Technische Universität Ilmenau

Fakultät für Elektrotechnik und Informationstechnik



# Application of Speech Recognition Algorithms to Singing

PhD Thesis at Fraunhofer Institute for Digital Media Technology

| | |
|---|---|
| **Submitted by:** | Anna Marie Kruspe |
| **Submitted on:** | |
| **Course of study:** | Media Technology |
| **Matriculation Number:** | 39909 |

**Advisor:**   Prof. Dr.-Ing. Dr. rer. nat. h.c. mult. Karlheinz Brandenburg

**Abstract**

The Higgs boson or Higgs particle is an elementary particle initially theorised in 1964,[6][7] and tentatively confirmed to exist on 14 March 2013.[8] The discovery has been called "monumental"[9][10] because it appears to confirm the existence of the Higgs field,[11][12] which is pivotal to the Standard Model and other theories within particle physics. In this discipline, it explains why some fundamental particles have mass when the symmetries controlling their interactions should require them to be massless, and?linked to this?why the weak force has a much shorter range than the electromagnetic force.

**Kurzfassung**

Das Higgs-Teilchen gehört zum Higgs-Mechanismus, einer schon in den 1960er Jahren vorgeschlagenen Theorie, nach der alle fundamentalen Elementarteilchen (beispielsweise das Elektron) ihre Masse erst durch die Wechselwirkung mit dem allgegenwärtigen Higgs-Feld erhalten. Als einziges Teilchen des Standardmodells ist das Higgs-Boson experimentell noch nicht vollständig gesichert.

### Acknowledgements

Thanks to Leonard Hofstadter and thanks to my mee-maw.

# Table of Contents

# 1 Introduction

This is my introduction...

# 2 State of the art

## 2.1 From speech to singing

Singing presents a number of challenges for language identification when compared to pure speech [1]. To mention a few examples:

**Larger pitch fluctuations** A singing voice varies its pitch to a much higher degree than a speaking voice. It often also has very different spectral properties.

**Higher pronunciation variation** Singers are often forced by the music to pronounce certain sounds and words differently than if they were speaking them.

**Larger time variations** In singing, sounds are often prolonged for a certain amount of time to fit them to the music. Conversely, they can also be shortened or left out completely.

**Different vocabulary** In musical lyrics, words and phrases often differ from normal conversation texts. Certain words and phrases have different probabilities (e.g. higher focus on emotional topics in singing).

**Background music** adds irrelevant data (for language identification) to the signal, which acts as an interfering factor to the algorithms. It therefore should be removed or suppressed prior to the language identification, e.g. by source separation algorithms.

Most of the experiments in this work were performed on unaccompanied singing in order to remove this last difficulty for the moment.

## 2.2 Phoneme recognition

### 2.2.1 Phoneme recognition in speech

### 2.2.2 Phoneme recognition in singing

As described in [2], [1], and [3], there are significant differences between speech and singing audio, such as pitch and harmonics, vibrato, phoneme durations and pronunciation. These factors make phoneme recognition on singing more difficult than on speech. It has only been a topic of research for the past few years.

Fujihara et al. first presented an approach using Probabilistic Spectral Templates to model phonemes in [4]. The phoneme models are gender-specific and only model five vowels, but also work for singing with instrumental accompaniment. The best result is 65% correctly classified frames.

In [5], Gruhne et al. describe a classical approach that employs feature extraction and various machine learning algorithms to classify singing into 15 phoneme classes. It also includes a step that removes non-harmonic components from the signal. The best result of 58% correctly classified frames is achieved with Support Vector Machine (SVM) classifiers. The approach is expanded upon in [6].

