

STATISTICAL MODELLING OF PRONUNCIATION: IT'S NOT THE MODEL, IT'S THE DATA

Florian Schiel¹

Andreas Kipp²

H. G. Tillmann²

¹Bavarian Archive for Speech Signals (BAS), Munich, Germany

²Department of Phonetics and Speech Communication

University of Munich, Germany

bas@phonetik.uni-muenchen.de

ABSTRACT

In this paper we describe a method to model pronunciation for ASR in the German VERBMOBIL task. Our findings suggest that a simple model, i.e. pronunciation variants modelled by SAM-PA units and weighted with a-posteriori probabilities, can be used successfully for ASR, if there is a *sufficient amount of reliably transcribed speech data* available. Manual segmentation and labelling of speech (especially spontaneous speech, as in the scheduling task of VERBMOBIL) is very expensive and time consuming and requires carefully trained experts and supervisors. Even with considerable effort it is not possible to produce broad phonetic transcripts for more than a small part of today customary speech databases. Therefore, as a first step in our approach we developed the fully automatic segmentation and labelling tool MAUS ('Munich AUtomatic Segmentation') for spontaneous German speech. The first part of our presentation will give a concise description of the MAUS method as well as an evaluation by comparing the results of MAUS with inter-labeller agreements of three expert phoneticians on the same data. The results show that MAUS operates within the range of human experts in terms of transcription while the timing information still lacks the quality of human segmenters. In a second step we used the MAUS system to segment and label 32h of speech in the 1996 VERBMOBIL acoustic evaluation to obtain more 320.000 transcribed words from the scheduling task. A simple counting, pruning and discounting technique (similar to that used for language modelling) is used to derive a probabilistic model of pronunciation. It provides a varying number of pronunciation variants per lexical entity together with the a-posteriori probability $P(V|W)$ that a variant V is uttered given the lexical entity W . A baseline system using HTK was set up for the 1996 VERBMOBIL evaluation task using monophones and a 'most likely' pronunciation dictionary (the 'most likelihood' was judged by a human expert NOT by empiric data). A second system with statistical modelling of pronunciation together with a proper re-training of the acoustic models showed significant better results on the same task in terms of word accuracy. From these findings we conclude that there's more to be done to achieve reliable and precisely labelled and segmented speech data than to investigate into very complex models which are usually prone to over-generalisation and lexical ambiguity.

1. AUTOMATIC SEGMENTATION WITH MAUS

1.1. Principle

The MAUS system was developed at the Bavarian Archive for Speech Signals (BAS) to facilitate the otherwise very time-consuming manual labelling and segmentation of speech corpora into phonetic units. Initially funded by the German government within the VERBMOBIL I project, MAUS is now further extended by BAS with the aim to automatically improve all BAS speech corpora by means of complete broad phonetic transcriptions and segmentations. The basic motivation for MAUS is the hypothesis that automatic speech recognition (ASR) of conversational speech as well as high quality 'concept-to-speech' systems will require huge amounts of carefully labelled and segmented speech data for their successful progress.

Traditionally a small part of a speech corpus is transcribed and segmented by hand to yield bootstrap data for ASR or basic units for concatenative speech synthesis (e.g. PSOLA). Examples for such corpora are the PhonDat I and II corpus (read speech) and the VERBMOBIL corpus (spontaneous speech). However, since these labellings and segmentations are done manually, the required time is about 800 times the duration of the utterance itself, e.g. to label and segment an utterance of 10 sec length a skilled phonetician spends about 2 h and 13 min at the computer. It is clear that with such an enormous effort it is impossible to annotate large corpora like the VERBMOBIL corpus with over 33 h of speech. On the other hand such large databases are needed urgently for empirical investigations on the phonological and lexical level.

Input to the MAUS system is the digitized speech wave and any kind of orthographic representation that reflects the chain of words in the utterance. Optionally there might be markers for non-speech events as well, but this is not essential for MAUS. The output of MAUS is a sequence of phonetic/phonemic symbols from the extended German SAM Phonetic Alphabet ([3]) together with the time position within the corresponding speech signal.

