

Bachelor thesis - Better accuracy of automatic  
lecture transcriptions by using context  
information from slide contents

Jonathan Werner

# Contents

<b>Introduction</b>	<b>4</b>
Structure of this thesis . . . . .	5
<b>1 Research questions</b>	<b>7</b>
<b>2 Background</b>	<b>8</b>
2.1 The field of Automatic Speech Recognition . . . . .	8
2.2 Dimensions of speech recognition . . . . .	8
2.3 Concepts . . . . .	9
2.3.1 Phonemes . . . . .	10
2.3.2 Phonetic dictionaries . . . . .	11
2.3.3 Acoustic models . . . . .	12
2.3.4 Language Models . . . . .	12
2.4 Work done on ASR for lecture transcription . . . . .	14
2.4.1 Generalization approaches . . . . .	14
2.4.2 Specialization approaches . . . . .	15
2.5 Metrics . . . . .	15
<b>3 Test data</b>	<b>15</b>
3.1 Materials overview . . . . .	16
3.1.1 Conclusion . . . . .	17
<b>4 The LM-Interpolation approach</b>	<b>18</b>
4.1 Sphinx 4 . . . . .	18
4.2 Implementation . . . . .	18
<b>5 Analysis</b>	<b>25</b>
5.1 Approaching a good metric . . . . .	25
5.1.1 Lecture-scoped WER excluding $top_X$ words . . . . .	25
5.1.2 Lemmas . . . . .	25
5.1.3 Proposed metric: KWER-x . . . . .	26

5.1.4	Secondary metrics . . . . .	27
5.2	Results . . . . .	28
5.3	Interpretation . . . . .	29
5.3.1	Qualitative interpretation . . . . .	30
<b>6</b>	<b>Visualization for scannability</b>	<b>32</b>
<b>7</b>	<b>Summary, improvements</b>	<b>35</b>
	<b>References</b>	<b>35</b>

## Introduction

Scannability is crucial for academic research: you have to be able to quickly evaluate the usefulness of a given resource by skimming the content and looking for the parts that are specifically relevant to the task at hand.

The medium in which those resources are available is very centered on textual representation. Spoken content, hereinafter called *speech media* (audio- or audiovisual media that mainly consist of spoken language) doesn't make it possible to scan its contents. You are "stabbing in the dark" when looking for something specific in a medium like this and have to consume it like a linear narrative.

This means that although lectures and conference talks are a central element to science they are much more challenging and tedious to use for research work.

Being able to a) efficiently search and b) look at the temporal distribution of important keywords in a visually dense way would increase the usefulness of speech media in the scientific context immensely.

One approach to accomplish these goals is to utilize Automatic Speech Recognition (ASR) in order to transcribe speech to text and also get timing information for the recognized words. This makes it possible to derive information about the density of given words at a given point of time in the talk, which in turn allows to compute word occurrence density maxima. This opens up possibilities for compact visual representation of the interesting keywords, thus allowing the user to scan.

The main challenge when using ASR for this task is the recognition accuracy of technical terms. Most of them are not included in the language models that are available as these are broad and generic so as to optimize accuracy over a wide topic spectrum. But when they are not included in the language model they have a very small chance to be correctly recognized at all.

So the usefulness of applying ASR with a generic language model to the problem is very small, as the intersection of interesting keywords with those technical terms that can not be recognized is very big.

The central goal of this thesis is to explore an approach to overcome this problem. This approach consists of using words from lecture slides or other notes to generate a lecture-specific language model. This is then interpolated with a generic language model. Finally the results are compared with the 'baseline' accuracy of the generic model.

## Structure of this thesis

The structure of this thesis is as follows:

(1) **Research questions**

I will state the research questions.

(2) **Scientific Background**

- (a) I will start by giving an overview over the state of the art of ASR and the most prevalent approaches.
- (b) I will explain the *concepts* which are fundamental for the understanding of speech recognition.
- (c) I will then examine the *scientific work* that has been done on applying ASR to the problem of lectures transcriptions.
- (d) Finally i will summarize the *metrics* that have been used to assess the quality of the improvements in different approaches.

(3) **Motivation**

Here i will motivate why it is necessary to improve on the baseline performance of ASR in our context.

I will talk about the role of keywords and technical terms and why they are not being detected and how that diminishes the usefulness of ASR for the purposes of scannability.

(4) **Test data**

I will use the openly available *Open Yale Courses* (*Open Yale Courses Website*, n.d.), which provide a diverse selection of audio and video recordings of university lectures at Yale, additionally supplying quality manual transcriptions and course notes or slides.

I will present the chosen courses, their selection criteria and discuss the range of types of lecture material.

(5) **The LM-Interpolation approach**

(a) **Technical basis**

I will introduce the open source speech recognition framework *Sphinx 4*. This is the software that is used for performing the actual recognition.

(b) **Process overview**

I will then give a overview of the design and architecture of our approach.

(c) **Implementation**

Finally i will describe the technical implementation by which the lecture material is compiled into a specialized language model and recognition is performed using a *interpolated* language model.

(6) **Analysis**

(a) **Methods**

I will discuss how to analyze the results and develop metrics that assess how well the given goals are met with our approach.

(b) **Analysis**

I will then perform quantitative analysis on our test dataset with the metrics we developed before.

(c) **Discussion, Finding and Conclusions**

I will discuss the findings and draw conclusions from the quantitative analysis concerning the effectiveness of our approach.

(7) **Visualization for Scannability**

I will present a prototype visualization method that uses the results from our approach to present a condensed representation of the keyword content from lectures with the goal of providing a quick, interactive way to search and scan speech media.

(8) **Improvements, Open Ends**

I will discuss possible improvements and open ends that were out of the scope of this thesis but would be interesting to explore further.

(9) **Summary**

I will end by summarizing the goals, the proposed approach, the design and implementation, the analysis and the results.

# 1 Research questions

The central research questions I want to investigate in this thesis can be formulated as follows:

- (1) When we apply ASR to university lectures, what is the advantage of using an approach that consists of creating a lecture-specific language model and interpolating it with a generic language model, given that we are interested in improving the recognition accuracy of *interesting keywords* for the sake of searchability and scannability?
- (2) What metric is useful for quantifying this advantage?

A secondary question is: How can we *use* the results from our approach to provide graphical interfaces for improving the user's ability to search and scan the given speech medium?

The exploration of this question will not be at the center of this thesis, but it will provide practical motivation for the results of our approach.

## 2 Background

### 2.1 The field of Automatic Speech Recognition

Automatic Speech Recognition (ASR) can be defined as the process by which a computer maps an acoustic speech signal to text (*comp.speech Frequently Asked Questions*, n.d.).

Rabiner & Juang (1993) date the first research on ASR back to the early 1950s, when Bell Labs built a system for single-speaker digit recognition. Since then the field has seen three major approaches, which Marquard (2012) summarizes as follows:

1. The *acoustic-phonetic approach* aimed to identify features of speech such as vowels directly through their acoustic properties, and from there build up words based on their constituent phonetic elements.
2. The *statistical pattern-recognition approach* measures features of the acoustic signal, and compares these to existing patterns established from a range of reference sources to produce similarity scores which may be used to establish the best match.
3. *Artificial intelligence (AI) approaches* have been used to integrate different types of knowledge sources (such as acoustic, lexical, syntactic, semantic and pragmatic knowledge) to influence the output from a pattern-recognition system to select the most likely match.

The most prevalent approach today is the *statistical pattern-recognition approach*, as it produces results with much higher accuracy compared to the acoustic-phonetic approach. The use of Hidden Markov Models (HMM) has been playing a key role in this approach, as it allows recognizers to use a statistical model of a given pattern rather than a fixed representation.

In the last years there has been a resurgence of AI approaches, specifically *deep learning approaches* (Hinton et al., 2012). The ASR paradigm we will use for this thesis will be limited to the former, however.