Mesaros presented a complex approach that is based on Hidden Markov Models which are trained on Mel-Frequency Cepstral Coefficients (MFCCs) and then adapted to singing using three phoneme classes separately [7][8]. The approach also employs language modeling and has options for vocal separation and gender and voice adaptation. The achieved phoneme error rate on unaccompanied singing is 1.06 without adaptation and 0.8 with singing adaptation using 40 phonemes (the error rate greater than one means that there were more insertion, deletion, or substitution errors than phoneme instances). The results also improve when using gender-specific adaptation (to an average of 0.81%) and even more when language modeling is included (to 0.67%).

Hansen presents a system in [9] which combines the results of two Multilayer Perceptrons (MLPs), one using MFCC features and one using TRAP (Temporal Pattern) features. Training is done with a small amount of singing data. Viterbi decoding is then performed on the resulting posterior probabilities. On a set of 27 phonemes, this approach achieves a recall of up to 48%.

## 2.3 Forced alignment

### 2.3.1 Forced alignment in speech

### 2.3.2 Forced alignment in singing

## 2.4 Language identification

### 2.4.1 Language identification in speech

Language identification has been extensively researched in the field of Automatic Speech Recognition since the 1980's. A number of successful algorithms have been developed over the years. An overview over the fundamental techniques is given by Zissman in [**?**].

Fundamentally, four properties of languages can be used to discriminate between them:

**Phonetics** The unique sounds that are used in a given language.

**Phonotactics** The probabilities of certain phonemes and phoneme sequences.

**Prosody** The "melody" of the spoken language.

**Vocabulary** The possible words made up by the phonemes and the probabilities of certain combinations of words.

Even modern system mostly focus on phonetics and phonotactics as the distinguishing factors between languages. Vocabulary is sometimes exploited in the shape of language models.

Zissman mentions Parallel Phone Recognition followed by Language Modeling (PPRLM) as one of the basic techniques. It requires audio data, language annotations, and phoneme annotations for each utterance. In order to make use of vocabulary characteristics, full sentence annotations and word-to-phoneme dictionaries are also necessary. Using the audio and phoneme data, acoustic models are trained. They describe the probabilities of certain sound and sound sequences occurring. This is done separately for each considered language. Similarly, language models are generated using the sentence annotations and the dictionary. These models describe the probabilities of certain words and phrases. Again, this is done for each language.

New audio examples are then run through all pairs of acoustic and language models, and the likelihoods produced by each model are retained. The highest acoustic likelihood, the highest language likelihood, or the highest combined likelihood are then

considered to determine the language. This approach achieves up to 79% accuracy for ten languages [**?**].

Another approach uses the idea to train Gaussian Mixture Models for each language. This technique can be considered a "bag of frames" approach, i.e. the single data frames are considered to be statistically independent of each other. The generated GMMs then describe probability densities for certain characteristics of each language. Using these, the language of new audio examples can be easily determined.

GMM approaches used to perform worse than their PPRLM counterparts, but the development of new features has made the difference negligible [**?**]. They are in general easier to implement since only audio examples and their language annotations are required. Allen et al. [**?**] report results of up to 76.4% accuracy for ten languages. Different backend classifiers, such as Multi-Layer Perceptrons (MLPs) and Support Vector Machines (SVMs) [**?**] have also been used succesfully instead of GMMs.

## 2.4.2 Language identification in singing

So far, only a few approaches to perform language identification on singing have been proposed.

Schwenninger et al. [10] use MFCC features, but do not mention how they perform their actual model training. They test different pre-processing techniques, such as vocal/non-vocal segmentation, distortion reduction, and azimuth discrimination. None of these techniques seem to improve the over-all results. They achieve an accuracy of 68% on a-capella music for two languages (English and German).

The approach of Tsai and Wang [11] follows a traditional PPRLM flow. After vocal/non-vocal segmentation using GMMs, they run their data through acoustic models using vector tokenization. One acoustic model for each language is used. The results are then processed by bigram language models, again for each language. The language model score is used for a maximum likelihood decision to determine the language. They achieve results of 70% accuracy for two languages (English and Mandarin) on pop music.