Example:

Input:
Speech Wave + 'bis morgen wiederhoeren'

Output:

```
MAU: 0 479 -1 <p:>  
MAU: 480 480 0 b  
MAU: 961 478 0 I
```

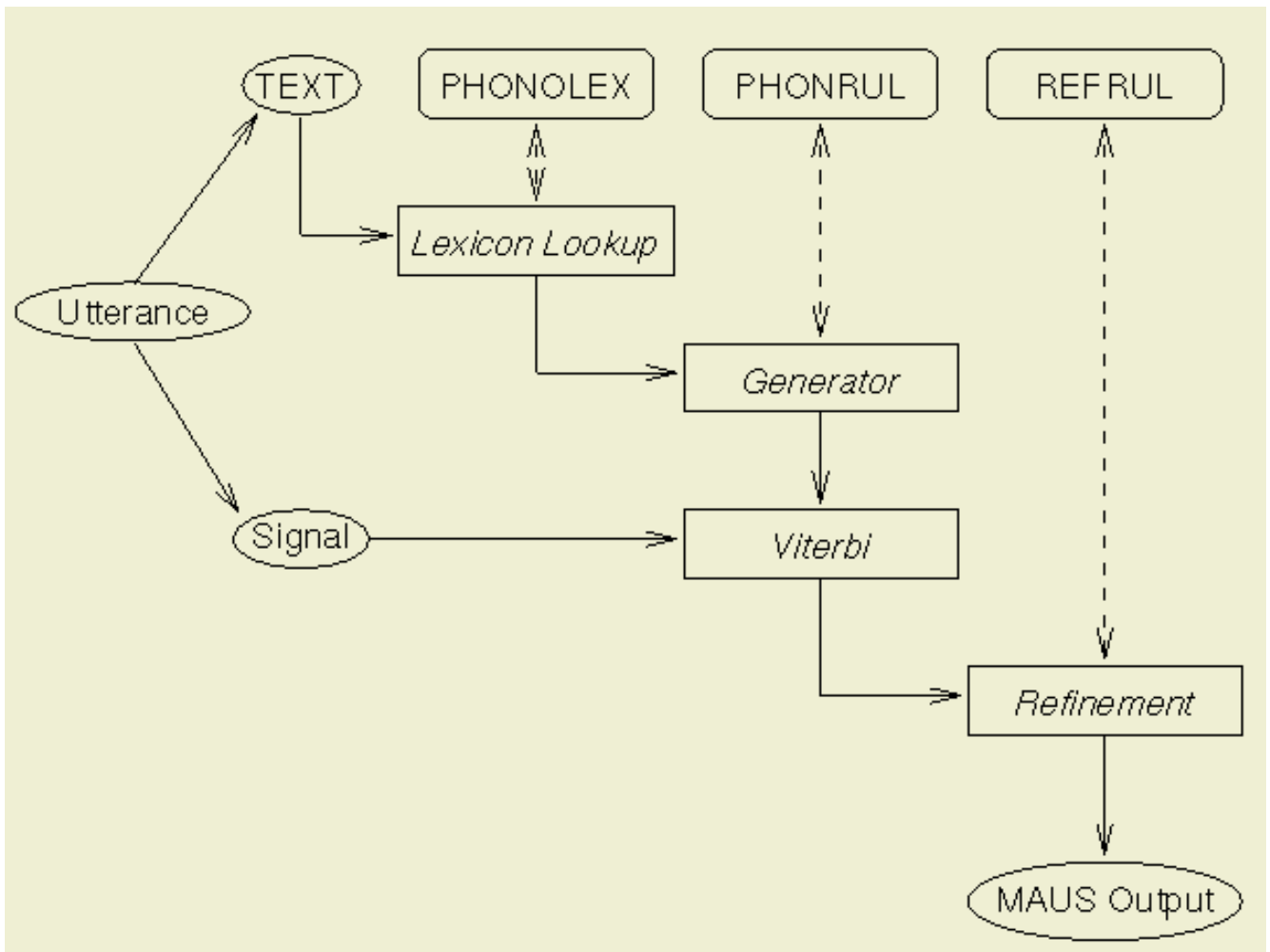


Figure 1. The MAUS system - block diagram

MAU: 1440 1758 0 s
 MAU: 2720 959 1 m
 MAU: 3680 799 1 0
 MAU: 4480 2399 1 6
 MAU: 6880 2079 1 N
 MAU: 8960 799 2 v
 MAU: 9760 959 2 i:
 MAU: 10720 479 2 d
 MAU: 11200 2239 2 6
 MAU: 13440 799 2 h
 MAU: 14240 639 2 2:
 MAU: 14880 1439 2 6
 MAU: 16320 1599 2 n
 MAU: 17920 1759 -1 <p:>