### 2.2 Dimensions of speech recognition

There are three dimensions which serve to classify different applications of speech recognition (*comp.speech Frequently Asked Questions*, n.d., Marquard (2012)):



- (1) **Dependent vs. independent.** Dependent recognition systems are developed to be used by one speaker, whereas independent systems are developed to be used by *any* speaker of a particular type, i.e North-American speakers. **Adaptive** systems lie between these poles, they are able to adapt to a particular speaker through training.
- (2) **Small vs. large vocabulary.** Small vocabularies contain only up to a few hundred words and might be modeled by an explicit grammar, whereas large vocabularies contain tens of thousands of words so as to be able to model general purpose spoken language over a variety of domains.
- (3) **Continuous vs. isolated speech.** Isolated speech consists of single words that are spoken with pauses in between them, whereas continuous speech consists of words that are spoken in a connected way. Continuous speech is significantly more difficult to recognize, as it is a) more difficult to find the start and end of words and b) the pronunciation of words changes in relation to their surrounding words.

With these three dimensions we can for example classify the application areas command and control systems, dictation and lecture transcription (Marquard, 2012):

Table 1: Three application areas

Application	Speaker	Vocabulary	Duration
Dictation	Dependent	Large	Connected
Command and control system	Independent	Small	Isolated
Lecture transcription	Independent	Large	Connected

The task of automatic lecture transcription can thus be characterized as speaker-independent (SI) large continuous speech recognition (LVCSR).

## 2.3 Concepts

Speech recognition in the *statistical pattern-recognition approach* paradigm has three major concepts that are necessary for its understanding:

- phonemes and phonetic dictionaries
- acoustic models (AM)
- language models (LM)

### 2.3.1 Phonemes

A *phoneme* is “the smallest contrastive linguistic unit which may bring about a change of meaning” (Cruttenden, 2014, p. 43). Phonemes are the smallest unit of sound in speech which are combined to form words. The word *sun* for example can be represented by the phonemes /s/, /u/ and /n/; the word *table* by /t/, /a/ and /bɪ/.

A language with a specific accent can be described by the set of phonemes that it consists of. Figure 1 uses symbols from the International Phonetic Alphabet (IPA) to display the 44 phonemes that are being used in Received Pronunciation (RP), which is regarded as the “standard accent” in the South of the United Kingdom (Stevenson & Waite, 2011).

VOWELS		monophthongs				diphthongs		Phonemic Chart voiced unvoiced	
		i:	ɪ	ʊ	u:	ɪə	eɪ		
		sheep	ship	good	shoot	here	wait		
		e	ə	ɜ:	ɔ:	ʊə	ɔɪ	əʊ	
		bed	teacher	bird	door	tourist	boy	show	
		æ	ʌ	ɑ:	ɒ	eə	aɪ	aʊ	
		cat	up	far	on	hair	my	cow	
CONSONANTS		p	b	t	d	tʃ	dʒ	k	g
		pea	boat	tea	dog	cheese	June	car	go
		f	v	θ	ð	s	z	ʃ	ʒ
		fly	video	think	this	see	zoo	shall	television
		m	n	ŋ	h	l	r	w	j
		man	now	sing	hat	love	red	wet	yes

The 44 phonemes of Standard British English based on the popular Adon Underhill's list

adapted by EnglishClub

The 44 phonemes of Received Pronunciation based on the popular Adrian Underhill layout

adapted by EnglishClub.com

Figure 1: Phonemic Chart representing 44 phonemes used in RP British English

To be able to use phonemes in software an ASCII representation is more suitable. The standard for General American English is the *Arpabet*. Here each phoneme is mapped to one or two capital letters. The digits 0, 1 and 2 signify stress markers: no stress, primary and secondary stress respectively. A comparison of the IPA format and the arphabet format can be seen in Figure 2, an excerpt that just shows the *monophthongs*.<sup>1</sup>

<sup>1</sup>pure vowel sounds with relatively fixed articulation at the start and the end that don't glide towards a new position of articulation

Arpabet	IPA	Word examples
AO	ɔ	off (AO1 F); fall (F AO1 L); frost (F R AO1 S T)
AA	ɑ	father (F AA1 DH ER), cot (K AA1 T)
IY	i	bee (B IY1); she (SH IY1)
UW	u	you (Y UW1); new (N UW1); food (F UW1 D)
EH	ɛ	red (R EH1 D); men (M EH1 N)
IH	ɪ	big (B IH1 G); win (W IH1 N)
UH	ʊ	should (SH UH1 D), could (K UH1 D)
AH	ʌ	but (B AH1 T), sun (S AH1 N)
	ə	sofa (S OW1 F AH0), alone (AH0 L OW1 N)
AX		discus (D IH1 S K AX0 S); note distinction from discuss (D IH0 S K AH1 S)
AE	æ	at (AE1 T); fast (F AE1 S T)

Figure 2: Excerpt from the Arpabet (*English Wikipedia Arpabet article*, 2015 (accessed 22.8.15))

### 2.3.2 Phonetic dictionaries

Phonetic dictionaries map words to one or more versions of phoneme sequences.

A phonetic representation of a word is specified manually based on the knowledge of how written words *actually sound* when spoken.

An excerpt from the dictionary `cmudict-en-us.dict` (*CMU EN-US Pronouncing Dictionary* (`cmudict-en-us.dict`), 2015) looks like this:

```
...
abdollah AE B D AA L AH
abdomen AE B D OW M AH N
abdomen(2) AE B D AH M AH N
abdominal AE B D AA M AH N AH L
abdominal(2) AH B D AA M AH N AH L
...
```

The dictionary has 133.425 entries. Generally only words that are in the phonetic dictionary being used can be recognized during speech recognition. *Grapheme<sup>2</sup>-to-Phoneme converters* (G2P) however make it possible to get phoneme sequence hypotheses for arbitrary words (i.e arbitrary sequences of graphemes). While

<sup>2</sup>“The smallest unit used in describing the writing system of a language” Florian (1996), p.174

these results are on average less accurate than manually created variants, they play a vital role in texts with many technical terms as these are often not included in phonetic dictionaries.

### 2.3.3 Acoustic models

An acoustic model (AM) describes the relation between an audio signal and the probability that this signal represents a given phoneme.

Acoustic models are created by *training* them on a *corpus* of audio recordings and matching transcripts. When being used in the context of speaker-independent recognition, these models are trained with a variety of speakers that represent a broad spectrum of the language/accents that the acoustic model should represent.

During the *decoding* phase the acoustic model and a phonetic dictionary are used to match sequences of small audio “slices” to possible phonemes and those phonemes to possible word sequence hypotheses.

However, acoustic models alone are not sufficient for speech recognition as they do not have the “higher-level” linguistic information necessary to distinguish e.g. between homonyms and similar-sounding phrases such as “wreck a nice beach” and “recognize speech” (Marquard, 2012, p. 11). This information is provided by *language models*.

### 2.3.4 Language Models

Language models (LM) guide and constrain the search process that a speech recognition system performs by assigning probabilities to sequences of words. They are trained by applying statistical methods on a text corpus. Analogous to acoustic models, generic language models use huge text corpora with a broad variety of topics. It is however possible to train language models on small and specialized text corpora, which is the central technical foundation for the approach discussed in this thesis.

The most commonly used form of language models are *n-gram language models*. In the context of a language model an *n-gram* is a sequence of  $n$  words. 1-grams are called *unigrams*, 2-grams are called *bigrams* and 3-grams are called *trigrams*. An *n-gram language model* maps a set of *n-grams* to probabilities that they occur in a given piece of text.

A key idea in modelling language like this is the *independence assumption*, which says that the probability of a given word is only dependent on the last  $n - 1$  words. This assumption significantly decreases the statistical complexity and thus makes it computationally feasible.

N-gram language models do not need to be constrained to one type of n-gram. The *Generic US English Language Model (CMUSphinx US English Generic*

Language Model (*cmusphinx-5.0-en-us.lm*), 2015) from CMUSphinx we will use as the baseline for our approach consists of 1-, 2, and 3-grams, for example.