Mehrabani and Hansen [12] also use a PPRLM system, with the difference that all combinations of acoustic and language models are tested. Their scores are combined by a classifier to determine the final language. This results in a score of 78% for a-capella music in three languages (English, Hindi, and Mandarin). Combining this technique with prosodic data improved the result even further.

Finally, Chandrasekhar et al.[13] try to determine the language for music videos using

both audio and video features. They achieve accuracies of close to 50% for 25 languages. It is interesting to note that European languages seem to achieve much lower accuracies than Asian and Arabic ones. English, French, German, Spanish and Italian rank below 40%, while languages like Nepali, Arabic, and Pashto achieve accuracies above 60%.

## 2.5 Keyword spotting

To the best of the author's knowledge, no keyword spotting systems for singing existed prior to this work.

# 3 Technical Background

## 3.1 General processing chain

## 3.2 Audio features

### 3.2.1 Perceptive Linear Predictive features (PLPs)

PLP features, first introduced in [**?**], are among the most frequently used features in speech processing. They are based on the idea to use knowledge about human perception to emphasize important speech information in spectra while minimizing the differences between speakers. We use a model order of 13 in two experiments and one of 32 in another. Deltas and double deltas between frames are also calculated. We test PLPs with and without RASTA pre-processing [**?**].

### 3.2.2 Mel-Frequency Cepstral Coefficients (MFCCs)

Just like PLPs, MFCCs are frequently used in all disciplines of automatic speech recognition [**?**]. We kept 20 cepstral coefficients for model training. Additionally, we calculated deltas and double deltas.

### 3.2.3 Shifted Delta Cepstrum (SDCs)

Shifted Delta Cepstrum features were first described in [**?**] and have since been successfully used for speaker verification and language identification tasks on pure speech data [**?**] [**?**] [**?**]. They are calculated on MFCC vectors and take their temporal evolution into account. Their configuration is described by the four parameter $N - d - P - k$, where $N$ is the number of cepstral coefficients for each frame, $d$ is the time context (in frames) for the delta calculation, $k$ is the number of delta blocks to use, and $P$ is the shift between consecutive blocks. The delta cepstrals are then calculated as:

$$\Delta c(t) = c(t + iP + d) + c(t + iP - d), 0 <= i <= k \qquad (3.1)$$

with $c \in [0, N-1]$ as the previously extracted cepstral coefficients. The resulting $k$ delta cepstrals for each frame are concatenated to form a single SDC vector of the length $kN$. We used the common parameter combination $N = 7, d = 1, P = 3, k = 7$.

### 3.2.4 TempoRal Patterns (TRAP)

## 3.3 Machine learning algorithms

This section describes the various machine learning algorithms employed throughout this thesis. Gaussian Mixture Models (GMMs), Hidden Markov Models (HMMs), and Support Vector Machines (SVMs) are three traditional approaches that are used as the basis of many new approaches, and were used for several starting experiments. i-Vector processing is a relatively new, more sophisticated approach that bundles several other machine learning techniques.

In recent years, Deep Learning has become the standard for machine learning applications []. This chapter also describes two of those new approaches that were used in this work: Deep Neural Networks (DNNs) and Deep Belief Networks (DBNs).

### 3.3.1 Gaussian Mixture Models

### 3.3.2 Hidden Markov Models

### 3.3.3 Support Vector Machines

### 3.3.4 i-Vector processing

I-Vector (identity vector) extraction was first introduced in [**?**] and has since become a state-of-the-art technique for various speech processing tasks, such as speaker verification, speaker recognition, and language identification [**?**]. To our knowledge, it has not been used for any Music Information Retrieval tasks before.

The main idea behind i-vectors is that all training utterances contain some common trends, which effectively add irrelevance to the data in respect to training. Using i-vector extraction, this irrelevance can be filtered out, while only the unique parts of the data relevant to the task at hand remain. The dimensionality of the training data is massively reduced, which also makes the training less computationally expensive. As a side effect, all feature matrices are transformed to i-vectors of equal length, eliminating problems that are caused by varying utterance lengths.

