipe_test.html

3. Acoustic model training

This chapter presents new Kaldi acoustic modelling scripts for free Czech and English "Vystadial" data. The scripts were developed as part of this thesis, they are licensed under the Apache 2.0 license and are publicly available in the Kaldi repository¹. The Acoustic Model (AM) trained using these scripts can be used for both batch speech recognition with common Kaldi decoders and for our *OnlineLatgenRecognizer*, which performs on-line decoding described in Chapter 4.

The first Section 3.1 describes the used data. The chapter continues by presenting the AMs training in Section 3.2. Later, in Section 3.3 we evaluate trained AMs and also compare them to generative HTK AMs which are trained using state of art HTK scripts.

3.1 Vystadial acoustic data

The data were collected in Vystadial project², and they are released under the Creative Commons Share-alike (CC-BY-SA 3.0) license. The Czech³ and English⁴ data are available online in the Lindat repository^{5,6}.

The English acoustic data consists of recorded phone calls among humans and the Spoken Dialog System, which was designed to provide the user with information on a suitable dining venue in the town. Most of the data was spoken in American English. The typical sentences recorded from users were queries for the dialogue system e.g.,

```
I NEED A CHINESE TAKE AWAY RESTAURANT IN THE CHEAP PRICE RANGE
I'M LOOKING FOR AN INTERNATIONAL RESTAURANT
I NEED TO FIND A PUB IT SHOULD ALLOW CHILDREN AND HAVE A TELEVISION
```

On the other hand, the Czech recordings were collected in three different ways[17]:

1. using a free Call Friend phone service
2. using the Repeat After Me speech data collecting process,
3. from the telephone interactions with the Alex SDS in a public transport domain.

In the Call Friend service native Czech speakers were invited to make free calls. In Repeat After Me process volunteers called a number where they were asked to repeat sentences synthesized by a Text to Speech (TTS).

The user language differs significantly in dialogues with Alex system and the other two settings. The sentences in Alex public transport domain, as seen in the first paragraph, are shorter and contain noises. The speech is spontaneous and proper names are frequently used. On the other hand, the other two recording

¹<http://sourceforge.net/p/kaldi/code/HEAD/tree/sandbox/oplatek2/egs/vystadial/>

²<http://ufal.mff.cuni.cz/grants/vystadial>

³Czech data: <http://hdl.handle.net/11858/00-097C-0000-0023-4670-6>

⁴English data: <http://hdl.handle.net/11858/00-097C-0000-0023-4671-4>

⁵<http://lindat.mff.cuni.cz/repository/>

⁶A previous version of our training scripts is published with the data in the Lindat repository and described in work [17].

tasks, as seen in the second paragraph, have much broader vocabulary with less named entities, and the ideas are expressed in longer sentences.

A DALŠÍ
NOISE
JO DĚKUJU MOC TO JSEM CHTĚL VĚDĚT
ZE ZASTÁVKY DEJVICKÁ

PRYČ S TYRANY A ZRÁDCI VŠEMI
UTRHNĚ SI KVĚT Z KYTICE A ODCHÁZÍ
DYŤ TO JE HORŠÍ NEŽ ZVÍŘE
O LIBERALIZMU TEHDY NEBYLO ŘEČI
CO BY TAM S TEBOU DĚLALI

The AMs for Czech are trained on acoustic data from all the three very different domains, because there is only two hours of in-domain data available in the Alex’s public transport domain. The evaluation for Czech data in Section 5.3 is performed on a Vystadial test set combined from all three domains. The English AMs are trained and tested on the data collected from the Venue domain using SDS. The summary of audio sizes in training, development and test set are presented in Table 3.1. Both Czech and English orthographic speech transcriptions were transcribed by humans.

dataset	audio[hour]	# sentences	# words
English			
training	41:30	47,463	178,110
development	01:45	2,000	7,376
test	01:46	2,000	7,772
Czech			
training	15:25	22,567	126,333
development	01:23	2,000	11,478
test	01:22	2,000	11,204

Table 3.1: Size of the data: length of the audio (hours:minutes), number of sentences (which is the same as the number of recordings), number of words in the transcriptions.[17]

3.2 Acoustic modelling scripts

We search for the best non-speaker adaptive AMs in our scripts for AM training. In this section, the explored methods and their settings are described, and the Section 5.3 presents the results for both Czech and English datasets. The Czech and English training scripts differ only in using a different phonetic dictionary, but otherwise the scripts remains exactly the same.

The AMs are trained via Viterbi training. The recordings and their transcriptions from training dataset are used for acoustic modelling. The estimated AMs are evaluated on the test set. The decoding of the test utterances is performed always with the same parameters, so that different AMs can be compared. The Figure 3.1 lists all acoustic models trained in our scripts. An advanced AM is

always initiated by audio alignments (respectively acoustic features alignments) using a simpler AM.

In paragraphs below, the organisation of acoustic model training is described. The used methods are listed in Figure 3.1 together with their hierarchy. The hierarchy shows that a more advanced method typically reuses initial values from previously trained simpler AM.

At first, a mono-phone model is trained from flat start using the MFCCs, Δ and $\Delta\Delta$ features. We force-align the feature vectors to HMM states using utterances' transcriptions. Secondly, we retrain the triphone AM (*tri1a*). One branch of experiments finishes by training MFCC $\Delta + \Delta\Delta$ triphone AM (*tri2a*).

On the other hand, the second branch instead of $\Delta + \Delta\Delta$ transformation uses LDA+MLLT to train AM (*tri2b*). Using the AM *tri2b* three AMs are discriminatively trained using the following objective functions:

1. Maximum Mutual Information[7]⁷. The model *tri2b_mmi* is trained in four loops.
2. Boosted Maximum Mutual Information[29]. The model *tri2b_bmmi* is trained in four loops with parameter 0.05.
3. Minimum Phone Error[26]. The model *tri2b_mpe* is also retrained in four loops.

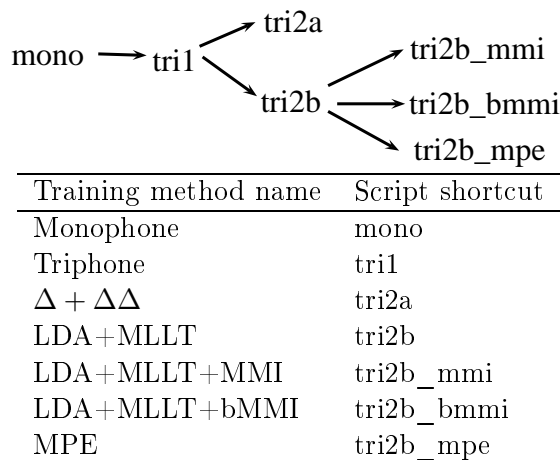


Figure 3.1: Training partial order among AM in our training scripts

The acoustic models *mono*, *tri1*, *tri2a* and *tri2b* are trained generatively. The models *tri2b_mmi*, *tri2b_bmmi* and *tri2b_mpe* are trained discriminatively in four iterations. The discriminative models yield better results than generative models if enough data is available. See Figure 3.2 for evidence.

The discriminative models from last iterations may over-fit to the training data, so that they perform worse on test dataset than models from previous iterations. Each iteration use a unigram LM estimated on training dataset in order to compute their objective function. As a result, each iteration adapts more to the training data, but may yield worse performance on the test dataset. We used only for iterations for discriminative models training and fortunately, we have not experienced such behaviour.

⁷Note the Maximum Mutual Information (MMI) function is implemented as bMMI with boosted parameter set to 0.

Setup for feature transformations

We explore not only AM training methods, but we also experiment with two feature transformation techniques. First, the $\Delta + \Delta\Delta$ triples the number of 13 MFCC features by computing also the first and the second derivatives from MFCC coefficients, resulting in 39 features per frame.

Second, the combination of LDA and MLLT is computed from 9 spliced frames consisting of 13 MFCC features. The default context window of 9 frames takes current frame, four frames from the left context and four frames from the right context. The LDA and MLLT feature transformation gains substantial improvement over $\Delta + \Delta\Delta$ transformation. See Figure 3.2.

Decoding setup

We use the trained AMs described above for decoding the utterances from the test dataset. For each trained AM we use the same speech parametrisation and feature transformation method as was used for the given AM at training time. We experiment with all trained AMs with both zerogram and bigram LM.

The default bigram and zerogram LMs for are built from orthographic transcriptions. The bigram LM is estimated from the training data transcriptions. Consequently, in a test set appear unknown words, so called Out of Vocabulary Word. The zerogram is extracted from a test set transcriptions. The zerogram is a list of words with probabilities uniformly distributed, so it helps decoding just by limiting the vocabulary size. The bigram LM contains 17433 unigrams and 79333 bigrams for Czech and 936 unigrams and 5521 bigrams for English. The zerogram LM is limited to 2944 words for Czech and to 302 words for English.

The speech recognition parameters are set to default values; the exceptions are decoding parameters: *beam=12.0*, *lattice-beam=6.0*, *max-active-states=14000* and Language Model Weight. The LMW parameter sets the weight of a LM, i.e., it regulates how much the LM is used to help AM in decoding. The LMW value is estimated on the development set and the best value is used for decoding on the test dataset. The details about *beam=12.0*, *lattice-beam=6.0* and *max-active-states=14000* can be found in Subsection 4.1.2. Section 5.3 evaluates the ASR performance for this parameters.