A toy example of a language model with 1- and 2-grams when represented in ARPA-format (as used by CMUSphinx) looks like follows (*CMUSphinx ARPA Language models*, 2015 (accessed 23.8.15)):

```
\data\
ngram 1=7
ngram 2=7

\1-grams:
-1.0000 <UNK>      -0.2553
-98.9366 <s>        -0.3064
-1.0000 </s>        0.0000
-0.6990 wood       -0.2553
-0.6990 cindy      -0.2553
-0.6990 pittsburgh -0.2553
-0.6990 jean       -0.1973

\2-grams:
-0.2553 <UNK> wood
-0.2553 <s> <UNK>
-0.2553 wood pittsburgh
-0.2553 cindy jean
-0.2553 pittsburgh cindy
-0.5563 jean </s>
-0.5563 jean wood

\end\
```

Here the first number in a row is the probability of the given n-gram in  $\log_{10}$  format. This means that the unigram *wood* has a probability of  $10^{-0.6990} \approx 0.2 = 20\%$  and the probability of the words “wood pittsburg” occurring in sequence is  $10^{-0.2553} \approx 0.55 = 55\%$ .

The optional third numeric column in a row is called *backoff weight*. Backoff weights make it possible to calculate n-grams that are not listed by applying the formula

$$P(\text{word}_N \mid \text{word}_{\{N-1\}}, \text{word}_{\{N-2\}}, \dots, \text{word}_1) = \\ P(\text{word}_N \mid \text{word}_{\{N-1\}}, \text{word}_{\{N-2\}}, \dots, \text{word}_2) * \\ \text{backoff-weight}(\text{word}_{\{N-1\}} \mid \text{word}_{\{N-2\}}, \dots, \text{word}_1)$$

With the side condition that missing entries for  $\text{word}_{\{N-1\}} \mid \text{word}_{\{N-2\}}, \dots, \text{word}_1$  are replaced by 1.0.

So if the text to be recognized would contain the sequence “wood cindy”, which does not appear as a bigram in the LM, the probability for this bigram could be calculated by  $P(\text{wood}|\text{cindy}) = P(\text{wood}) * \text{Bwt}(\text{cindy})$ .

Finally, the overall probability of a sentence with the words  $w_1, \dots, w_n$  can be approximated as follows:

$$P(w_1, \dots, w_n) = \prod_{n=1}^m P(w_i | w_1, \dots, w_{i-1})$$

An example approximation with a bigram model for the sentence “I saw the red house” (*English Wikipedia Language Model article*, 2015 (accessed 23.8.15)) represented as  $P(\text{I}, \text{saw}, \text{the}, \text{red}, \text{house})$  would look like

$$P(\text{I} | \langle s \rangle) \times P(\text{saw} | \text{I}) \times P(\text{the} | \text{saw}) \times P(\text{red} | \text{the}) \times P(\text{house} | \text{red}) \times P(\langle s \rangle | \text{house})$$

## 2.4 Work done on ASR for lecture transcription

I will now give an overview over the scientific work done on lecture transcription, using Marquard (2012) as a guiding reference.

The research for speech recognition on lectures can be partitioned into three general approaches: generalization approaches, specialization approaches and approaches involving the user for manual correction and improvements.

### 2.4.1 Generalization approaches

Generalization approaches try to create models that capture common characteristics of lectures. Those characteristics include highly spontaneous presentation style and “strong coarticulation effects, non-grammatical constructions, hesitations, repetitions, and filled pauses” (Yamazaki, Iwano, Shinoda, Furui, & Yokota, 2007). Glass, Hazen, Hetherington, & Wang (2004) note the “colloquial nature” of lectures as well as the “poor planning at the sentence level [and] higher structural levels”.

The generalization approach has been applied on the acoustic model level: Cettolo, Brugnara, & Federico (2004) have examined adapting a generic acoustic model to account for spontaneous speech phenomena (“filler sounds”).

While a subfield of ASR called “speaker diarization” tries to account for the interactivity between lecturers and students by identifying multiple speakers, most research treats lectures as single speaker events with the audience as background noise.

Generalization approaches at the language model level try to model common linguistic traits of the lecture genre (this can be called the “macro level”). Kato,

Nanjo, & Kawahara (2000) investigate topic-independent language modeling by creating a large corpus of text from lecture transcripts and panel discussions and then removing topic-specific keywords.<sup>3</sup>

#### 2.4.2 Specialization approaches

Specialization approaches try to use context specific to a single lecture (“meso level”) or parts of a single lecture (“micro level”<sup>4</sup>).

Methods used for creating LMs from context information can be categorized into two approaches: direct usage of lecture slides and notes for the creation of LMs versus usage of “derived” data from these materials. Deriving data by using keywords found in slides, using them as web search query terms and using the found documents as the basis for LM creation is explored in Munteanu, Penn, & Baecker (2007), Kawahara, Nemoto, & Akita (2008) and Marquard (2012).

Using the whole text from lecture slides has been explored by Yamazaki et al. (2007). They compare the *meso level* with the *micro level* by dynamically adapting the LM to the speech corresponding to a particular slide. Kawahara et al. (2008) also examine dynamic local slide-by-slide adaption and compare it to global topic adaption using Probabilistic Latent Semantic Analysis (PLSA)<sup>5</sup> and web text collection, concluding that the latter performs worse than the former because of a worse orientation to topic words. .

### 2.5 Metrics

## 3 Test data

The test data I will use for evaluating our approach will be from *Open Yale Courses*<sup>6</sup>, which is a selection of openly available lectures from Yale university. It consists of 42 courses from 25 departments. Each course has about 20-25 sessions that have an average length of 50 minutes. Each lecture is provided with good quality audio and video recordings, precise manual transcripts and lecture material when available. Only about 20% of the lectures have lecture notes or slides at all and most materials from the natural and formal science departments (physics, astronomy, mathematics) consist of hand-written notes, making them unsuitable for our approach. All talks are in English.

---

<sup>3</sup>In a second step they combine this generalization technique with a specialization technique by adapting the resulting LM with a lecture-specific language model by using preprint papers of a given lecture.

<sup>4</sup>The three levels are taken from Marquard (2012).

<sup>5</sup>Latent Semantic Analysis is an approach to document comparison and retrieval which relies on a numeric analysis of word frequency and proximity.

<sup>6</sup><http://oyc.yale.edu/>

I have chosen the following lectures: (Department, Course, Lecture Number - Title, abbreviation)

- *Biomedical Engineering*: Frontiers of Biomedical Engineering, 1 - What is Biomedical Engineering? (**biomed-eng-1**)
- *Environmental Studies*: Environmental Politics and Law, 8 - Chemically Dependent Agriculture (**environmental-8**)
- *Geology & Geophysics*: The atmosphere, the ocean, and environmental change, 8 - Horizontal transport (**geology-8**)
- *Philosophy*: Philosophy and the science of human nature, 8 - Flourishing and Detachment (**human-nature-8**)
- *Psychology*: Introduction to Psychology, 14 - What Motivates Us: Sex (**psy-14**)
- *Psychology*: Introduction to Psychology, 5 - What Is It Like to Be a Baby: The Development of Thought (**psy-5**)

The main selection criterion here was topical diversity, the challenge being that the majority of talks with computer-parsable notes was from the humanities.

### 3.1 Materials overview

The available material is very heterogeneous. I will now give an overview with excerpts which will serve as a basis for examining at a later point if the quality and quantity of the supplied material is correlated with the amount of improvement of our approach.

**geology-8** supplies a 2-page exercise sheet.

“Mars has a radius of  $3.39 \times 10^6$  m and a surface gravity of  $3.73 \text{ ms}^{-2}$ . Calculate the escape velocity for Mars and the typical speed of a  $\text{CO}_2$  molecule (assume  $T = 250 \text{ K}$ ). How can Mars retain its  $\text{CO}_2$  atmosphere? (Hint: the molecular weight of carbon dioxide is 44. Use the formulae given in class.) [...]”