Mathematically, this assumption can be expressed as:

$$M(u) = m + Tw \tag{3.2}$$

In this equation, $M(u)$ is the GMM supervector for utterance $u$. The supervector approach was first presented in [**?**] and has since been successfully applied to a number of speech recognition problems. A music example can be found in [**?**]. $m$ represents the language- and channel-independent component of $u$ and is estimated using a Universal Background Model (UBM). $T$ is a low-rank matrix modeling the relevant language- and channel-related variability, the so-called Total Variability Matrix. Finally, $w$ is a normally distributed latent variable vector: The i-vector for utterance $u$.

**Step 1: UBM training** A Universal Background Model (UBM) is trained using Gaussian Mixture Models (GMMs) from all utterances. This UBM models the characteristics that are common to all of them.

**Step 2: Statistics extraction** 0th and 1st order Baum-Welch statistics are calculated for each of the utterances from the UBM according to:

$$N_c(u) = \sum_{t=1}^{L} P(c|y_t, \Omega) \tag{3.3}$$

$$\widetilde{F}_c(u) = \sum_{t=1}^{L} P(c|y_t, \Omega)(y_t - m_c) \tag{3.4}$$

where $u = y_1, y_2, ..., y_L$ denotes an utterance with $L$ frames, $c = 1, ..., C$ denotes the index of the Gaussian component, $\Omega$ denotes the UBM, $m_c$ is the mean of the UBM mixture component $c$, and $P(c|y_t, \Omega)$ denotes the posterior probability that the frame $y_t$ was generated by mixture component $c$. As the equation shows, the 1st order statistics are centered around the mean of each mixture component.

**Step 3: T matrix training** Using the Baum-Welch statistics for all utterances, the Total Variability Matrix $T$ is now trained iteratively according to:

$$w = (I + T^t \Sigma^{-1} N(u) T)^{-1} T^t \Sigma^{-1} \widetilde{F}(u) \tag{3.5}$$

using Expectation Maximization.

**Step 4: Actual i-vector extraction**  Finally, an i-vector $w$ can be extracted for each utterance using equation 3.5 again. This can also be done for unseen utterances, using a previously trained $T$.

### 3.3.5 Artificial Neural Networks

#### 3.3.5.1 Deep Neural Networks

#### 3.3.5.2 Deep Belief Networks

## 3.4 Evaluation

### 3.4.1 Evaluation of phoneme recognition and alignment tasks

### 3.4.2 Evaluation of language identification tasks

### 3.4.3 Evaluation of keyword spotting tasks

## 3.5 Common application systems

### 3.5.1 Systems for phoneme recognition

### 3.5.2 Systems for forced alignment

### 3.5.3 Systems for language identification

### 3.5.4 Systems for keyword spotting

# 4 Data sets

This chapter contains descriptions of all the data sets (or corpora) used over the course of this thesis. They are grouped into speech-only data sets, data sets of unaccompanied (=a-capella) singing, and data sets of full musical pieces with singing ("real-world" data sets).

## 4.1 Speech data sets

### 4.1.1 TIMIT

*TIMIT* is, presumably, the most widely used corpus in speech recognition research [14]. It was developed in ??? and consists of ??? English-language audio recordings of native speakers with annotations on the phoneme, word, and sentence levels. The corpus is split into a training and a test section, with the training section containing 4620 utterances, and the test section containing 1680. Each of those utterances has a duration of a few seconds.

The phoneme annotations follow the ??? model of ??? phonemes.

### 4.1.2 NIST Language identification corpora

### 4.1.3 OGI Language identification corpus

For comparison, the algorithms in this work were also tested on the *OGI Multi-language Telephone Speech Corpus (OGIMultilang)* [?], using all recordings for the three previously mentioned languages. There are 3,177 utterances in sum with more varying durations (1-60 seconds). For experiments on longer recordings, results on these individual utterances were aggregated for each speaker, producing 118 documents per language (354 in sum).