The *gmm-latgen-faster* decoder is used for the evaluation on testing data. It generates a word level lattice for each utterance and the one-best hypothesis is extracted from the decoded lattice and evaluated by WER and SER metrics against the reference transcription.

Note, that we are able to exactly reproduce the results of *gmm-latgen-faster* decoder with our *OnlineLatgenRecogniser*. The *gmm-latgen-faster* was used for evaluation in the scripts, so the Kaldi users do not have to install our extension.

3.3 Evaluation

The experiments focus on comparing the quality of ASR hypothesis measured by WER on AMs trained by different methods. We are not interested in absolute numbers since we model the language using a weak LM focusing on the acoustic modelling. By training only simple bigram LM we let the AM influence the recognition quality more significantly. The same motivation lead us to use zerogram LM which just limits vocabulary in the decoding task, and does not advise the

decoding search more probable phrases as higher order LM does. Consequently, the best words are chosen among all hypotheses only by acoustic similarity.

We concentrate on acoustic modelling since we believe that; if two AMs am_1 , am_2 are trained with the same weak Language Model lm_{weak} and the first AM gains lower WER than the second one ($wer_1^{weak} < wer_2^{weak}$), then in the same experiment with a richer LM lm_{rich} will still gain lower WER for the first AM ($wer_1^{rich} < wer_2^{rich}$).

First, we show how the data size influence the quality of AMs measured by WER. Second, the best results on full data is presented. Finally in Subsection 3.3.2, the best Kaldi results are compared against the results obtained by well-written HTK scripts by Keith Vertanen and improved by Matěj Korvas [17] on the same Vystadial dataset.

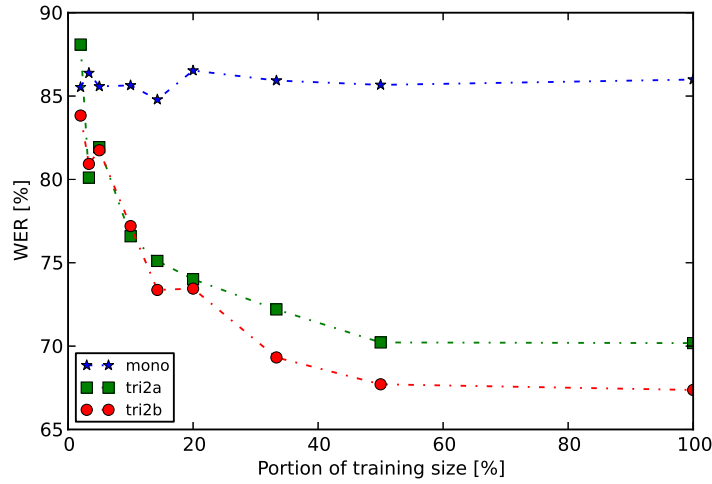


Figure 3.2: The figure displays improving performance of Czech generative AMs based on growing size of training data for acoustic modelling. The zerogram LM allows to evaluate only acoustic modelling, but causes a high WER.

The Figure 3.2 describes how the amount of acoustic data influences the WER. We illustrate that even with small datasets like Vystadial the high quality AM can be trained. The WER decreases significantly if new data are added to small dataset, but only small WER reduction is achieved when the last 50% of data is added. One can also see that the $\Delta + \Delta\Delta$ feature transformation is clearly outperformed on full data by LDA+MLLT setup. Note also that the monophone AM is typically used for the initialisation of triphone models and requires small portion of data to reach its limit. The WER is rather high due to the use of zerogram LM. We evaluate only generative LMs since we would have to have a fixed LM for discriminative methods and we do not have any obvious choice how to build one.

It may seem that more acoustic data is not needed for this domain, but discriminative training methods require more training data, and with more transcribed data a better LM adaptation can be achieved. The Figure 3.3 shows the effect of in-domain data size for LM on quality of speech decoding. The AM *tri2b_bmmi* and decoding parameters were fixed. The experiments were performed with different LMs which differ only in the training size used for their estimation. Note that

this experiment was run by Ondřej Dušek.⁸ The experiment was run on different test set from Public Transport Information (PTI) domain and the LMs were built also from that in-domain data.

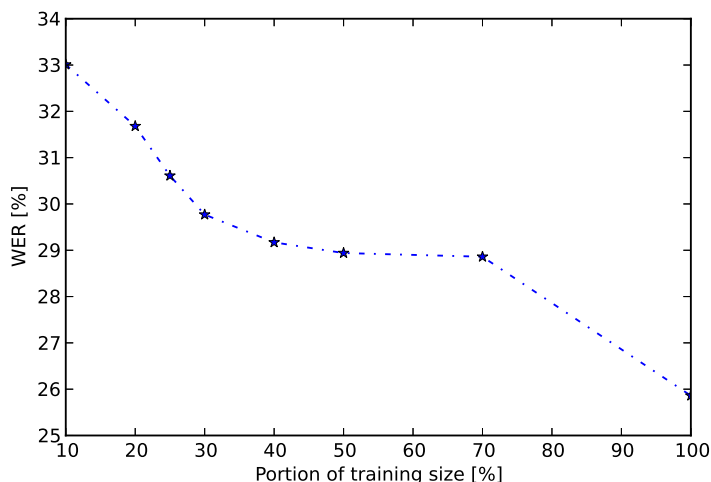


Figure 3.3: Influence of in-domain text size of LM on speech recognition quality. The AM *tri2b_bmmi* and parameters are fixed and only LM training size varies.

To conclude we are able to train reasonable AM with relatively small dataset such as Vystadial. On the other hand, additional data should improve speech recognition accuracy because

- the language domain changes in time and the new data reflect the differences,
- more data still may improve the best discriminatively trained AM,
- and last but not least the speech recogniser is more robust to new speakers.

3.3.1 Results

In this section we present the results of different acoustic training methods and we choose the best non-speaker adaptive setup. The Table 3.2 presents AMs results.

The complexity of the Czech data is clearly much larger than the complexity of English data. The high WER on the Czech dataset may be explained by following reasons:

- The mix of a very different domain and recording conditions is difficult to model by both AM and LM.
- The *Call Friend* and *Repeat After Me* collections task have a really broad domain which affect language modelling.
- The flexive languages such as Czech have larger vocabulary and higher Out of Vocabulary Words (OOVs) since one word may have several variations.

⁸Ondřej Dušek used our scripts developed for Alex dialogue system for the Public Transport Information (PTI) domain for the experiment.

language/method	zerogram	bigram
Czech		
tri $\Delta + \Delta\Delta$	70.7	56.6
tri LDA+MLLT	68.2	53.9
tri LDA+MLLT+MMI	65.3	49.5
tri LDA+MLLT+bMMI	65.3	49.3
tri LDA+MLLT+MPE	63.8	49.2
English		
tri $\Delta + \Delta\Delta$	35.7	16.2
tri LDA+MLLT	33.28	15.8
tri LDA+MLLT+MMI	25.01	10.4
tri LDA+MLLT+bMMI	23.9	10.2
tri LDA+MLLT+MPE	22.41	11.1

Table 3.2: Word error rates for zerogram and bigram LM for different training triphone methods. The ‘tri $\Delta + \Delta\Delta$ ’ row shows results for a generative model which is comparable to the model trained using the HTK scripts.

Nevertheless, the training scripts for the Czech data are very important since there is no other Czech acoustic data available, at least according our knowledge.

The WER on the Vystadial English data is lower than 20% for discriminative methods, which is reasonable, given the broad domain.

The discriminative training methods clearly outperformed the generative AMs, and also the LDA+MLLT is more effective feature transformation than using $\Delta + \Delta\Delta$ features. On the other hand, there are subtle differences among the three discriminatively trained AM in terms of performance. As a result, we choose AM (*tri2b_bmmi*) discriminatively trained by Boosted Maximum Mutual Information (bMMI) with MFCC, and LDA+MLLT preprocessing because informal experiments shows that decoding with Minimum Phone Error (MPE) Acoustic Model is slightly slower than with bMMI AM.

3.3.2 Kaldi and previous HTK results comparison

We compared Kaldi and HTK on the Vystadial Czech and English datasets, since we were not sure if the Kaldi toolkit is good enough alternative for HTK toolkit. In addition, by using state of the art HTK scripts we saw that the complexity of the Vystadial datasets is higher than in other datasets trained with the HTK scripts.⁹

We present results for triphone AM estimated using Baum-Welsh iterative training on zerogram and bigram LMs. The *HVite* HTK decoder was used to perform the decoding with the same LMs as used in Kaldi scripts. The training procedure is further described in work [17].

The results suggest that Kaldi achieves similar WER compared to HTK when using standard generative training methods and bigram LMs. Furthermore, one obtains a substantial decrease in WER by using more advanced discriminative training methods.

The experiment using MFCC, LDA & MLLT and bMMI discriminative training is a state of the art set up for speaker independent speech recognition[22] and

⁹Unfortunately, the dataset is not publicly available.

language/method	zerogram	bigram
Czech		
tri $\Delta + \Delta\Delta$	64.5	60.4
English		
tri $\Delta + \Delta\Delta$	50.0	17.5

Table 3.3: HTK results: Word error rates on test set are obtained by both a zerogram and a bigram LM. [17]. The AMs can be compared with the basic *tri $\Delta + \Delta\Delta$* Kaldi setup in Table 3.2.

outperforms HTK models.