**biomed-eng-1** provides a 7-page glossary of technical terms.

“[...] active transport - the transport of molecules in an energetically unfavorable direction across a membrane coupled to the hydrolysis of ATP or other source of energy

ATP (adenosine 5'-triphosphate) - a nucleotide that is the most important molecule for capturing and transferring free energy in



cells. Hydrolysis of each of the two high-energy phosphoanhydride bonds in ATP is accompanied by a large free-energy change (“G”) of 7 kcal/mole

aquaporin – a water channel protein which allows water molecules to cross the cell membrane much more rapidly than through the phospholipid bilayer [...]”

**human-nature-8** provides reading assignments for four books with short summaries each.

“[A] Epictetus, The Handbook

Background information about the Stoic philosopher Epictetus (c. 50-130 CE) and his famous work *Encheiridion* (The Handbook) appears in Nicholas White’s introduction to our translation. White has also added footnotes that explain points of potential confusion.

As the title indicates, The Handbook is intended as a tidy introduction to a more complex philosophical outlook. It is written in an accessible and engaging style.

The Stoic movement originated around 300 BCE and flourished for over five hundred years. The Stoics believed that the external world is deterministic: its state at any time is completely determined by its prior states. So, they maintained, it is pointless to wish for things to be different because to do so is to wish for something impossible. A wise person would, therefore, accept whatever befalls them without desiring that things go otherwise – hence the English word ‘stoic.’

Passages to focus on/passages to skim

I encourage you to read the text in full, at a steady reading pace. [...]”

**psy-14/5** and **enviromental-8** provide ~10-page slides with a typical amount of text.

### 3.1.1 Conclusion

Only about 20% of the courses have lecture material at all; only about 20% of these courses actually have typical “slides” – the rest provides heterogenous other kinds of material. While it cannot be inferred from this dataset that this is a general condition, it nevertheless shows a clear “real-world” disadvantage of an approach only relying on those materials. We will look at the impact of the varying quality and quantity in the analysis later on.

## 4 The LM-Interpolation approach

I will now describe the LM-Interpolation approach. The high level overview is as follows: we will use the open source speech recognition framework Sphinx 4<sup>7</sup> as the software for performing speech recognition. Sphinx 4 has a modular architecture which allows specifying components of the whole process per configuration. It provides multiple implementations of LMs<sup>8</sup>, the default one being an n-gram model.

It also provides an `InterpolatedLanguageModel`<sup>9</sup> (ILM) which allows to specify multiple LMs and weights and interpolate the probabilities for a given n-gram from all models' probabilities ( $p = w_1 * p_1 + w_2 * p_2 + \dots$  where  $w_n$  are the weights ( $\sum_{i=1}^n (w_i) = 1$ ) and  $p_i$  is the probability for a given n-gram in  $LM_i$ ).

The purpose of the ILM in our approach is to factor in the importance of keywords. These keywords have to be supplied in the form of an n-gram language model. For this we extract text content from the lecture material, preprocess it and create an n-gram LM from the resulting corpus. Sphinx 4 is then a) run with a generic English n-gram LM only and b) with the ILM configured to use the generic English LM and the keyword language model in a 50/50 weighting. Finally the two resulting transcriptions are compared with a selection of metrics.

As an example, the 1-gram *sex* has a probability of 2.82% in the keyword model of `psy-14`, but a probability of 0.012% in the generic English LM.<sup>10</sup> When applying 50/50 interpolation, the result is  $2.82\% * 0.5 + 0.012\% * 0.5 = 1.416\%$ , which is an increase by the factor  $\sim 117$  over the generic probability.

### 4.1 Sphinx 4

### 4.2 Implementation

The pipeline is implemented with a collection of standalone command line tools and a set of Bash and Python scripts<sup>11</sup>.

The tasks are the following, in chronological order:

#### 1. Prepare the input

- The audio file is converted into Sphinx 4 compatible format (16khz, 16bit mono little-endian).

---

<sup>7</sup>Homepage: <http://cmusphinx.sourceforge.net/wiki/sphinx4:webhome>

<sup>8</sup>Overview: <http://cmusphinx.sourceforge.net/doc/sphinx4/edu/cmu/sphinx/linguist/language/ng4/LanguageModel.html>

<sup>9</sup>Javadoc: <http://cmusphinx.sourceforge.net/doc/sphinx4/edu/cmu/sphinx/linguist/language/ng4/InterpolatedLanguageModel.html>

<sup>10</sup>(*CMUSphinx US English Generic Language Model (cmusphinx-5.0-en-us.lm)*, 2015)

<sup>11</sup>The source code is available here: <https://github.com/jonathanewerner/bachelor/tree/master/bin>

- A testcase folder with a given shortname (e.g. **psy-15**) is created in the **results**-directory<sup>12</sup> of the source code repository.
- The reference transcript, the material (PDF format is required) and the converted audio file are moved into a **resources** subfolder of the testcase folder.

## 2. Create a keyword LM from lecture material

- `pdftohtml -i -xml` is applied on the given material PDF. The XML output representation is input to `pdfreflow`<sup>13</sup>. Compared to the tool `pdftotext` the combination of these 2 tools preserves paragraphs correctly, whereas `pdftotext` represents each line break in the input PDF as a new paragraph in the output text file. This is a significant disadvantage for the LM creation step, as a newline in the input file there has the semantic “end of sentence” – so that a sentence split into 4 lines by `pdftotext` would count as 4 sentences in the LM.
- The HTML output from `pdfreflow` is filtered by taking only relevant HTML-tags such as `<p>`'s (paragraphs) and `<blockquote>`'s, further improving the content-to-noise ratio.
- The resulting text is then preprocessed for optimal compatibility with the LM creation tool by removing punctuation and superfluous whitespace<sup>14</sup>.
- The resulting corpus is input to `estimate-ngram`, a LM creation tool from the MIT Language Modeling Toolkit<sup>15</sup> (MITLMT).

For clarification intermediate results from this step follow as an example. They are from the test case **psy-5**<sup>16</sup>. Figure 3 shows an example slide.

When using `pdftotext` the result looks like this for the given slide:

```
Piaget's Theory of
Cognitive Development
• Piaget believed that "children are active
thinkers, constantly trying to construct more
advanced understandings of the world"
• Little scientists
• These "understandings" are in the form of
structures he called schemas
```

Notice how each newline in the slide maps to a newline in the output. When using the combination of `pdftohtml` and `pdfreflow` the result looks like this:

<sup>12</sup><https://github.com/jonathanewerner/bachelor/tree/master/results>

<sup>13</sup>`pdftohtml` and `pdfreflow` are open source linux command line utilities

<sup>14</sup>I use a combination of command line text processing (`sed`) and a perl script from Stephen Marquard here.

<sup>15</sup><https://code.google.com/p/mitlm/wiki/EstimateNgram>

<sup>16</sup>“Introduction to Psychology, 5 - What Is It Like to Be a Baby: The Development of Thought”

## Piaget's Theory of Cognitive Development

- Piaget believed that "children are active thinkers, constantly trying to construct more advanced understandings of the world"
- Little scientists
- These "understandings" are in the form of structures he called *schemas*

Figure 3: Slide from lecture psy-5

```
<p class="p9">Piaget's Theory of  
Cognitive Development </p>  
<p class="p10">• Piaget believed that "children are active  
thinkers, constantly trying to construct more  
advanced understandings of the world" </p>  
<blockquote class="b9">• Little scientists </blockquote>  
<p class="p10">• These "understandings" are in the form of  
structures he called <i>schemas</i> </p>
```

Notice how a paragraph is captured in a `<p>`-tag. This allows extracting a sentence as one line in the corpus. After applying the preprocessing described above the final corpus for the slide looks like this (where “..” marks a line continuation):

```
piagets theory of cognitive development  
piaget believed that children are active thinkers constantly  
.. trying to construct more advanced understandings of the world  
little scientists  
these understandings are in the form of structures he called schemas
```

### 3. Convert transcript to reference corpus

The transcript from Open Yale is supplied as HTML. We apply processing steps to transform it to a corpus ready to be consumed by the WER analysis tool (no punctuation, all lowercase). As these are specific to just the format chosen by Open Yale Courses, the details are omitted, as they have no general use.