(The output is written as a tier in the new BAS Partitur format. 'MAU:' is a label to identify the MAUS tier; the first integer gives the start of the segment in samples counted from the beginning of the utterance; the second integer gives the length of the segment in samples; the third number is the word order and the final string is the labelling of the segment in extended German SAM-PA. See [13] for a detailed description of the BAS Partitur format)

MAUS is a three-staged system (see figure 1):

In a first step the orthographic string of the utterance

is looked up in a canonical pronunciation dictionary (e.g. PHONOLEX, see [14]) and processed into a Markov chain (represented as a directed acyclic graph) containing all possible alternative pronunciations using either a set of data driven microrules or using the phonetic expert system PHONRUL.

A microrule set describes possible alterations of the canonical pronunciation within the context of ± 1 segments together with the probability of such a variant. The microrules are automatically derived from manually segmented parts of the corpus. Hence, these rules are corpus dependent and contain no a priori knowledge about German pronunciation. Depending on the pruning factor (very seldom observations are discarded) and the size of the manually segmented data the microrule set consists of 500 to 2000 rules. In this paper we use a set of approx. 1200 rules derived from 72 manually segmented VERBMOBIL dialogs of The Kiel Corpus of spontaneous Speech ([15]). Details about this method can be found in [9].

The expert system PHONRUL consists of a rule set of over 6000 rules with unlimited context. The rules were compiled by an experienced phonetician on the basis of literature and generalised observations in manually transcribed data. There is no statistical information within this rule set;

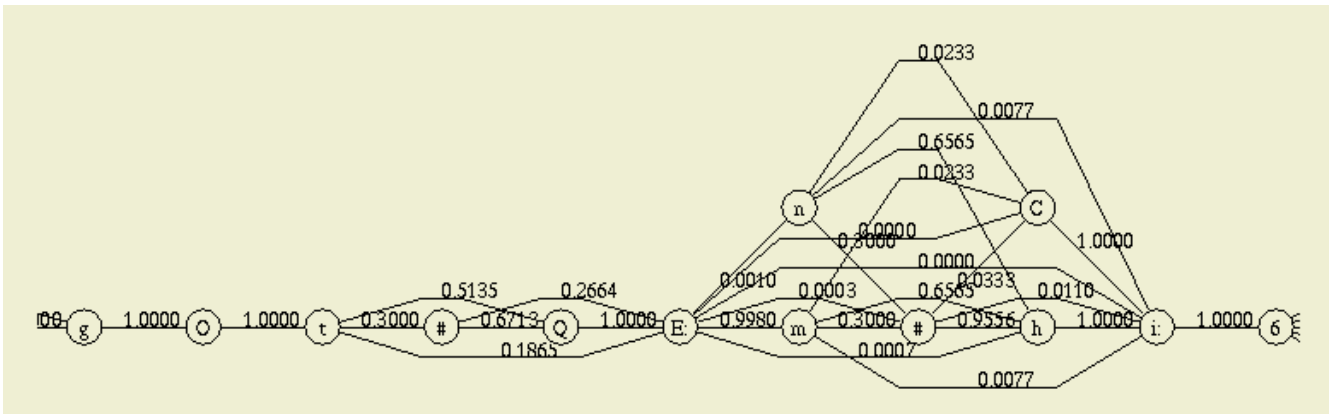


Figure 2. Acyclic graph of the utterance "Gott... ähm... hier..." with possible pronunciation variants

all rules are treated with equal probability. PHONRUL is therefore a generic model and should be considered independent of the analysed speech corpus. A more detailed description of PHONRUL can be found in [8].