Furthermore, in Chapter 5, we evaluate the trained Czech AMs on Public Transport Information domain on a different test set with a fine tuned LM and the best AM from list in Figure 3.1. The best AM is selected based on the results in Section 3.3.

4. Real time recogniser

This chapter presents the *OnlineLatgenRecogniser*, the new on-line Kaldi recogniser which can be used in real-time applications. Section 4.1 describes the implementation of the *OnlineLatgenRecogniser*. Next Section 4.2 introduces *PyOnlineLatgenRecogniser*, a Python extension of C++ *OnlineLatgenRecogniser*. Finally, Section 4.3 summarizes properties of the new implemented Kaldi recogniser.

We implemented a lightweight modification of the *LatticeFasterDecoder* from the Kaldi toolkit, improved on-line speech parametrisation and feature processing in order to create an *OnlineLatgenRecogniser*. The Kaldi *OnlineLatgenRecogniser* implements on-line interface which allows incremental speech processing, and it is able to process the incoming speech in small chunks incrementally. As a result, the real-time speech decoding can be performed while a user is speaking and the ASR output is obtained with a minimal latency.

The implementation of the recogniser was motivated by the lack of an on-line recognition support in Kaldi toolkit. Therefore, the toolkit decoders could not be used in applications such as spoken dialogue systems. Although Kaldi included an on-line recognition application; hard-wired timeout exceptions, audio source fixed to a sound card, and a specialised 1-best decoder limit its use only to demonstration of Kaldi recognition capabilities.

Our on-line recogniser uses acoustic models trained using the state-of-the-art techniques, such as Linear Discriminant Analysis (LDA), Maximum Likelihood Linear Transform (MLLT), Boosted Maximum Mutual Information (BMMI) and Minimum Phone Error (MPE). It produces word posterior lattices which can be easily converted into high quality n-best lists.

The recogniser's speed and latency can be effectively controlled off-line by optimising a language model. At runtime the speed of decoding is controlled by a beam threshold. The latency depends on the amount of time spent on word posterior lattice extraction from the recogniser, which can be regulated by a level of approximations used during the word lattice creation.

4.1 OnlineLatgenRecogniser

The standard Kaldi executables which implements speech parametrisation, feature transformations and decoder are using a batch file interface. Each executable loads the input from a file, processes a whole utterance and saves its output to another file. However, in real-time applications one would like to take advantage of the fact that an acoustic signal of an utterance is recorded in small chunks and can be processed incrementally.

We reimplemented speech parameterisation and feature transformations in order to fit on-line *OnlineLatgenRecogniser*'s interface, which can process audio features incrementally. In addition, we subclassed *LatticeFasterDecoder* and reorganized its original batch interface, so that it supports on-line decoding. Such implementation almost eliminates latency of a recogniser since almost all of the decoding can be performed while the user is still speaking.

First, we present the public on-line interface of *OnlineLatgenRecogniser* and in next subsections we introduce its components. The Subsection 4.1.2 describes the decoder, the core component. Subsection 4.1.3 introduces on-line speech parametrisation and feature transformations and the Subsection 4.1.4 discusses word posterior lattice extraction.

4.1.1 *OnlineLatgenRecogniser* interface

The *OnlineLatgenRecogniser* makes use of the incremental speech pre-processing and modified *LatticeFasterDecoder* in order to provide the following speech recognition interface:

- *AudioIn* – queueing new audio for pre-processing,
- *Decode* – decoding a fixed number of audio frames,
- *PruneFinal* – preparing internal data structures for lattice extraction,
- *GetLattice* – extracting a word posterior lattice and returning log likelihood of processed audio,
- *GetBestPath* – extracting a one best word sequence,
- *Reset* – preparing the recogniser for a new utterance,

The interface is influenced by the decoder interface and the preprocessing of the utterance is completely hidden for the user of *OnlineLatgenRecogniser*. The *AudioIn* is the only method which is not related to the decoder functionality.

The C++ example in Listing 4.1 shows a typical use of *OnlineLatgenRecogniser*. When audio data becomes available, it is queued into the recogniser’s buffer (line 11) and immediately decoded (lines 12-14). If the audio data is supplied in sufficiently small chunks, the decoding of queued data is finished before new data arrives. When the recognition is finished, the recogniser prepares for lattice extraction (line 16). Line 20 shows how to obtain word posterior lattice as an OpenFST object. The auxiliary *getAudio()* function represents a separate process supplying speech data. Please note that the recogniser’s latency is mainly determined by the time spent in the *GetLattice* function since the whole loop is processed while the user is speaking.

Listing 4.1: Example of the decoder usage

```
1 OnlineLatgenRecogniser rec;
2 rec.Setup(...);
3
4 size_t decoded_now = 0;
5 size_t max_decode = 10;
6 char *audio_array = NULL;
7
8 while (recognitionOn())
9 {
10     size_t audio_len = getAudio(audio_array);
11     rec.AudioIn(audio_array, audio_len);
12     do {
13         decoded_now = rec.Decode(max_decode);
14     } while(decoded_now > 0);
15 }
16 rec.PruneFinal();
17
18 double tot_lik;
19 fst::VectorFst<fst::LogArc> word_post_lat;
20 rec.GetLattice(&word_post_lat, &tot_lik);
21
22 rec.Reset();
```

We designed the interface with following criteria in mind:

- Passing the audio in the recogniser should accept any size of audio input.

- Decoding should return a number of actually decoded frames. The recogniser may decode less frames than requested if not enough audio is available.
- Decoding should be called frequently on small chunks, which guaranties quick response times of the *Decode* method.

Consequently, *OnlineLatgenRecogniser* does not block a process to either load audio or decode an utterance. The loading of audio and the decoding can be easily alternated back and forth. Obviously, the decoding of single utterance can be separated into number of parts and other tasks can be run in a single process with speech recognition in order to allow an application to stay responsive. We are able to decode the utterance while the user speaks.

On the other hand, extracting the word posterior lattice may block the process since it is very computationally demanding. It lasts several tens of milliseconds. However, it is called only at the end of each utterance. Extracting one best word sequence is much faster and can be called at any time.

4.1.2 *OnlLatticeFasterDecoder*

We did not implement any new functionality *OnlLatticeFasterDecoder*, but we only reorganised the code of base class *LatticeFasterDecoder*. We split the *LatticeFasterDecoder::Decode* method which performed several tasks into more methods. The *LatticeFasterDecoder::Decode* function runs a beam search from frame 0 to the end of each utterance. In addition, a pruning is triggered periodically in the function. In *OnlineLatgenRecogniser*, we control the beam search by the following functions:

- *Decode* – decoding a fixed number of audio frames instead of decoding whole utterance, pruning is triggered periodically,
- *PruneFinal* – run final pruning and so prepare the internal data structures for lattice extraction,
- *Reset* – preparing the recogniser for a new utterance.

In the *PruneFinal* function, which is called at the end of an utterance, the states are pruned by beam search with the knowledge that no further search will be performed, so more states can be safely discarded.

The decoding is performed on request by calling the *Decode* method with a parameter (int *max_frames*) which limits the number of decoded frames. It returns the number of frames which were actually decoded, which is always smaller or equal to *max_frames* value. The *OnlLatticeFasterDecoder::Decode* method performs decoding frame by frame using the Viterbi beam search. The speed of the Viterbi search is highly predictable for fixed settings of the recogniser. As a result, the *max_frames* parameter effectively limits the amount of time in the *Decode* method. Repeated calls of *Decode* with small values of *max_frames* keep the recognition responsive as implemented in Listing 4.1.

The ASR output is extracted by the original methods of *LatticeFasterDecoder*:

- *GetRawLattice* returns state-level lattice,
- *GetLattice* extracts from state-level lattice word lattice which is returned,
- *GetBestPath* returns just one-best path hypothesis.

The state-level lattice, which is returned from the *GetRawLattice* method, can be understood as lattice on triphone level. In the state-level lattice, a single word hypothesis is typically can be obtained from multiple state-level hypotheses due to different word alignments, i.e., the same words sequences were pronounced with different timing.

The decoding of *LatticeFasterDecoder* as well as lattice extraction can be controlled by several parameters. We mention the most important parameters which affect both speed and the ASR output quality. The parameters either increase speed and decrease ASR output quality or vice versa.

The *beam* and *max-active-states* parameters directly affect the speed of decoding. The *beam* parameter affects the speed of all utterances, whereas the *max-active-states* parameter plays its role for utterances with uncertainty in beam search due to noises. In fact, the *max-active-states* is a threshold for worst case scenarios. The *lattice-beam* influences speed of lattice extraction.

The properties of the parameters and its relationship to ASR output quality is described in detail in Section 5.3 where we evaluate the recogniser.

4.1.3 On-line feature pre-processing

This section describes audio signal buffering, MFCC feature extraction and feature transformation. The resulting acoustic features are then used with an AM in *OnlLatticeFasterDecoder* to obtain likelihood of each state explored by Viterbi search. *OnlineLatgenRecogniser* only uses the likelihood to run Viterbi search. The likelihood itself is extracted from AM based on the acoustic features by *DecodableInterface*.