#### 4. Run Sphinx 4 in baseline and interpolated mode

`bin/sphinx-interpolated.py`<sup>17</sup> supplies a wrapper for interfacing with Sphinx 4. The Java API of Sphinx 4 is exposed for command line usage by a JAR package which bundles the Sphinx 4 libraries and a small Main class. This class uses command line arguments supplied from `bin/sphinx-interpolated.py` to correctly configure Sphinx 4 and start the actual recognition.

Each testcase folder has a configuration file which specifies which models the test run should use:

```
{
  "acousticModelPath": "en-new/cmuspinx-en-us-5.2",
  "dictionaryPath": "en-new/cmudict-en-us.dict",

  "languageModelPath": "en-new/cmuspinx-5.0-en-us.lm",
  "keywordModelPath": "model.lm",
  "g2pModelPath": "en-new/en_us_nostress/model.fst.ser",

  "resultsFolder": "biomed-eng-1"
}
```

`bin/sphinx-interpolated.py` interprets the “global” models relative to the repository root-folder `models`, the `resultsFolder` relative to the root folder `results` and the `keywordModelPath` relative to the `resultsFolder`. It then supplies the absolute paths to the JAR. It also supplies absolute output file paths for the transcription result and transcription word timings results.

This setup ensures reproducible results, as the environment of a given testcase is exactly specified (as long as the same binaries and script versions are assumed).

`bin/sphinx-interpolated.py` can now be used to run the baseline or/and interpolated version.

#### 5. Analyze and compare the results

Finally the results from the two recognition runs are analyzed and compared by running `bin/hotword-analyze <testcase folder name>`. This performs two things: a) WER comparison and metrics generation and b) keyword visualization.

##### 5.1 WER comparison and metrics generation

This first calls `bin/wer.py`<sup>18</sup> on each run, which will calculate the WER and show a summary of substituted (SUB), inserted (INS) and deleted

<sup>17</sup><https://github.com/jonathanewerner/bachelor/blob/master/bin/sphinx-interpolated.py>

<sup>18</sup>`wer.py` has been adapted from <http://progfruits.blogspot.de/2014/02/word-error-rate-wer-and-word.html>

(DEL) words when comparing the reference (REF) to the hypothesis (HYP):

```

OP   | REF      | HYP
INS  | ****     | this
INS  | ****     | is
INS  | ****     | that
OK   | this     | this
OK   | is       | is
OK   | a        | a
OK   | course   | course
SUB  | a        | that
SUB  | version  | aversion
OK   | of       | of
OK   | which    | which
SUB  | i've     | i
OK   | taught   | taught
INS  | ****     | him
...
...
{'Sub': 1230, 'Ins': 674, 'WER': 0.316, 'Del': 324, 'Cor': 5492}

```

In a second step it compares the two WER result files with `bin/compare-wer.py`.

The result is an HTML file with a) shows a WER comparison table and b) various statistical measures which will be explored later. The table (shown in Figure 4) colors correctly recognized words as green and incorrect words as red. It also marks words that have been improved in the interpolated version with a green border and words that have been worsened with a red border.

## 5.2 Keyword visualization

Data from the reference and the recognition results is compiled into a format suitable for consumption by a visualisation module, which will be discussed in chapter 6.

All intermediate steps from the pipeline are represented as files in the testcase folder. Table 2 gives an overview of the files created by each pipeline step.

Table 2: File results of a testcase run

File	Description
<b>Step 1: Prepare the input</b>	
<code>resources/audio.mp3</code>	original audio
<code>resources/audio.wav</code>	converted audio

File	Description
resources/slides.pdf	lecture material
resources/transcript.html	lecture transcript
config.json	run configuration
<b>Step 2: Create a keyword LM</b>	
slides.corpus.txt	lecture material corpus
model.lm	keyword LM
<b>Step 3: Convert reference to corpus</b>	
reference.corpus.txt	reference transcription corpus
reference_wordcounts.json	reference transcription word counts <sup>19</sup>
<b>Step 4: Run Sphinx 4</b>	
sphinx_log_baseline.txt	Sphinx 4 logging output
sphinx_log_interpolated.txt	
sphinx_result_baseline.txt	Sphinx 4 transcription
sphinx_result_interpolated.txt	
sphinx_word_times_baseline.txt	Sphinx 4 word times
sphinx_word_times_interpolated.txt	
<b>Step 5.1: WER comparison / metrics generation</b>	
results.json	run metrics in json format <sup>20</sup>
wer_baseline.txt	WER table / metrics
wer_interpolated.txt	
wer_comparison.html	rich WER comparison + metrics
<b>Step 5.2: Keyword visualization</b>	
cloud_baseline.json	data representation for visualization
cloud_interpolated.json	

<sup>19</sup>They are needed for the visualization later.

<sup>20</sup>This eases parsability for aggregating multiple testcase results later.

Reference	baseline	interpolated
this	this	this
is	is	is
a	a	a
course	course	course
a	that	ab
version	aversion	version
of	of	of
which	which	which
i've	i	i
taught	taught	taught
almost	almost	almost
every	every	every
year	year	year
for	for	for
the	the	the
last	last	last
twenty	the	the
years	years	years
and	and	added
it	it	a

Figure 4: WER comparison



## 5 Analysis

I will now discuss how to evaluate the usefulness of the LM-Interpolation approach in light of the goal to improve recognition accuracy of interesting keywords.

### 5.1 Approaching a good metric

We want to find a metric that describes if and how much the interpolated version improves upon the baseline version. Comparing the generic WER of the two runs does not help to answer the question of how much our approach improves the accuracy of interesting keywords.

#### 5.1.1 Lecture-scoped WER excluding $top_X$ words

In the interpolated approach we have included the LM created from the lecture material. The basic question to ask when assessing the effectiveness of this approach is: how much better is the WER when *just looking at the words from the lecture material LM*? This is only a starting point however. The lecture material corpus includes a substantial amount of words that would not be classified as “interesting keywords”: filler words and very common words. One approach to sort them out is subtracting a set of top  $x$  most common words (“ $top_X$  words”) from this list. The resulting metric can then be parameterized on the given  $x$ . This is an idea that Marquard uses when he proposes the metric “Ranked Word Correct Rate” (RWCR-n):

“RWCR-n is defined as the Word Correct Rate for all words in the document which are not found in the first n words in a given general English word dictionary with words ranked from most to least frequent.” (Marquard, 2012, p. 71)

#### 5.1.2 Lemmas

When searching for a specific term the user is interested in the *lemma* for a given word: when he wants to find occurrences of *child* in the given lecture, occurrences of “children”, “child’s”, “children’s” etc. would also be relevant. This implies two things: 1) when looking at the “atomic” level of improvements and degradations it is more relevant to have lemmas as atoms and not words and 2) the exact matching (a hypothesis word is only “correct” if it exactly matches the reference word) of the WER algorithm should be “loosened” to also mark hypothesis words as correct if their lemmatized version matches the reference.

The same principle holds for the  $top_X$  words: we only want to capture words for which the *lemma* is not in the  $top_X$  words.

### 5.1.3 Proposed metric: KWER-x

We can distill these concerns into a definition of a metric called  $KWER-x$ , which expresses the “Keyword Error Rate”, where a keyword is defined as the lemma of a word occurring in a given lecture material corpus given that this lemma is not present in the  $top_x$  list of most common words of the given language.

The value of  $x$  has to be determined empirically: how many of the top words should be filtered out? There has to be a balance between not accidentally excluding keywords (i.e “sex” is in the in the  $top_{500}$  words) and filtering out enough filler words.

After experimenting with some values I went with  $x = 500$  for my measurements. It is hard to find a less ad hoc approach to determining the “best”  $x$ , as you have no “meta”-metric that assesses how well a given  $x$  captures the goal of accurately describing the detection accuracy of keywords; it necessarily is a “best guess”.  $x = 500$  was chosen by looking at the relationship from  $x$  to  $\Delta KWER$  (the improvement in KWER-500) as shown in figure 5. For values of  $x$  below  $\sim 200$  the growth of improvement is explained by the gradual removal of filler words from the set of keywords. Their recognition accuracy is not improved by our approach which is why they prevent the accuracy improvement of actual keywords to be visible. Above  $\sim 200$  this factor is ruled out and the chosen  $x$  value has only miniscule influence on the resulting KWER.