The second stage of MAUS is a standard HMM Viterbi alignment where the search space is constrained by the directed acyclic graph from the first stage (see figure 2 for an example). Currently we use the HTK 2.0 as the aligner ([12]) with the following preprocessing: 12 MFCCs + log Energy, Delta, Delta-delta every 10 msec. Models are left-to-right, 3 to 5 states and 5 mixtures per state. No tying of parameters was applied to keep the model as sharp as possible. The models were trained to manually segmented speech only (no embedded re-estimation).

The outcome of the alignment is a transcript and a segmentation of 10 msec accuracy, which is quite broad. Therefore in a third stage REFINE the segmentation is refined by a rule-based system working on the speech wave as well as on other fine-grained features. However, the third stage cannot alter the transcript itself, only the individual segment boundaries.

The general drawback of the MAUS approach is, of course, that MAUS cannot detect variants that are not 'foreseen' by the first stage of the process. However, we found that using the microrule method the number of distinct rules converges after a relatively small sub-portion of the whole corpus. This indicates that the number of non-canonical pronunciations occurring in a certain domain such as the VERBMOBIL corpus is in fact limited and therefore treatable by a limited number of rules.

1.2. Evaluation

The output of MAUS can be separated into two different classes: the transcript (the chain of symbols) and the corresponding segmental information (begin and end of each segment).

Unlike in an ASR task the evaluation of a phonetic/phonemic segmentation of arbitrary utterances has a great disadvantage: there is no reference. Even very experienced phoneticians will not produce the same segmentation, not even the same transcript on the same speech wave.

We tried to circumvent this general problem by first comparing the results of three experienced human transcribers on the same corpus with each other to get a feeling for what is possible and set an upper limit for MAUS. We used standard Dynamic Programming techniques as used in ASR evaluations (e.g. [12]) to calculate the inter-labeller agree-

| | felix | marion | microrule | PHONRUL |
|--------|-------|--------|-----------|---------|
| dani | 82.6 | 78.8 | 80.2 | 76.7 |
| felix | - | 79.9 | 80.3 | 77.2 |
| marion | - | - | 74.9 | 72.5 |

Table 1. Comparison between 3 manual segmentations (dani, felix, marion), an automatic segmentation with the statistical microrule set and an automatic segmentation with the pronunciation model PHONRUL.

ment between different transcripts (see table 1). We found that the coverage of the three human transcribers ranges from 78.8% to 82.6% (on the basis of approx. 5000 segments). We then calculated the accuracy for the MAUS output with regard to each set of human results and found values ranging from 74.9% to 80.3% using the microrule method and 72.5% to 77.2% using PHONRUL. Not surprisingly, the worst and best coverage were correlated in all three experiments. This means that if we set the upper limit to the best match within human transcription results (82.6%) and compare this to the average agreement of MAUS with these two human transcribers, we'll end up with a relative performance of 97.2% for MAUS. (Note that this relative performance measure might be higher than 100% at some distant point in the future!)

For a more detailed discussion about the problem of evaluation as well as a more accurate analysis of the MAUS output (applied to read speech) please refer to [1].

In terms of accuracy of segment boundaries the comparison between manual segmentations shows a high agreement: on average 93% of all corresponding segment boundaries deviate less than 20msec from each other. The average percentage of corresponding segment boundaries in a MAUS versus a manual segmentation is only 84%. This yields a relative performance of 90.3%. We hope that a further improvement of the third stage of MAUS will increase these already encouraging results.

2. PRONUNCIATION MODEL FOR ASR

Aside from the many other uses of the MAUS output for this paper we'll show how to derive a simple but effective probabilistic pronunciation model for ASR from the data. There are two obvious ways to use the MAUS results for this purpose:

- use direct statistics of the observed variants

- use generalised statistics in form of microrules

In the following we will discuss both approaches.

2.1. Direct Statistics

Since in the MAUS output each segment is assigned to a word reference level (Partitur Format, see [13]), it is quite easy to derive all observed pronunciation variants from a corpus and collect them in a PHONOLEX ([14]) style dictionary. The analysis of the training set of the 1996 VERB-MOBIL evaluation (volumes 1-5,7,12) led to a collection of approx. 230.000 observations.