When a decoder is asked to perform decoding it needs to estimate likelihood for the states which should be explored, and so it requests the *DecodableInterface*. We implemented on-line version of *DecodableInterface* which let the decoder ask for likelihoods of new acoustic features frame by frame. The decoder, the pre-processing pipeline and the data flow between the components are illustrated in Figure 4.1. We briefly describe one step of Viterbi search:

- Audio is extracted from a buffer.
- The MFCC features are computed on overlapping audio window. The new audio is used for shifting the audio window.
- Applying feature transformation on top of MFCC features.
 - $\Delta + \Delta\Delta$ requires at least two previous frames, if available the acoustic features a are returned.
 - The $LDA + MLLT$ is computed using context, which by default is set to four previous and four future frames. If context is available, the acoustic features a are returned.

Note, that the $LDA + MLLT$ and the $\Delta + \Delta\Delta$ transformations are complementary.

- The *OnlDecodableDiagGmmScaled* queries the AM for the likelihood of acoustic features and given state.
- The decoder itself performs the search in state level space having the probabilities from the *Decodable* interface.

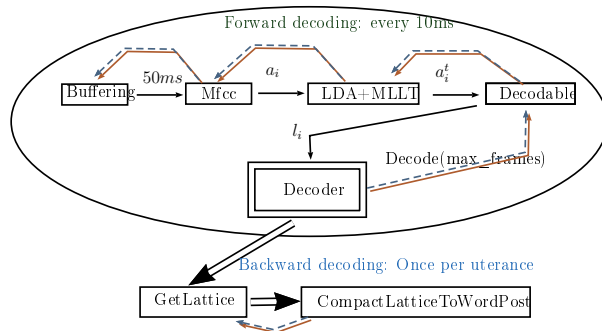


Figure 4.1: Components for on-line decoding

Each step in Figure 4.1 is implemented as a separate C++ class. *OnlineLatticeRecogniser* instantiate each class during the setup.

The on-line implementation *OnlDecodableDiagGmmScaled* of *DecodableInterface* easily handles missing audio data. If *OnlDecodableDiagGmmScaled* has no acoustic features precomputed, it asks the feature transformation to compute new features. If the feature transformation can not provide the features, so it returns default empty value. The empty value signals to *OnlDecodableDiagGmmScaled* that no features are available at the moment. Finally, *OnlDecodableDiagGmmScaled* produces also default empty value and the decoder's method *Decode(int max_frames)* returns zero indicating that no frames were decoded.

We have not experimented speech parametrisation settings. We used the recommended values, which are tested in tens of Kaldi recipes. The list the most important parameters:

- The frame width (set to 25 ms),
- the frame shift (set to 10 ms),
- and the frame splicing used for LDA+MLLT (nine frames are spliced).

4.1.4 Post-processing the lattice

The *OnlineLatgenRecogniser* not only extracts word lattice using *OnlLatticeFasterDecoder::GetLattice* function, but also computes posterior probabilities for the word lattice. The *OnlLatticeFasterDecoder* returns word lattice with alignments in form of *CompactLattice*. The *CompactLattice* determinised at state level still may contain multiple paths for each word sequence encoded in the lattice. The *CompactLattice* distinguishes each path not only according to the word labels on the path, but also according to the alignments. In order to obtain only the word lattice, we discard the alignments.

The steps of converting *CompactLattice* to word posterior lattice are listed below. For the implementation details see Listing 4.2:

- Joining multiple word sequences which differer in word alignments is performed in two steps:
 - Discarding the alignments from *CompactLattice*.
 - Converting the lattice to its minimal lattice representation with no alternatives for one word hypothesis.

- The computing of the posterior probabilities through a standard forward-backward algorithm, which is implemented in two steps:
 - Computing α and β data structures for which a Kaldi implementation is reused.
 - Updating the lattice weights from likelihood to posterior probabilities based on α and β , which we implemented in the *MovePostToArcs* function.

Listing 4.2: Converting *CompactLattice* to posterior word lattice

```

1 double CompactLatticeToWordsPost(CompactLattice &clat,
2                                 fst::VectorFst<fst::LogArc> *pst) {
3     {
4         Lattice lat;
5         fst::VectorFst<fst::StdArc> t_std;
6         RemoveAlignmentsFromCompactLattice(&clat); // remove the alignments
7         ConvertLattice(clat, &lat); // convert to non-compact form.. no new
            ↪ states
8         ConvertLattice(lat, &t_std); // this adds up the (lm,acoustic) costs
9         fst::Cast(t_std, pst); // reinterpret the inner implementations
10    }
11    fst::Project(pst, fst::PROJECT_OUTPUT);
12    fst::Minimize(pst);
13    fst::ArcMap(pst, fst::SuperFinalMapper<fst::LogArc>());
14    fst::TopSort(pst);
15    std::vector<double> alpha, beta;
16    double tot_lik = ComputeLatticeAlphasAndBetas(*pst, &alpha, &beta);
17    MovePostToArcs(pst, alpha, beta);
18    return tot_lik;
19 }

```

The word posterior probability is converted from the likelihood of the words in word lattice. The word lattice obviously contains alternatives which were explored by the beam search during decoding the utterance. Consequently, the posterior probability is an approximation because the very low probable alternatives discarded by beam search are not considered. On the other hand, the discarded alternatives are so improbable so they almost do not influence the posterior probability.

Presumably, the word posterior values are more impacted by inaccurate likelihood values taken from the Acoustic Model. Generative models are improved so the likelihood match the reality as much as possible. On the other hand, the discriminative AM models deliberately favour the most probable hypothesis by boosting the likelihood of the most probable hypothesis. As a result, the word posterior probability for the best hypothesis is artificially boosted. At the moment, we do not calibrate the word posterior probabilities in extracted lattices.

4.2 PyOnlineLatgenRecogniser

We also developed a Python extension, *PyOnlineLatgenRecogniser*, exporting the *OnlineLatgenRecogniser* C++ interface to Python. It can be used as an example of bringing Kaldi's on-line speech recognition functionality to higher-level programming languages. We extended also PyFST library[6], which interfaces OpenFST C++ template library into Python because we need to process further the OpenFST lattices produced by *PyOnlineLatgenRecogniser* in Python. Consequently, the

recogniser as well as its input and output can be seamlessly used both from C++ and Python.

PyOnlineLatgenRecogniser is a thin wrapper around *OnlineLatgenRecogniser* implemented using Cython[4]. The Cython compiler is well known for generating fast code when interfacing Python and C++ and the wrapper causes no measurable overhead.

We implemented conversion of the word posterior lattices to an n-best list. The implementation is efficient since the OpenFST shortest path algorithm is used on small lattices.

The minimalistic Python example in Listing 4.3 shows usage of the *PyOnlineLatgenRecogniser* and the decoding of a single utterance.

The audio is passed to the recogniser in small chunks (line 4), so the decoding (line 5 and 8) can be performed while the user is speaking. When no more audio data is available a likelihood and a word posterior lattice is extracted from the recogniser (line 10).

Listing 4.3: Fully functional example of the *PyOnlineLatgenRecogniser* interface

```

1 | d = PyGmmLatgenWrapper()
2 | d.setup(argv)
3 | while audio_to_process():
4 |     d.frame_in(get_raw_pcm_audio())
5 |     dec_t = d.decode(max_frames=10)
6 |     while dec_t > 0:
7 |         decoded_frames += dec_t
8 |         dec_t = d.decode(max_frames=10)
9 | d.prune_final()
10 | lik, lat = d.get_lattice()
```

Note that *PyOnlineLatgenRecogniser* and *OnlineLatgenRecogniser* are initialised by string vector of arguments in command line format. The parameters are parsed using Kaldi’s command line parser and options affect behaviour speech parametrisation, feature transformations and the *OnlLatticeFasterDecoder*. In addition, exactly the same parameters can be parsed by standard Kaldi utilities. We created demos¹ which use the same parameters for speech recognition using:

- standard Kaldi executables and scripts
- *PyOnlineLatgenRecogniser*
- *OnlineLatgenRecogniser*

The alternatives produce exactly the same results.

4.3 Summary

The *OnlLatticeFasterDecoder* performs the on-line speech recognition. We suggest exploiting the *OnlineLatgenRecogniser* and decoding utterances in small chunks and pass the audio to the recogniser immediately as it is available. The speech recognition parameters are initialized with reasonable default values and the parameters are the same as used in Kaldi executables. As a result, one can use the parameters from any Kaldi recipe to obtain the exactly same high quality results in the on-line speech recognition setting.

¹https://github.com/UFAL-DSG/pykaldi/tree/master/egs/vystadial/online_demo

The implemented minimal on-line interface which supports MFCC speech parametrisation, $\Delta - \Delta\Delta$ feature transformation or LDA+MLLT and both generative training and discriminative training using bMMI and MPE. The MFCC, LDA+MLLT and bMMI is one of the best setup for the speaker independent speech recognition. To conclude, we reimplemented Kaldi batch speech recognition, so that it can perform on-line real-time speech recognition and still maintain its high quality. The next Chapter 5 evaluates in detail the recognisers' real-time performance in the Alex Dialogue Systems Framework.

5. Kaldi ASR in Alex SDS

This chapter discuss the details of deploying *OnlineLatgenRecogniser* into Alex dialogue system. The *OnlineLatgenRecogniser* is used in Alex dialogue system for Czech Public Transport Information (PTI) domain available on public toll-free (+420) 800 899 998 line.

First, the architecture of Alex Spoken Dialog System (SDS) is described. Second, Section 5.2 presents how the wrapper *PyOnlineLatgenRecogniser* is integrated into SDS Alex. Finally, Section 5.3 evaluates the decoder in Alex dialogue system on Czech PTI domain.