Table 4.1: *Amounts of data in the three used data sets: Sum duration on top, number of utterances in italics.*

| hh:mm:ss #*Utterances* | NIST2003LRE | OGIMultilang | YTAcap |
|:---:|:---:|:---:|:---:|
| English | 00:59:08 *240* | 05:13:17 *1912* | 08:04:25 *1975* |
| German | 00:59:35 *240* | 02:52:27 *1059* | 04:18:57 *1052* |
| Spanish | 00:59:44 *240* | 03:05:45 *1151* | 07:21:55 *1810* |

## 4.2 A-Capella singing data sets

### 4.2.1 YouTube data set

As opposed to the speech case, there are no standardized corpora for sung language identification. For the sung language identification experiments, A-capella audio files were therefore extracted from *YouTube*[1] videos. This was done for three languages: English, German, and Spanish. The author collected between 116 (258min) and 196 (480min) examples per language. These were mostly videos of amateur singers freely performing songs without accompaniment. Therefore, they are of highly varying quality and often contain background noise. Most of the performers contributed only a single song, with just a few providing up to three. In this way, we aim to avoid effects where the classifier recognizes the singer's voice instead of the language.

Special attention was paid to musical style. Rap, opera singing, and other specific singing styles were excluded. All the songs performed in these videos were pop songs. Different musical styles can have a high impact on language classification results. The author tried to limit this influence as much as possible by choosing recordings of pop music instead of language-specific genres (such as latin american music).

### 4.2.2 Hansen's vocal track data set

This is one of the data sets used for keyword spotting and phoneme recognition. It was first presented in [9]. It consists of the vocal tracks of 19 commercial English-language

---

[1]`http://www.youtube.com`, Last checked: 05/16/13

pop songs. They are studio quality with some post-processing applied (EQ, compression, reverb). Some of them contain choir singing. These 19 songs are split up into ??? clips that roughly represent lines in the song lyrics.

Twelve of the songs were annotated with time-aligned phonemes. The phoneme set is the one used in CMU Sphinx[2] and TIMIT [14] and contains 39 phonemes. All of the songs were annotated with word-level transcriptions. This is the only one of the singing data sets that has full manual annotations, which are assumed to be reliable and can be used as ground truth.

For comparison, recordings of spoken recitations of all song lyrics were also made. These were all performed by the same speaker (the author).

### 4.2.3 DAMP data set

As described, Hansen's data set is very small and therefore not suited to training phoneme models for singing. As a much larger source of unaccompanied singing, the *DAMP* data set, which is freely available from Stanford University[3][15], was employed. This data set contains more than 34,000 recordings of amateur singing of full songs with no background music, which were obtained from the *Smule Sing!* karaoke app. Each performance is labeled with metadata such as the gender of the singer, the region of origin, the song title, etc. The singers performed 301 English-language pop songs. The recordings have good sound quality with little background noise, but come from a lot of different recording conditions.

No lyrics annotations are available for this data set, but the textual lyrics can be obtained from the *Smule Sing!* website[4]. These are, however, not aligned in any way. Such an alignment was performed automatically on the word and phoneme levels (see section ??).

### 4.2.4 Aji's synthesized singing data set

Since it was not feasible to hand-annotate a large data set over the course of this work, another approach was the automatic generation of sung audio. The advantage of this approach is that the results can be assumed to be perfectly aligned to the given phonemes.

---

[2] http://cmusphinx.sourceforge.net/
[3] https://ccrma.stanford.edu/damp/
[4] http://www.smule.com/songs

For the generation of this data set, ???  recordings from the previously described *DAMP* data set were selected. Their phonemes were automatically aligned, and an automatic transcription of the melody was performed. These two sources of data were then aligned to each other. This alignment did not need to be perfect, it just needed to produce a plausible combination of melody line and phonemes.