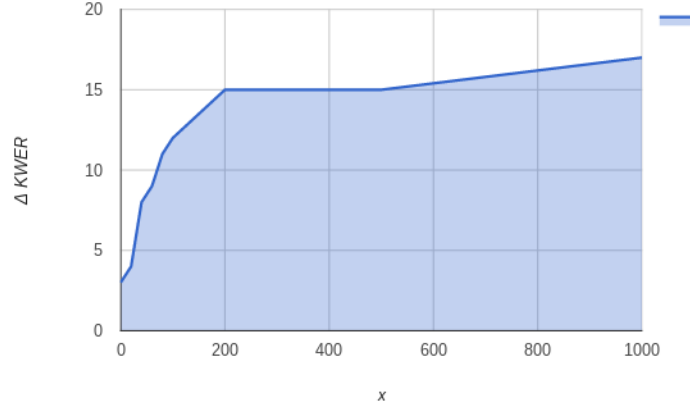


Figure 5: Relation of  $x$  to  $\Delta KWER$

This strategy, while depending on choosing an “ad hoc” value, was sufficient to validate our approach because it was possible to manually evaluate the metric’s “precision” by observing the resulting word sets. Another approach would have been to take the  $tf-idf$  (Term Frequency - Inverse Document Frequency) as a criterion for “keyword-ness” of words.  $tf-idf$  computes the “relevance” of a word in the context of a document by taking into account the occurrences of the

word in the document offset by the word’s frequency in a broader corpus. This way common words are rated lower although they occur frequently in the given document. This way there is no need for the arbitrary aspect of choosing an value of  $x$  and the negative side effect of accidentally excluding a keyword. On the other hand there would be the need to choose a threshold *tf-idf* score which would have to be met for inclusion into the keyword set.

#### 5.1.4 Secondary metrics

The following secondary or derived metrics are also evaluated:

- $W$ : Number of words
- $KW$ : Number of keywords
- $WER_{A|B}$ : WER of baseline (A) / interpolated version (B)
- $KWER_{A|B500}$ : KWER-500 of baseline (A) / interpolated version (B)
- $W_{worse|improved}$ : Proportion of worsened/improved words
- $KW_{worse|improved}$ : Proportion of worsened/improved keywords
- $W_{worse|improved}(K)$ : Proportion of worsened/improved words that are keywords
- $E$ :  $W_{improved}(K) - W_{worse}(K)$ : A percentage score for “effectiveness” of version B

Their use can be exemplarily demonstrated on the lecture **human-nature-8**: The lecture has 5342 words overall, of which 376 are keywords. When looking at the general WER, run A and B both have a score of 43%. This can be “explained” by looking at  $W_{worse|improved}$ , which is 4% each, meaning that 4% (223/5342) of the words have been improved from run A to B, but 4% (227/5342) them have been worsened, which sums up to 0% difference in WER.

Secondly, the KWER-500 of A is 48% (182/376 keywords) versus 32% for B (121/376 keywords). This improvement of 16% can analogously be explained by looking at  $KW_{worse|improved}$ : when looking at the 376 keywords, 2% (6/376) of them have been worsened while 18% (67/376) have been improved.  $18 - 2 = 16\%$  explains the improvement from 48% to 32%.

The last metric of  $W_{worse|improved}(K)$  looks at the overall worsened/improved words and informs about the proportion of words that were keywords. As mentioned,  $W_{worse}$  is 4% (223 of the overall 5342 words have been worsened). What is the proportion of keywords in this number? Analogously, what is the proportion of keywords when looking at the overall improved words? This metric is key in identifying the *effectiveness* (E) of our approach: the  $W_{improved}(K)$  value answers the question how well our approach is targeted towards improving the words we are interested in, the  $W_{worse}(K)$  value answers the question how big the “side effect” of worsening keywords is. In the example,  $W_{worse}(K)$  is 3% (6/227) and  $W_{improved}(K)$  is 30% (67/223). This is great: of the 227 overall worsened words only 6 were relevant given our goals. In essence, we

can interpret  $W_{improved}(K) - W_{worse}(K)$  as an **effectiveness score**, the same way we interpret the difference between  $W/KW_{improved}$  and  $W/KW_{worse}$  as “singular” metrics (WER and KWER respectively). We can say that our example had an effectiveness of  $30 - 3 = 27\%$ . An effectiveness of 100% would mean that *all* words that were improved had been keywords and *none* of the worsened words would have been keywords.

## 5.2 Results

The results for the test lectures described above (chapter 3) are as follows<sup>21</sup>:

Table 3: Results

Metric	1	2	3	4	5
W	5342	7233	7618	7142	7046
KW	376	715	974	607	518
$WER_A$	43%	30%	34%	37%	22%
$WER_B$	43%	30%	34%	37%	22%
$W_{improved}$	4%	4%	5%	5%	3%
$W_{worse}$	4%	4%	5%	5%	3%
$\Delta WER$ <sup>22</sup>	<b>0%</b>	<b>0%</b>	<b>0%</b>	<b>0%</b>	<b>0%</b>
$KWER_A$ <sup>23</sup>	48%	34%	33%	40%	32%
$KWER_B$	32%	18%	17%	19%	17%
$KW_{improved}$	18%	16%	17%	22%	16%
$KW_{worse}$	2%	1%	1%	0%	1%
$\Delta KWER$	<b>16%</b>	<b>15%</b>	<b>16%</b>	<b>22%</b>	<b>15%</b>
$W_{improved}(K)$	30%	39%	41%	40%	44%
$W_{worse}(K)$	3%	2%	2%	0%	2%
E	<b>27%</b>	<b>37%</b>	<b>39%</b>	<b>40%</b>	<b>42%</b>

The means are:

```
```{.table type="pipe" aligns="MMMMM" caption="Results" header="yes"}
Metric, Mean in %
$WER$, 33.2%
$KWER$, TODO
```
```

<sup>21</sup>Column 2-x represent the lectures, the numbers refer to the following lectures: 1: human-nature-8, 2: environmental-8, 3: psy-14, 4: psy-5, 5: biomed-eng-1.

<sup>22</sup> $\Delta$  refers to the improvement from version A to B in this context.

<sup>23</sup>KWER means KWER-500 for brevity if not noted otherwise.

The mean  $WER$  is 33.2%, the  $ma\Delta WER$  is 0.0%, for  $\Delta KWER$  it is 16.8%, for E it is 37%.

### 5.3 Interpretation

Several things are notable. The  $WER$  as well as  $W_{improved}$  and  $W_{worse}$  nearly don't change at all, the differences are only zero-digit absolute amounts. It is interesting that the results are so unambiguous in this respect; it is also unexpected that  $W_{improved}$  and  $W_{worse}$  always cancel each other out completely.

Assessing the  $\Delta KWER$  presents the challenge that no comparison is available that uses the exact same metric. However it is possible to “fuzzily” compare the performance by looking at metrics with the same basic idea.

The metric “RWCR-n” used by Marquard (2012) mentioned above is comparable, as it also uses the concept of filtering out the  $top_n$  most frequent words; it differs by not taking the lemmatized word version as their atomic unit. With that said, the average improvement in RWCR-10k over 13 lectures also taken from Open Yale Courses is 9.0%, while their average  $WER$  decreases by 0.8%.

Kawahara et al. (2008) use a metric called “Keyword Detection Rate”, where keywords are defined as content words (nouns and verbs excluding numbers and pronouns) that appear in the slide text. They then compute the f-measure (the “mean of the recall rate of keywords included in utterances and the precision of keywords detected in ASR results.”). They report improvements of 7.5% and 3.0% (for two test sets) in detection rate over the baseline accuracy, while the increase in  $WER$  is 2.2% and 1.3% over the baseline respectively<sup>24</sup>.