5.1 Alex dialogue system architecture

The Alex dialogue system has a speech to speech user interface. The Alex dialogue system is developed in Python programming language and consists of six major components.

1. Voice Activity Detection (VAD)
2. Automatic Speech Recognition (ASR)
3. Spoken Language Understanding (SLU)
4. Dialogue Manager (DM)
5. Natural Language Generation (NLG)
6. Text to Speech (TTS)

The system interacts with the user in *turns*. The schema in Figure 5.1 illustrates how the user's input is processed in single turn. The spoken input is passed to ASR component which generates corresponding textual representation. SLU extracts semantic features from the text and DM decides which response to present. The NLG component generates textual response from an internal representation of DM and finally the TTS read the text with human voice.

Each of the Alex's component runs in separate process in order parallelize the input data processing and output data generation. The components communicates among themselves through system pipes.

In order to prepare ASR unit for *PyOnlineLatgenRecogniser* we have implemented not only the wrapper itself, but also preparation scripts and useful utilities. Let us introduce the Alex SDS framework organisation, so we can better explain how our scripts are used. The framework is separated into several logical parts:

- The core library is located at *alex/components/*. The library is domain and language independent. All components in Figure 5.1 are implemented in this core library.
- Settings and scripts for specific domain applications are located in *alex/applications/*. For example, application for PTI domain can be found in *alex/applications/PublicTransportInfoCS/* directory.
- The scripts which use external tools or data can be found in:
 - *alex/corpus_tools/* directory which focuses on formatting and organising the collected data,

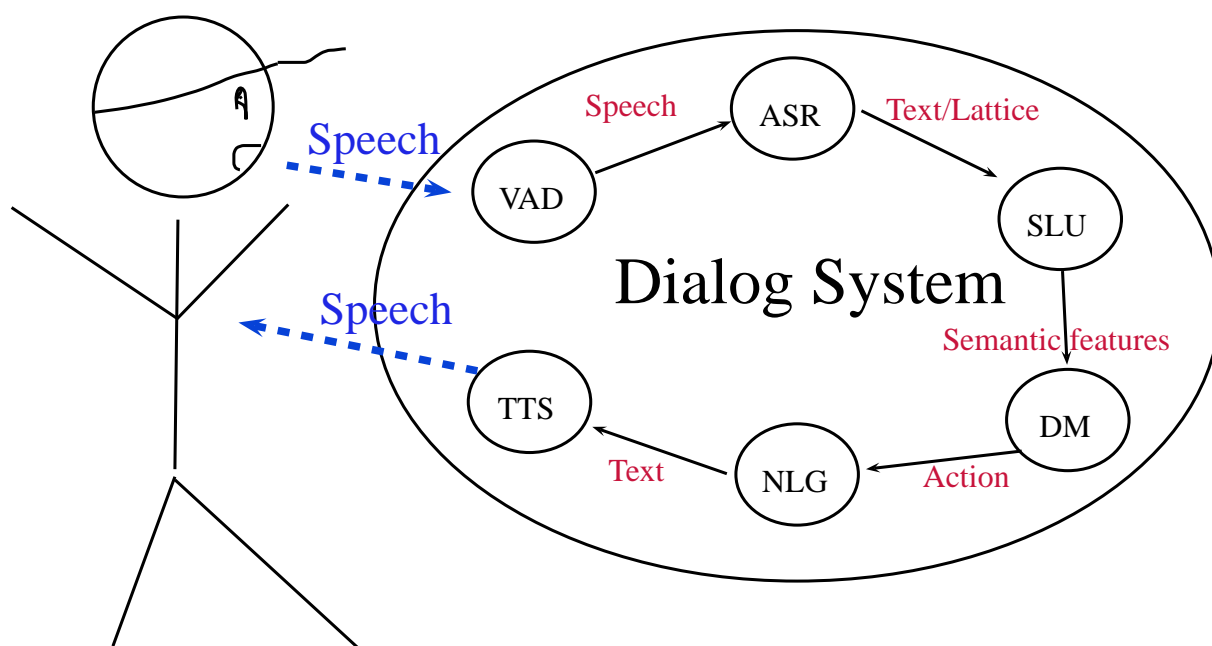


Figure 5.1: Single turn in Alex dialogue system

– and `alex/tools/` directory which stores code for modelling VAD, ASR, *SIP* client, etc.

- Integration tests are stored in `alex/tests/`.
- The `alex/utils/` directory contains simple utilities for various purposes.

The components depicted in Figure 5.1 are represented as Python modules under the `alex/components/` directory. The source code of the components is very modular, so each component may support multiple implementations. For example the ASR component currently supports several ASR recognisers. The recognisers implement a common base class *ASRInterface* which is presented in Listing 5.1. The supported speech recognisers are:

- OpenJulius (`alex/components/asr/julius.py`) interfaces OpenJulius decoder through sockets for on-line recognition.
- Google (`alex/components/asr/google.py`) uses cloud service for batch decoding.
- Kaldi (`alex/components/asr/kaldi.py`) imports *PyOnlineLatgenRecogniser* class and uses its functionality for on-line decoding.

One can easily choose an ASR recogniser in Alex configuration file. The configuration file is also the right place to specify AM and LM and the speech recognition parameters if necessary.

In order to prepare a specific application with *PyOnlineLatgenRecogniser* one need to train AM and LM. One need to always train SLU unit based on the ASR unit outputs. The LM model training and consequently SLU training is very domain specific so the scripts are deployed for each application separately. For example, the scripts for LM and SLU model training for PTI domain are located under directory `alex/applications/PublicTransportInfoCS/` in directories `lm/` and `slu/`.

5.2 Kaldi integration into Alex’s Spoken Dialog System framework

Integration of the Kaldi real-time recognizer into Alex’s framework requires implementing following features:

1. The *kaldi.py* module which exploits functionality of *PyOnlineLatgenRecogniser* and implements the abstract *ASRInterface*.
2. The training scripts for AMs.
3. The scripts for building custom decoding graph *HCLG.HCLG* graph is a Kaldi effective representation of AM and LM used for decoding.
4. Evaluation of the ASR recogniser in Alex, so the best speech recogniser can be selected.

The *PyOnlineLatgenRecogniser* integration is described in Subsection 5.2.1. The training scripts for training AMs were described in Chapter 3. However, note that we adjusted their directory structure and copied them into *alex/tools/kaldi* directory so they nicely in the Alex SDS. The scripts for building the *HCLG* decoding graph are introduced in Subsection 5.2.2. They are stored at *alex/applications/PublicTransportInfoCS/hclg/*. Finally, the Section 5.3 evaluates the performance of *PyOnlineLatgenRecogniser* in Alex SDS and briefly compares it with Google speech recognition service used through Python module *alex/components/asr/google.py*.

5.2.1 *PyOnlineLatgenRecogniser* in Alex

The ASR component in the Alex dialogue system runs as separate process, and the speech recognition is triggered based on VAD decisions.

If VAD detects start of speech in the input audio stream, it sends the speech signal to ASR component and the *rec_in* method is called. The *rec_in* method is a part of Alex abstract *ASRInterface* illustrated in Listing 5.1. See Listing 5.1. In Kaldi implementation of *rec_in*, the audio is decoded using beam search while the user is speaking, i.e., the method *rec_in* gradually adds the new audio to *PyOnlineLatgenRecogniser*’s buffer and immediately decodes it.¹

If VAD recognises end of speech, no more data are sent to *PyOnlineLatgenRecogniser* engine and *hyp_out* method is called in order to extracted word posterior lattice. Then, the word posterior lattice is converted to an n-best list.²

The *flush* method is used only if the speech recogniser wants to throw away the buffered audio input and reset the decoding.

The method *rec_wav* from Alex’s *ASRInterface* nicely illustrates how the two methods *rec_in* and *hyp_out* are used for decoding. Since the method is used only for testing purposes, it sends all input audio to the speech recogniser at once. However, in real-time application the audio is passed to *PyOnlineLatgenRecogniser* in small chunks, so the decoding can run as a user speaks.

In the on-line Kaldi settings, latency of the ASR unit depends mostly on the time spent in *hyp_out* method. In the *hyp_out* method a word posterior lattice

¹If the ASR component is busy with decoding the audio just waits in VAD buffer instead of in *PyOnlineLatgenRecogniser*’s buffer.

²We would like to implement direct keyword spotting from *pyfst* lattices in Alex SLU unit in future.

Listing 5.1: ASRInterface

```

1 class ASRInterface(object):
2
3     def rec_in(self, frame):
4
5     def flush(self):
6
7     def hyp_out(self):
8
9     def rec_wav(self, pcm):
10         self.rec_in(pcm)
11         return self.hyp_out()
```

is extracted using the *PyOnlineLatgenRecogniser::GetLattice* method as described in Subsection 4.2. For most cases the latency is well below 200 ms for our settings as illustrated in Figure 5.3.

The Alex dialogue system frequently handles several spoken requests immediately one after another. It may seem that a user should wait to response of the Alex dialogue system before he speaks again and consequently new audio is buffered to recognition. However in practice, users speak spontaneously; as first utterance is pronounced and users immediately start speaking in order to change his request. Nevertheless, at the end of each utterance the *hyp_out* method is called and the ASR hypothesis is extracted. Since the user already speaks when the lattice is extracted, the processor time which is meant for lattice extraction can not be used for decoding. Consequently, the *rec_in* should decode the buffered audio faster in shorter time than the user speaks.