The result of this step was then fed into the $Sinsy[]^5$ singing synthesizer to generate new singing recordings. This synthesizer provided one female singing voice. The resulting recordings are good in quality and relatively natural sounding, but appear to have a slight accent.

The whole generation process of this data set was performed in collaboration with Adam Aji.

### 4.2.5 Choosing keywords

## 4.3 "Real-world" data sets

### 4.3.1 QMUL Expletive data set

This data set consists of 80 popular songs which were collected at Queen Mary University, most of them Hip Hop. 711 instances of 48 expletives were annotated on these songs. In addition, the matching textual, unaligned lyrics were retrieved from the internet.

### 4.3.2 "69 Love Songs" data set

"69 Love Songs" is a 3-CD album by the band "The Magnetic Fields", which was released in ???  and named one of the ???. It contains 69 songs in various musical styles and instrumentations, performed by a variety of musicians, including 4???  vocalists. The total duration is ???. The data set is interesting for the purposes of this work because the songs' lyrics all cover a similar theme - namely, love. A word count on the lyrics shows, for example, that the word "love" itself occurs $\tilde{2}25$ times in these songs.

Unaligned lyrics were retrieved from ???. A thorough semantic analysis can be found in [**?**].

---

5

# 5 Singing phoneme recognition and alignment

## 5.1 Phoneme recognition using models trained on speech

### 5.1.1 State-of-the-art processing chain

### 5.1.2 Results

## 5.2 Phoneme recognition using models trained on "songified" speech

### 5.2.1 Modifications to the TIMIT data set

### 5.2.2 Results

## 5.3 Phoneme recognition using models trained on a-capella singing

### 5.3.1 Modifications to the training process

### 5.3.2 Results

## 5.4 Phoneme recognition on synthesized singing

### 5.4.1 Approach

### 5.4.2 Results

## 5.5 Forced alignment

## 5.6 Conclusion

# 6 Language identification

## 6.1 LID in singing using GMMs

### 6.1.1 Processing chain

### 6.1.2 Results

## 6.2 LID in singing using i-Vectors and GMMs

### 6.2.1 i-Vector implementation

### 6.2.2 i-Vector processing chain

### 6.2.3 Results

## 6.3 LID in singing using phoneme recognition posteriors

### 6.3.1 Phoneme recognition for LID

### 6.3.2 Post-processing

### 6.3.3 Results

## 6.4 Conclusion

# 7 Sung keyword spotting experiments and results

## 7.1 Keyword spotting using keyword-filler HMMs

### 7.1.1 Phoneme posterior extraction and further processing

### 7.1.2 Implementation of keyword-filler HMMs

### 7.1.3 Results on speech and music

## 7.2 Keyword spotting using duration-informed keyword-filler HMMs

### 7.2.1 Duration modeling approaches

### 7.2.2 Implementation of duration modeling approaches for keyword-filler HMMs

### 7.2.3 Results on speech and music

## 7.3 Improving keyword spotting using specified phoneme models

### 7.3.1 Improving phoneme models

### 7.3.2 Post-processing

### 7.3.3 Results on speech and music

## 7.4 Conclusion

# 8 Applications

## 8.1 Experiments on the QMUL Expletive data set

### 8.1.1 Implementation

### 8.1.2 Results

## 8.2 Experiments on the the "69 Love Songs" data set

### 8.2.1 Implementation

### 8.2.2 Results

## 8.3 Automatic lyrics retrieval

### 8.3.1 Implementation

### 8.3.2 Results

# 9 Conclusion

# 10 Future work

# Bibliography

[1] H. Fujihara and M. Goto, *Multimodal Music Processing*, chapter Lyrics-to-audio alignment and its applications, Dagstuhl Follow-Ups, 2012.