Miranda, Neto, & Black (2013) do not use a custom metric and report a  $WER$  improvement of 3.6%, when interpolating the LM with slide text contents; they achieve an improvement of 5.9%  $WER$  when using their proposed method of integrating the speech input with synchronized slide content.

While comparing  $WER$  performance has the discussed disadvantage of low relevance to the given evaluation goals and the non-standardized spectrum of custom metrics disallows an objective comparison of the different approaches, it yet gives an impression how our approach's performance relates to other work: the  $\Delta KWER$  of 16.8% seems like a good indicator that our approach is a viable solution for the goal of improving speech recognition for searchability and scannability. Additionally, the *effectiveness score* demonstrates that the approach nearly does not worsen keywords at all and 38.8% of the improved words are actually keywords.

In general, the uniform distribution of results over the various topic domains with their very different types of provided materials is also surprising. The results seem to suggest that the form and supposed “quality” of material (e.g. exercise sheet

---

<sup>24</sup>The mentioned results refer to the combined method of global and local adaptation.

versus lecture slides) does not correlate with the improvement in KWER. The initial assumption that lectures from the natural and formal sciences would be harder to recognize, based on the “naive” presumption that words like “adenosine 5'-triphosphate” would be impossible to recognize, seems to be invalid as well – apparently the combination of preprocessing, G2P and adapted weighting in the LM makes it possible to detect complicated technical terms like this as well.

### 5.3.1 Qualitative interpretation

While representing the performance of our approach with a set of metrics allows (at least internal) comparability of results, it can not convey a holistic impression of what would actually change for a user of a hypothetical speech media search/scan interface when using data generated with our approach versus the baseline approach.

This impression can be given by looking at the following detailed results of the run on the `biomed-eng-1` lecture.

#### Normal words improved

(of, 8) (that, 7) (the, 6) (or, 6) (and, 5) (a, 5) (in, 4) (to, 4) (is, 4) (it, 3) (course, 3) (into, 2) (an, 2) (your, 2) (from, 2) (than, 2) (one, 2) (those, 2) (this, 2) (talk, 2) (bridge, 1) (set, 1) (don't, 1) (some, 1) (are, 1) (annoying, 1) (really, 1) (again, 1) (there's, 1) (would, 1) (it's, 1) (there, 1) (how, 1) (version, 1) (we're, 1) (which, 1) (you, 1) (more, 1) (week, 1) (be, 1) (students, 1) (free, 1) (i've, 1) (with, 1) (by, 1) (distance, 1) (about, 1) (like, 1) (well, 1) (infectious, 1) (yale, 1) (very, 1) (where, 1) (engineers, 1)

#### Normal words worse:

(and, 15) (a, 10) (so, 9) (you, 8) (the, 8) (it, 7) (have, 6) (to, 5) (they're, 4) (of, 4) (that, 4) (are, 3) (can, 3) (be, 3) (we, 3) (on, 3) (at, 3) (in, 3) (how, 3) (online, 3) (that's, 3) (day, 2) (we'll, 2) (see, 2) (our, 2) (for, 2) (genes, 2) (could, 2) (it's, 2) (one, 2) (there, 2) (we're, 2) (but, 2) (is, 2) (as, 2) (if, 2) (two, 2) (principle, 2) (concept, 1) (office, 1) (years, 1) (london, 1) (go, 1) (just, 1) (had, 1) (easy, 1) (bridge, 1) (somebody, 1) (increased, 1) (very, 1) (familiar, 1) (safe, 1) (i've, 1) (every, 1) (they, 1) (now, 1) (organ, 1) (did, 1) (doctor's, 1) (because, 1) (old, 1) (some, 1) (really, 1) (what, 1) (said, 1) (lots, 1) (vessels, 1) (health, 1) (approach, 1) (patient, 1) (here, 1) (come, 1) (about, 1) (bow, 1) (or, 1) (cancer, 1) (point, 1) (period, 1) (long, 1) (apply, 1) (city, 1) (would, 1) (leading, 1) (three, 1) (been, 1) (their, 1) (way, 1) (was, 1) (tell, 1) (life, 1) (buy, 1) (posted, 1) (physician, 1) (these, 1) (say, 1) (us, 1) (patient's, 1) (thin, 1) (were, 1) (heart, 1) (an, 1) (heard, 1) (get, 1) (other, 1) (details, 1) (week, 1) (kinds, 1) (i, 1) (mechanical, 1)

#### KW improved:

(biomedical, 35) (dna, 7) (cells, 7) (engineering, 6) (biochemistry, 3) (cell, 3) (polymer, 2) (graph, 2) (gibbs, 2) (certain, 1) (energy, 1) (site, 1) (occur,

1) (plot, 1) (due, 1) (specifically, 1) (membrane, 1) (answer, 1) (has, 1)  
(higher, 1) (drugs, 1) (molecule, 1) (known, 1) (post, 1) (polymers, 1)  
(disease, 1) (order, 1)

**KW worse:**

(cells, 1) (maintain, 1) (beyond, 1) (genetic, 1) (due, 1)

You notice two things: a) the “exchange” of filler words from version A to B and vice versa, which is of no interest for searching and scanning, and b) interesting keywords that have substantial amounts of occurrences, that were not found before, while the amount of worsened KW is tiny. This is the important “qualitative”, high-level conclusion: the approach allows users to find technical terms in speech media which they weren’t able to find before and it works consistently over a broad spectrum of topics.

## 6 Visualization for scannability

We have shown that the LM-Interpolation approach is a viable tool for improving recognition accuracy of keywords on university lectures. The output data of our system are words with meta information: their associated timing and if they are keywords. How can this information be further used for helping a user with the task of scanning and searching through a given lecture? While it is technically possible to use the whole transcript and present the user an interface where the transcript is time-aligned with the lecture, that presentation is problematic as the *WER* of the transcript has not been improved and reading comprehension for texts with *WERs* above 30% is too low.

A better approach would be focusing the interface exclusively on the keywords in such a way that the provided timing meta information is transformed into a dense visual representation, thus making scanning possible. The user should be able to see the distribution of topics during the timeline of the lecture with *once glance*.

To this end I have developed a prototype implementation of such an interface. It features two views: the first one is a list of word timelines (Figure 6). A word timeline shows the distribution of occurrences of a given word over the time of the lecture. An occurrence is displayed as a dot; clicking the dot positions the corresponding lecture audio at the time the word occurrence is spoken. The timelines are vertically sorted by count of word occurrences. For analytical purposes the interface also shows the count of recognized occurrences in relation to the actual count of occurrences in the reference transcript, seen next to the word. It also overlays a graph which shows the *word density* at a given time point. The density function is calculated by performing a Gaussian Kernel Density Estimation (KDE) algorithm on the array of time positions for a given word. The red dots are local maxima of the function<sup>25</sup>, so that a word can have multiple maxima. The information about maxima is being used primarily in the second view.

The second view (Figure 7) is a *word cloud* with “semantic axes”, compared to regular word cloud visualizations where the axes don’t have meaning. The x-axis still is the time-axis of the lecture and the y-axis still is the keyword frequency. The central feature of this cloud is that it can show *multiple instances* of one keyword – one instance for each local maximum. The word instance is on the same point on the x-axis as the corresponding local maximum. The timeline for a word can be shown by clicking on it. The example shows “brain” in the activated state; the timeline shows up below the map. One can see the two instances of “brain” being horizontally aligned with the two local maxima below<sup>26</sup>. Clicking

---

<sup>25</sup>The local maxima are computed with the `scipy.signal.argrelextrema` function from the python `scipy` package and had some mildly surprising results, which were of no relevance for the interface prototyping task however.

<sup>26</sup>It is obvious here that the first local maximum for the word should rather be at about 30-40% of the word’s timeline, but that could be optimized.