We noticed the problem for chains of noises detected in VAD components as multiple short utterances. The (*rec_in* method) was extracting hypothesis while the audio input was still buffering the noisy speech and want to call *rec_in*. The (*hyp_out* method) was called so often that almost no decoding was performed. The problem was solved by improving VAD so the *hyp_out* method is not triggered so often.

5.2.2 Building in-domain decoding graph

A decoding graph is a graph represented as an OpenFst object it stores all the LM model information and part of information for acoustic modelling. The decoding graph is necessary for decoding with Kaldi decoders. We build the *HCLG* graph using standard OpenFst operations which are implemented in Kaldi utilities.

We designed our scripts so they automatically update newly built AMs and LMs and create all files necessary for decoding with *OnlineLatgenRecogniser* including *HCLG* graph. The same files can also be used with standard Kaldi decoders or *PyOnlineLatgenRecogniser*.

The *HCLG* build script requires:

- Language Model
- Acoustic Model
- Acoustic phonetic decision tree
- Phonetic dictionary

As a results, the *HCLG* graph is built. In addition, the script copies files generated by AM training which are required for decoding with *PyOnlineLatgenRecogniser*. To sum up, following files are necessary for decoding with Kaldi decoders:

- Decoding graph *HCLG*
- Acoustic Model
- Matrix which defines feature transformations
- Settings for speech parametrisation and feature transformations which match settings used for training the used AM.
- Word Symbol Table (WST) — mapping between integer labels

We also developed evaluation scripts which simply compute the statistics of measures which are evaluated given AM, LM and parameters in Section 5.3. Both the evaluation scripts and build *HCLG* script are located in the *alex/application-s/PublicTransportInfoCS/hclg/* directory.

Acoustic and language models for PTI domain

The *OnlineLatgenRecogniser* is evaluated on a corpus of audio data from the Public Transport Information (PTI) domain. In PTI, users can interact in Czech language with a telephone-based dialogue system to find public transport connections [39]. The PTI corpus consist of approximately 12,000 user utterances with a length varying between 0.4 s and 18 s with median around 3 s. The data were divided into training, development, and test data where the corresponding data sizes are 9496, 1188, 1188 utterances respectively. For evaluation, a domain specific class-based language model with a vocabulary size of approximately 52,000 and 559,000 n-grams was estimated from the training data. Named entities e.g., cities or bus stops, in class-based language model are expanded before building a decoding graph. The perplexity of the resulting language model evaluated on the development data is about 48.

Since the PTI acoustic data amounts to less then 5 hours, the acoustic training data was extended by additional 15 hours of telephone out-of-domain data from VYSTADIAL 2013 - Czech corpus [17]. The acoustic models were obtained by BMMI discriminative training with LDA and MLLT feature transformations. A detailed description of the training procedure is given in Chapter 3.

5.3 Evaluation of *PyOnlineLatgenRecogniser* in Alex

We focus on evaluating the speed of the *OnlineLatgenRecogniser* and its relationship with the accuracy of the decoder. We evaluate following measures:

- Real Time Factor (RTF) of decoding – the ratio of the recognition time to the duration of the audio input,
- Latency – the delay between utterance end and the availability of the recognition results,
- Word Error Rate (WER).

Accuracy and speed of the *OnlineLatgenRecogniser* are controlled by the *max-active-states*, *beam*, and *lattice-beam* parameters [27]. *Max-active-states* limits the maximum number of active tokens during decoding. *Beam* is used during graph search to prune ASR hypotheses at the state level. *Lattice-beam* is used

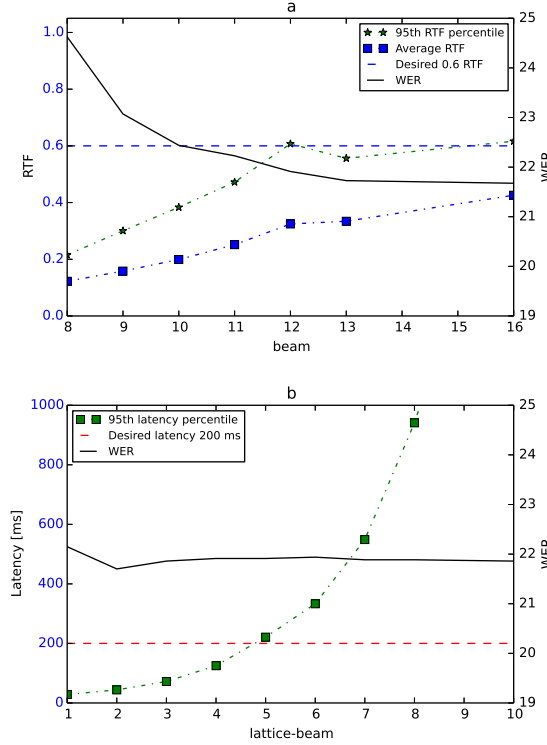


Figure 5.2: The upper graph (a) shows that WER decreases with increasing *beam* and the average RTF linearly grows with the beam. The growth of the 95th RTF percentile is limited at 0.6 by setting *max-active-states* to 2000, because the *max-active-states* parameters influence presumably the worst cases with large search space. The lower graph (b) shows latency growth in response to increasing *lattice-beam*.

when producing word level lattices after the decoding is finished. It is crucial to tune these parameters to obtain good results.

In general, one aims for a RTF smaller than 1.0. Moreover, it is useful in practice if the RTF is even smaller because other processes running on the machine can influence the amount of available computational resources. Therefore, we target the RTF with value of 0.6, which was estimated as sufficient by informal experiments.

We used grid search on the test set to identify the optimal parameters values. Figure 5.4 (a) shows the impact of the *beam* on the WER and RTF measures. In this case, we set *max-active-states* to 2000 in order to limit the worst case RTF to 0.6. Observing Figure 5.4 (a), we chose *beam* of value 13 for further experiments as this setting balances the WER. Figure 5.4 (b) shows the impact of the *lattice-beam* on WER and latency when *beam* is fixed to 13. We set *lattice-beam* to 5 based on Figure 5.4 (b) to obtain the 95th latency percentile of 200 ms, which is considered natural in a dialogue [38]. *Lattice-beam* does not affect WER, but larger *lattice-beam* improves the oracle WER of generated lattices [28]. Richer lattices may improve SLU performance.

Figure 5.3 shows the percentile graphs of the RTF and latency measures over the test set. The 95th percentile is the value of a measure such that 95% of the data has the measure below that value. One can see from Figure 5.3 that 95% of test

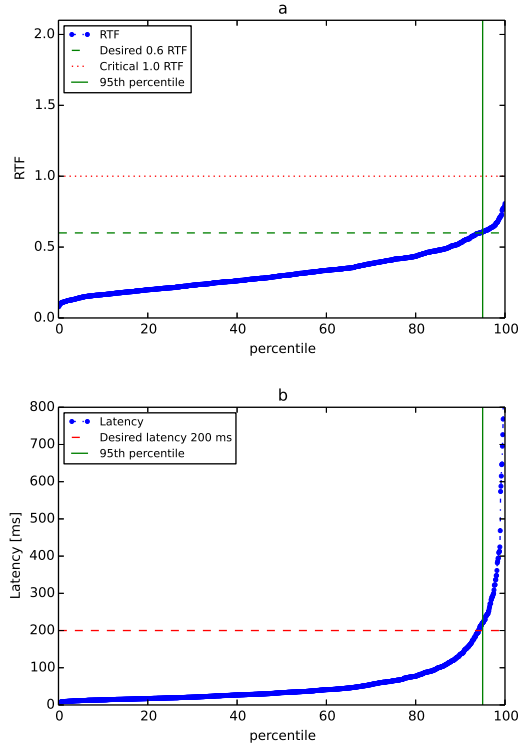


Figure 5.3: The percentile graphs show RTF and Latency scores for test data for $max-active-states=2000$, $beam=13$, $lattice-beam=5$. Note that 95 % of utterances were decoded with the latency lower than 200ms.

utterances is decoded with RTF under 0.6 and latency under 200 ms. The extreme values for 5% of test utterances are in most cases caused by decoding long noisy utterances where uncertainty in decoding increase the search space slows down the recogniser. Using $beam$ of 13, the $lattice-beam$ of 5 and 2000 $max-active-states$, the *OnlineLatgenRecogniser* decodes the test utterances with a WER of about 21%.

In addition, we have also evaluated Google ASR service as we used it previously in Alex SDS. The Google ASR service decoded the test utterances from the PTI domain with 95% latency percentile of 1900ms and it reached WER about 48%. The high latency is presumably caused by the batch processing of audio data and network latency, and the high WER is likely caused by a mismatch between Google’s acoustic and language models and the test data.

Results

To conclude, we implemented ASR component based on *OnlineLatgenRecogniser*. We also implemented scripts which allow easy AM training and testing, and LM evaluation for Kaldi speech recognition in Alex SDS.

Based on evaluation, we selected the best setup³ for ASR component in Alex Dialogue System Framework with WER under 22 %, latency less than 200 ms and RTF under 0.6 on PTI domain. As a result, the *OnlineLatgenRecogniser* performs significantly better than the previous ASR engines.

³Setup: $beam$ 12, $lattice-beam$ 5, $max-active-states$ 2000.