[2] A. Loscos, P. Cano, and J. Bonada, "Low-delay singing voice alignment to text," in *Proceedings of the ICMC*, 1999.

[3] A. M. Kruspe, "Keyword spotting in a-capella singing," in *15th International Conference on Music Information Retrieval (ISMIR)*, Taipei, Taiwan, 2014.

[4] H. Fujihara, M. Goto, and H. G. Okuno, "A novel framework for recognizing phonemes of singing voice in polyphonic music.," in *WASPAA*. 2009, pp. 17–20, IEEE.

[5] M. Gruhne, K. Schmidt, and C. Dittmar, "Phoneme recognition on popular music," in *Proceedings of the 8th International Conference on Music Information Retrieval (ISMIR)*, Vienna, Austria, September 2007.

[6] G. Szepannek, M. Gruhne, B. Bischl, S. Krey, T. Harczos, F. Klefenz, C. Dittmar, and C. Weihs, *Classification as a tool for research*, chapter Perceptually Based Phoneme Recognition in Popular Music, Springer, Heidelberg, 2010.

[7] A. Mesaros and T. Virtanen, "Recognition of phonemes and words in singing.," in *ICASSP*. 2010, pp. 2146–2149, IEEE.

[8] A. Mesaros and T. Virtanen, "Automatic recognition of lyrics in singing.," *EURASIP J. Audio, Speech and Music Processing*, vol. 2010, 2010.

[9] J. K. Hansen, "Recognition of phonemes in a-cappella recordings using temporal patterns and mel frequency cepstral coefficients," in *9th Sound and Music Computing Conference (SMC)*, Copenhagen, Denmark, 2012, pp. 494–499.

[10] J. Schwenninger, R. Brueckner, D. Willett, and M. E. Hennecke, "Language identification in vocal music," in *7th International Conference on Music Information Retrieval (ISMIR)*, Victoria, Canada, 2006, pp. 377–379.

[11] W.-H. Tsai and H.-M. Wang, "Towards automatic identification of singing language in popular music recordings," in *5th International Conference on Music Information Retrieval (ISMIR)*, Barcelona, Spain, 2004, pp. 568–576.

[12] M. Mehrabani and J. H. L. Hansen, "Language identification for singing," in *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, Prague, Czech Republic, 2011, pp. 4408–4411.

[13] V. Chandraskehar, M. E. Sargin, and D. A. Ross, "Automatic language identification in music videos with low level audio and visual features," in *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, Prague, Czech Republic, 2011, pp. 5724–5727.

[14] J. S. Garofolo et al., "TIMIT Acoustic-Phonetic Continuous Speech Corpus," Tech. Rep., Linguistic Data Consortium, Philadelphia, 1993.

[15] J. C. Smith, *Correlation analyses of encoded music performance*, Ph.D. thesis, Stanford University, 2013.

# List of Figures

# List of Tables

# A  Appendix

A lot of stuff that didn't fit into the main part ...

# B Eigenständigkeitserklärung

Die vorliegende Arbeit habe ich selbstständig ohne Benutzung anderer als der angegebenen Quellen angefertigt.

Alle Stellen, die wörtlich oder sinngemäß aus veröffentlichten Quellen entnommen wurden, sind als solche deutlich kenntlich gemacht. Die Arbeit ist in gleicher oder ähnlicher Form oder auszugsweise im Rahmen einer oder anderer Prüfungen noch nicht vorgelegt worden.

Ilmenau, 17.12.2013

Sheldon Cooper

# Thesis Summary

1. Scissors cuts paper, paper covers rock, rock crushes lizard, lizard poisons Spock, Spock smashes scissors, scissors decapitates lizard, lizard eats paper, paper disproves Spock, Spock vaporizes rock, and as it always has, rock crushes scissors.

2. I'm not insane, my mother had me tested!

3. All I need is a healthy ovum and I can grow my own Leonard Nimoy!

# Thesen

1. These 1

2. These 2

3. These 3