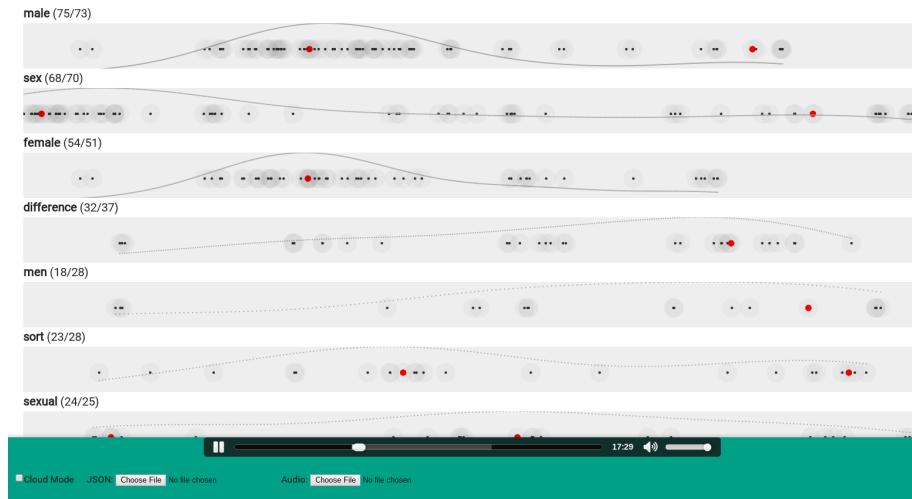


Figure 6: Word timelines

on the word also transports the audio to the position of the word next to the given local maximum. The font size of the word is computed by counting the word occurrences for which this maximum is the nearest. Additionally multiple instances of one keyword have the same color to further aid scanning by allowing the brain to pre-attentively process the representation.

This view allows an user to immediately scan the distribution of topics during the whole lecture. If particularly interested in the parts about the brain, he/she might click on “brain”, be immediately transported to the relevant audio position and additionally have a more in-depth view in the bottom timeline below the cloud, allowing him/her to intuitively grasp how long the relevant part might be, maybe skipping around by clicking on other instances of the word in the timeline.

You could imagine integrating this interface as a semi-transparent overlay view on a video player, for example on platforms like [lecture2go<sup>27</sup>](https://lecture2go.uni-hamburg.de), the lecture video streaming platform used by the University of Hamburg. When using a system that integrates many lectures in one database like this, it would also be possible to not only link to keyword instances in the same lecture but also on a broader scope, e.g the whole course or even other relevant courses & lectures. Another interesting extension point would be to integrate human intelligence by allowing to review/score the quality of keyword instances. This would allow filtering out false-positives and emphasize the keyword instances that students find helpful.

<sup>27</sup>[lecture2go.uni-hamburg.de](https://lecture2go.uni-hamburg.de)

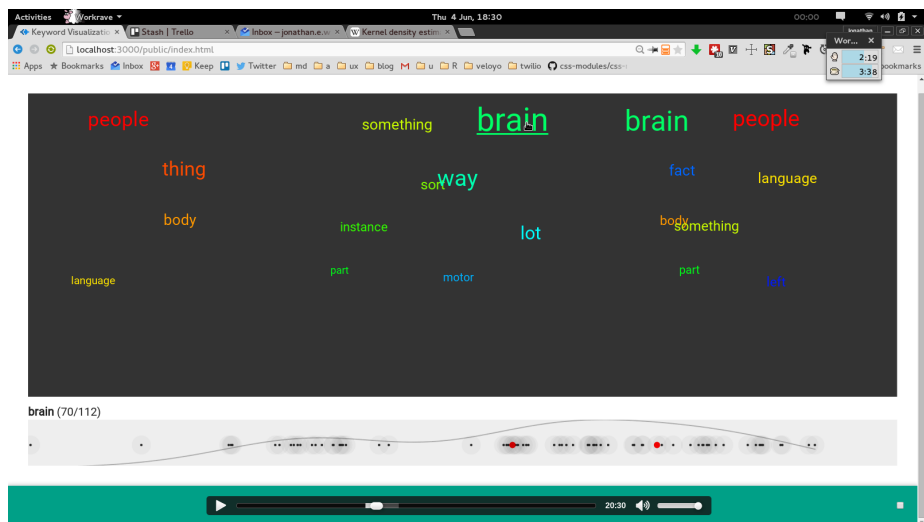


Figure 7: Word cloud

## 7 Summary, improvements

### References

- Cettolo, M., Brugnara, F., & Federico, M. (2004). Advances in the automatic transcription of lectures. In *Acoustics, speech, and signal processing, 2004. proceedings.(ICASSP'04). IEEE international conference on* (Vol. 1, pp. I-769). IEEE.
- CMU EN-US Pronouncing Dictionary (cmudict-en-us.dict)*. (2015). <https://github.com/cmusphinx/sphinx4/blob/master/sphinx4-data/src/main/resources/edu/cmu/sphinx/models/en-us/cmudict-en-us.dict>.
- CMUSphinx ARPA Language models*. (2015 (accessed 23.8.15)). <http://cmusphinx.sourceforge.net/wiki/sphinx4:standardgrammarformats>.
- CMUSphinx US English Generic Language Model (cmusphinx-5.0-en-us.lm)*. (2015). <http://sourceforge.net/projects/cmusphinx/files/Acoustic%20and%20Language%20Models/US%20English%20Generic%20Language%20Model/>.
- comp.speech Frequently Asked Questions*. (n.d.). <http://www.speech.cs.cmu.edu/comp.speech/Section6/Q6.1.html>.
- Cruttenden, A. (2014). *Gimson's pronunciation of English*. Routledge.
- English Wikipedia Arpabet article*. (2015 (accessed 22.8.15)). <https://en.wikipedia.org/wiki/Arpabet>.
- English Wikipedia Language Model article*. (2015 (accessed 23.8.15)). [https://en.wikipedia.org/wiki/Language\\_model](https://en.wikipedia.org/wiki/Language_model).
- Florian, C. (1996). The Blackwell encyclopedia of writing systems. *Oxford: Blackwell*.
- Glass, J., Hazen, T. J., Hetherington, L., & Wang, C. (2004). Analysis and processing of lecture audio data: Preliminary investigations. In *Proceedings of the workshop on interdisciplinary approaches to speech indexing and retrieval at hLT-nAACL 2004* (pp. 9–12). Association for Computational Linguistics.
- Hinton, G., Deng, L., Yu, D., Dahl, G. E., Mohamed, A.-r., Jaitly, N., Senior, A., et al. (2012). Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. *Signal Processing Magazine, IEEE*, 29(6), 82–97. IEEE.
- Kato, K., Nanjo, H., & Kawahara, T. (2000). Automatic transcription of lecture speech using topic-independent language modeling. In *Sixth international conference on spoken language processing*.
- Kawahara, T., Nemoto, Y., & Akita, Y. (2008). Automatic lecture transcription by exploiting presentation slide information for language model adaptation. In

*Acoustics, speech and signal processing, 2008. iCASSP 2008. IEEE international conference on* (pp. 4929–4932). IEEE.

Marquard, S. (2012). Improving searchability of automatically transcribed lectures through dynamic language modelling. University of Cape Town.

Miranda, J., Neto, J. P., & Black, A. W. (2013). Improving aSR by integrating lecture audio and slides. In *Acoustics, speech and signal processing (iCASSP), 2013 IEEE international conference on* (pp. 8131–8135). IEEE.

Munteanu, C., Penn, G., & Baecker, R. (2007). Web-based language modelling for automatic lecture transcription. In *INTERSPEECH* (pp. 2353–2356).

*Open Yale Courses Website*. (n.d.). <http://oyc.yale.edu/>.

Rabiner, L., & Juang, B.-H. (1993). Fundamentals of speech recognition. Prentice hall.

Stevenson, A., & Waite, M. (2011). *Concise Oxford English Dictionary: Book & CD-ROM Set*. Oxford University Press.

Yamazaki, H., Iwano, K., Shinoda, K., Furui, S., & Yokota, H. (2007). Dynamic language model adaptation using presentation slides for lecture speech recognition. *Proc. INTERSPEECH 2007*, 2349–2352.