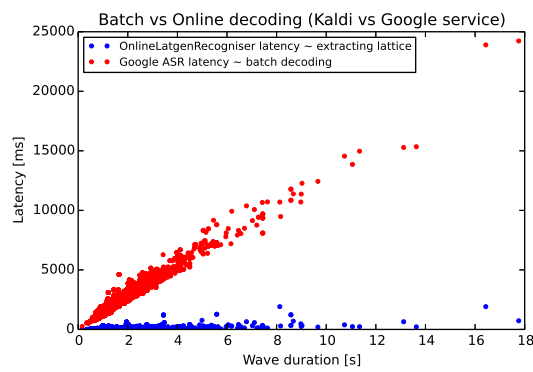


Figure 5.4: Almost constant latency of on-line decoder (OnlineLatgenRecogniser) and linearly growing latency of cloud based speech recogniser (Google ASR service) for increasing utterance length.

6. Conclusion

The Kaldi toolkit is a speech recognition toolkit distributed under a free license [27]. The toolkit is based on Finite State Transducers, implements state-of-the-art acoustic modelling techniques, is computationally efficient, and is already widely adapted among research groups. Its only major drawback was the lack of on-line recognition support. Therefore, it could not be used directly in applications such as spoken dialogue systems.

This work presented the *OnlineLatgenRecogniser*, an extension of the Kaldi automatic speech recognition toolkit. The *OnlineLatgenRecogniser* is distributed under the Apache 2.0 license, and therefore it is freely available for both research and commercial applications. The recogniser and its Python extension is stable and intensively used in a publicly available spoken dialogue system [39]. Thanks to the use of a standard Kaldi lattice decoder, the recogniser produces high quality word posterior lattices.

The training scripts as well as the source code of the *OnlineLatgenRecogniser* are currently merged into Kaldi repository¹². The Alex dialogue system and the integration of *OnlineLatgenRecogniser* is Apache, 2.0 licensed and freely available on Github³. The training scripts, the *OnlineLatgenRecogniser* and its Python wrapper *PyOnlineLatgenRecogniser* were developed also under Apache, 2.0 license on Github⁴.

The work specified by goals in introduction varies in many aspects. We successfully trained acoustic models, design and also implement decoders interface, improve real-time decoder and prepared experiments for evaluating results.

In addition to our implementation effort, we have also co-authored an article which uses AM training scripts described in Chapter 3. The article[17] describes the Czech and English Vystadial data sets as well as its acoustic modelling scripts in Kaldi and HTK. We also submitted an article about *OnlineLatgenRecogniser's* implementation and properties to the Sigdial conference⁵. The article is currently in a review process.

Future plans include implementing more sophisticated speech parameterisation interface and feature transformations, implementing normalisation of word posterior lattices and exploring acoustic modelling based on Deep Neural Networks.

Acknowledgments

This research was partly funded by the MEYS of the Czech Republic under the grant agreement LK11221 and core research funding of Charles University in Prague. The work described herein uses language resources hosted by the LINDAT/CLARIN repository, funded by the project LM2010013 of the MEYS of the Czech Republic. We would also like to thank Daniel Povey, Vassil Panayotov, Pavel Mencl, Ondřej Dušek, Matěj Korvas, Lukáš Žilka, David Marek and Tomáš Martinec for their useful comments and discussions.

¹<http://sourceforge.net/p/kaldi/code/HEAD/tree/sandbox/oplatek2/src/dec-wrap/>

²<http://sourceforge.net/p/kaldi/code/HEAD/tree/sandbox/oplatek2/egs/vystadial/>

³<https://github.com/UFAL-DSG/alex>

⁴<https://github.com/UFAL-DSG/pykaldi>, <https://github.com/UFAL-DSG/pyfst>

⁵<http://www.sigdial.org/>

A. Acronyms

DNN	Deep Neural Networks	4
SLU	Spoken Language Understanding	3
ASR	Automatic Speech Recognition	3
FST	Finite State Transducer	23
DFT	Discrete Fourier Transformation	9
GPU	Graphics Processing Unit	4
HTK	Hidden Markov Model Toolkit	1
EM	Expectation Maximization	12
PDF	Probability Density Function	24
OOV	Out of Vocabulary Word	30
RTF	Real Time Factor	20
HLDA	Heteroscedastic Linear Discriminant Analysis	10
HMM	Hidden Markov Model	7
LDA	Linear Discriminant Analysis	10
LM	Language Model	7
AM	Acoustic Model	4
MFCC	Mel Frequency Cepstral Coefficients	8
PLP	Perceptual Linear Prediction	8
PTI	Public Transport Information	30
IID	Independent and Identically Distributed	13
MLE	Maximum Likelihood Estimation	13
MLLT	Maximum Likelihood Linear Transform	10
CMVN	Cepstral Mean and Variance Normalisation	10
STC	Semi-Tied Covariance	10
ET	Exponential Transform	10
MMI	Maximum Mutual Information	27
bMMI	Boosted Maximum Mutual Information	31
PDF	Probability Density Function	24
MPE	Minimum Phone Error	31
PLP	Perceptual Linear Prediction	8
SER	Sentence Error Rate	20
SDS	Spoken Dialog System	4
WER	Word Error Rate	19
LMW	Language Model Weight	15
VAD	Voice Activity Detection	41
DM	Dialogue Manager	41
TTS	Text to Speech	25

NLG	Natural Language Generation	41
WST	Word Symbol Table	45

B. CD content

The CD content contains source code of projects developed, extended or modified as implementation part of this thesis. The thesis texts describes my work on projects listed below:

- Alex — Alex Dialogue System Framework where I added following files and directories:
 - *alex/components/asr/kaldi.py* — ASR component interfacing *PyOnlineLatgenRecogniser*
 - *alex/tools/kaldi/* — Kaldi training scripts modified for Alex
 - *alex/applications/PublicTransportInfoCs/hclg/* — Decoding graph (*HCLG*) scripts, and scripts for ASR evaluation.
- The Kaldi toolkit — Speech recognition toolkit where I added directories:
 - *src/onl-rec* — Implementation of *OnlineLatgenRecogniser* and utilities
 - *src/pykaldi* — Python wrapper *PyOnlineLatgenRecogniser* and utilities
 - *egs/vystadial/s5* — Training scripts for acoustic modelling¹
 - *egs/vystadial/online_demo* — Demos using using *OnlineLatgenRecogniser* and *PyOnlineLatgenRecogniser*.
- Pyfst — Python wrapper of OpenFst, where I improved installation and added several simple functions. Note I forked the original pyfst library.
- Pykaldi-eval — Repository for evaluation *OnlineLatgenRecogniser* written in IPython notebook. See interesting graphs.
- thesis.pdf
- Reference documentation for C++ code in *kaldi/src/onl-rec*.
- Reference documentation for Python code in *kaldi/src/pykaldi*.
- The reference documentation for my code in Alex.
- Related papers — papers where I am main author or co-author and are related to this work.

¹The same scripts were integrated into Kaldi svn trunk repository. However, the scripts are separated for Czech and English data. See http://sourceforge.net/p/kaldi/code/HEAD/tree/trunk/egs/vystadial_cz/ and http://sourceforge.net/p/kaldi/code/HEAD/tree/trunk/egs/vystadial_en/

Bibliography

- [1] ADSF, *The Alex Dialogue Systems Framework*, April 2014, <https://github.com/UFAL-DSG/alex>.
- [2] Lee Akinobu, *Open-Source Large Vocabulary CSR Engine Julius*, April 2014, http://julius.sourceforge.jp/en_index.php.
- [3] Cyril Allauzen, Michael Riley, Johan Schalkwyk, Wojciech Skut, and Mehryar Mohri, *Openfst: A general and efficient weighted finite-state transducer library*, Implementation and Application of Automata, Springer, 2007, pp. 11–23.
- [4] Stefan Behnel, Robert Bradshaw, Lisandro Dalcín, Mark Florisson, Vitja Makarov, and Dag Seljebotn, *Cython: C-Extensions for Python*, 2014, <http://cython.org/>.
- [5] Senaka Buthpitiya, Ian Lane, and Jike Chong, *A parallel implementation of Viterbi training for acoustic models using graphics processing units*, Innovative Parallel Computing (InPar), 2012, IEEE, 2012, pp. 1–10.
- [6] Victor Chahuneau and Ondrej Platek, *The PyFst library: OpenFst in Python*, 2014, <https://github.com/UFAL-DSG/pyfst>.
- [7] Y-L Chow, *Maximum mutual information estimation of HMM parameters for continuous speech recognition using the $\langle e1 \rangle N \langle /e1 \rangle$ -best algorithm*, Acoustics, Speech, and Signal Processing, 1990. ICASSP-90., 1990 International Conference on, IEEE, 1990, pp. 701–704.
- [8] Steven Davis and Paul Mermelstein, *Comparison of parametric representations for monosyllabic word recognition in continuously spoken sentences*, Acoustics, Speech and Signal Processing, IEEE Transactions on **28** (1980), no. 4, 357–366.
- [9] Mark JF Gales, *Semi-tied covariance matrices for hidden Markov models*, Speech and Audio Processing, IEEE Transactions on **7** (1999), no. 3, 272–281.
- [10] Zoubin Ghahramani, *Unsupervised learning*, Advanced Lectures on Machine Learning, Springer, 2004, pp. 72–112.
- [11] Ramesh A Gopinath, *Maximum likelihood modeling with Gaussian distributions for classification*, Acoustics, Speech and Signal Processing, 1998. Proceedings of the 1998 IEEE International Conference on, vol. 2, IEEE, 1998, pp. 661–664.
- [12] Hynek Hermansky, *Perceptual linear predictive (PLP) analysis of speech*, The Journal of the Acoustical Society of America **87** (1990), 1738.
- [13] Xuedong Huang, Alejandro Acero, Hsiao-Wuen Hon, et al., *Spoken language processing*, vol. 15, Prentice Hall PTR New Jersey, 2001.
- [14] David Huggins-Daines, Mohit Kumar, Arthur Chan, Alan W Black, Mosur Ravishankar, and Alex I Rudnicky, *Pocketsphinx: A free, real-time continuous speech recognition system for hand-held devices*, Acoustics, Speech and Signal Processing, 2006. ICASSP 2006 Proceedings. 2006 IEEE International Conference on, vol. 1, IEEE, 2006, pp. I–I.

- [15] Abdul J Jerri, *The shannon sampling theorem—its various extensions and applications: A tutorial review*, Proceedings of the IEEE **65** (1977), no. 11, 1565–1596.
- [16] Filip Jurčiček, *VYSTADIAL: Development of statistical methods for spoken dialogue systems*, April 2014, <http://ufal.mff.cuni.cz/grants/vystadial>.
- [17] Matěj Korvas, Ondřej Plátek, Ondřej Dušek, Lukáš Žilka, and Filip Jurčiček, *Free English and Czech telephone speech corpus shared under the CC-BY-SA 3.0 license*, Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC 2014), 2014, p. To Appear.
- [18] Akinobu Lee and Tatsuya Kawahara, *Recent development of open-source speech recognition engine julius*, 2009.
- [19] Mehryar Mohri, Fernando Pereira, and Michael Riley, *Weighted finite-state transducers in speech recognition*, Computer Speech & Language **16** (2002), no. 1, 69–88.
- [20] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar, *Foundations of machine learning*, The MIT Press, 2012.
- [21] Sirko Molau, Florian Hilger, and Hermann Ney, *Feature space normalization in adverse acoustic conditions*, Acoustics, Speech, and Signal Processing, 2003. Proceedings.(ICASSP’03). 2003 IEEE International Conference on, vol. 1, IEEE, 2003, pp. I–656.
- [22] Fabrizio Morbini, Kartik Audhkhasi, Kenji Sagae, Ron Artstein, Dogan Can, Panayiotis Georgiou, Shri Narayanan, Anton Leuski, and David Traum, *Which ASR should I choose for my dialogue system?*, Proceedings of the SIGDIAL 2013 Conference (Metz, France, 2013, pp. 394–403.
- [23] Hermann Ney, *Acoustic modeling of phoneme units for continuous speech recognition*, Proc. Fifth Europ. Signal Processing Conf, 1990, pp. 65–72.
- [24] Daniel Povey, *The Kaldi ASR toolkit*, April 2014, <http://sourceforge.net/projects/kaldi>.
- [25] Daniel Povey, Lukas Burget, Mohit Agarwal, Pinar Akyazi, Kai Feng, Arnab Ghoshal, Ondrej Glembek, Nagendra K Goel, Martin Karafiát, Ariya Rastrow, et al., *Subspace Gaussian mixture models for speech recognition*, Acoustics Speech and Signal Processing (ICASSP), 2010 IEEE International Conference on, IEEE, 2010, pp. 4330–4333.
- [26] Daniel Povey, Mark JF Gales, Do Yeong Kim, and Philip C Woodland, *MMI-MAP and MPE-MAP for acoustic model adaptation.*, INTERSPEECH, 2003.
- [27] Daniel Povey, Arnab Ghoshal, Gilles Boulianne, Lukas Burget, Ondrej Glembek, Nagendra Goel, Mirko Hannemann, Petr Motlicek, Yanmin Qian, Petr Schwarz, Jan Silovsky, Georg Stemmer, and Karel Vesely, *The Kaldi speech recognition toolkit*, IEEE 2011 Workshop on Automatic Speech Recognition and Understanding, IEEE Signal Processing Society, December 2011, IEEE Catalog No.: CFP11SRW-USB.

- [28] Daniel Povey, Mirko Hannemann, Gilles Boulianne, Lukas Burget, Arnab Ghoshal, Milos Janda, Martin Karafiát, Stefan Kombrink, Petr Motlicek, Yanmin Qian, et al., *Generating exact lattices in the wfst framework*, Acoustics, Speech and Signal Processing (ICASSP), 2012 IEEE International Conference on, IEEE, 2012, pp. 4213–4216.
- [29] Daniel Povey, Dimitri Kanevsky, Brian Kingsbury, Bhuvana Ramabhadran, George Saon, and Karthik Visweswariah, *Boosted MMI for model and feature-space discriminative training*, Acoustics, Speech and Signal Processing, 2008. ICASSP 2008. IEEE International Conference on, IEEE, 2008, pp. 4057–4060.
- [30] Daniel Povey and Brian Kingsbury, *Evaluation of proposed modifications to MPE for large scale discriminative training*, Acoustics, Speech and Signal Processing, 2007. ICASSP 2007. IEEE International Conference on, vol. 4, IEEE, 2007, pp. IV–321.
- [31] Daniel Povey, Brian Kingsbury, Lidia Mangu, George Saon, Hagen Soltau, and Geoffrey Zweig, *fMPE: Discriminatively trained features for speech recognition*, Proc. ICASSP, vol. 1, Philadelphia, 2005, pp. 961–964.
- [32] Daniel Povey, G. Zweig, and A. Acero, *The Exponential Transform as a generic substitute for VTLN*, IEEE ASRU, 2011.
- [33] Josef Psutka, *Benefit of maximum likelihood linear transform (MLLT) used at different levels of covariance matrices clustering in ASR systems*, Text, Speech and Dialogue, Springer, 2007, pp. 431–438.
- [34] Josef Psutka, Ludek Müller, and Josef V Psutka, *Comparison of MFCC and PLP parameterizations in the speaker independent continuous speech recognition task.*, INTERSPEECH, 2001, pp. 1813–1816.
- [35] Michael Riley, *OpenFst Quick Tour*, April 2014, <http://www.openfst.org/twiki/bin/view/FST/FstQuickTour>.
- [36] Luis Javier Rodríguez and Inés Torres, *Comparative study of the Baum-Welch and Viterbi training algorithms applied to read and spontaneous speech recognition*, Pattern Recognition and Image Analysis, Springer, 2003, pp. 847–857.
- [37] David Rybach, Stefan Hahn, Patrick Lehnert, David Nolden, Martin Sundermeyer, Zoltan Tüske, Siemon Wiesler, Ralf Schlüter, and Hermann Ney, *RASR-The RWTH Aachen University open source speech recognition toolkit*, Proc. IEEE Automatic Speech Recognition and Understanding Workshop, 2011.
- [38] Gabriel Skantze and David Schlangen, *Incremental dialogue processing in a micro-domain*, Proceedings of the 12th Conference of the European Chapter of the Association for Computational Linguistics, Association for Computational Linguistics, 2009, pp. 745–753.
- [39] UFAL-DSG, *The Alex Dialogue Systems Framework - Public Transport Information*, April 2014, <https://github.com/UFAL-DSG/alex>.
- [40] Keith Vertanen, *Baseline WSJ acoustic models for HTK and Sphinx: Training recipes and recognition experiments*, Tech. report, Cavendish Laboratory, University of Cambridge, 2006.

- [41] Karel Veselý, Arnab Ghoshal, Lukáš Burget, and Daniel Povey, *Sequencediscriminative training of deep neural networks*, Proc. INTERSPEECH, 2013, pp. 2345–2349.
- [42] R Weide, *The cmu pronunciation dictionary, release 0.7a*, Carnegie Mellon University, 1998.
- [43] Ian H Witten and Timothy Bell, *The zero-frequency problem: Estimating the probabilities of novel events in adaptive text compression*, Information Theory, IEEE Transactions on **37** (1991), no. 4, 1085–1094.
- [44] Xuchen Yao, Pravin Bhutada, Kallirroi Georgila, Kenji Sagae, Ron Artstein, and David R Traum, *Practical evaluation of speech recognizers for virtual human dialogue systems.*, LREC, Citeseer, 2010.
- [45] Jinjin Ye, *Speech recognition using time domain features from phase space reconstructions*, Ph.D. thesis, Marquette University Milwaukee, Wisconsin, 2004.
- [46] SJ Young, *The HTK Hidden Markov Model Toolkit: Design and Philosophy*, vol. 2, 1994, pp. 2–44.
- [47] Steve Young, Gunnar Evermann, Mark Gales, Thomas Hain, Dan Kershaw, Xunying Liu, Gareth Moore, Julian Odell, Dave Ollason, Dan Povey, et al., *The htk book (for htk version 3.4)*, Cambridge university engineering department **2** (2006), no. 2, 2–3.
- [48] Steve J. Young, Gunnar Evermann, Mark J. F. Gales, Thomas Hain, Dan Kershaw, Xunying Liu, Gareth Moore, Julian Odell, Dave Ollason, Dan Povey, Valtcho Valtchev, and Phil C. Woodland, *The HTK book, version 3.4*, Cambridge University Engineering Department, Cambridge, UK, 2006.
- [49] Xiaohui Zhang, Jan Trmal, Daniel Povey, and Sanjeev Khudanpur, *Improving deep neural network acoustic models using generalized maxout networks*, submitted to ICASSP (2014).