The following shows a random excerpt of the resulting dictionary:

```

terminlich
adj
t E 6 m i: n l I C
t E 6 m i: n I C      3
t @ m i: l I C 3
t E 6 m i: n l I C    10
t E 6 m i: l I C      1
t @ m i: n l I C      7
&
...
Karfreitag
nou
k a: 6 f r a I t a: k
k a: 6 f r a I t a: k    15
k a: 6 f r a I t a x    3
&
...
weil
par
v a I l
v a l    11
v a I    108
v a I l  207
&
...
siebenundzwanzigsten
adj
z i: b @ n U n t t s v a n t s I C s t @ n
z i: b @ n U n s v a n t s I s t @ n    1
z i: b m U n s v a n t s I k s t n      2
z i: b m U n s v a n s I C s t n        1
z i: b m U n s v a n t s I C s t @ n    1
z i: m U n s v a n t s s t @ n          1
z i: m U n s v a n t s s n              1
... (remaining 48 variants deleted)
&
...
Namen
nou
n a: m @ n
n a: m    30
n a: m @ n    15
&
...
Essen
nou
Q E s @ n
@ s n    2
E s n    16
E s @ n  6
s n      3
E s      1
Q E s @ n    7
Q E s      1
Q E s n    21
&

```

The above modified PHONOLEX format is defined as follows:

```

<orthography>
<comma separated list of
linguistic classes>
<canonical pronunciation>
<empiric pronunciation> <count>
&
...

```

Obviously many of the observations are not frequent enough for a statistical parameterisation. Therefore we prune the baseline dictionary in the following way:

- Observations with a total count of less than N per lexical item are discarded.
- From the remaining observations for each lexical word L the a-posteriori probabilities $P(V|L)$ that the variant V was observed are calculated. All variants that have less than $M\%$ of the total probability mass are discarded.
- The remaining variants are re-normalised to a total probability mass of 1.0.

Applied to the above example this yields the following more compact statistics (pruning parameters: $N = 20$, $M = 10$):

```

terminlich      0.434783
t E 6 m i: n l I C
terminlich      0.130435
t E 6 m i: n I C
terminlich      0.304348
t @ m i: n l I C
terminlich      0.130435
t @ m i: l I
Karfreitag      1.000000
k a: 6 f r a I t a: k
weil            0.342857
v a I
weil            0.657143
v a I l
siebenundzwanzigsten 0.509091
z i: b m U n s v a n t s I s t n
siebenundzwanzigsten 0.490909
z i: m U n s v a n t s I s t n
Namen          0.333333
n a: m @ n
Namen          0.666667
n a: m
Essen          0.320000
E s n
Essen          0.420000
Q E s n
Essen          0.120000
E s @ n
Essen          0.140000
Q E s @ n

```

where the second column contains the a-posteriori probabilities. This form can be directly used in a standard ASR system with multi pronunciation dictionary like HTK (version 2.1).

2.2. Generalised Statistics

The usage of direct statistics like in the previous section has the disadvantage that because of lack of data most of the words will be modelled by only one variant, which in many cases will be the canonical pronunciation. An easy way to generalise to less frequent words (or unseen words)

is to use not the statistics of the variants itself but the underlying rules that were applied during the segmentation process of MAUS. Note that this has nothing to do with the statistical weights of the microrules mentioned earlier in this paper; it's the number of appliances of these rules that counts. Since there is formally no distinction between microrules for segmentation in MAUS and probabilistic rules for recognition, we can use the same format and formalism for this approach as in MAUS. The step-by-step procedure is as follows:

- Derive a set of statistical microrules from a subset of manually segmented data or use the rule set PHONRUL as *labelling rule set* (see section 1.).
- Use the *labelling rule set* to label and segment the training corpus and count all appliances of each rule forming the statistics of the *recognition rule set*.
- Apply the *recognition rule set* during the ASR search to all intra-word and inter-word phoneme strings to create statistically weighted alternate paths in the search space

Note that the *recognition rule set* will very likely be a subset of the *labelling rule set*, if we use PHONRUL as the *labelling rule set*.

This approach has the advantage that the statistics are more compact, independent of the dictionary used for recognition (which for sure will contain words that were never seen in the training set) and generalise knowledge about pronunciation to unseen cases. However, the last point may be a source of uncertainty, since it cannot be foreseen whether the generalisation is valid to all cases where the context matches. We cannot be sure that the context we are using is sufficient to justify the usage of a certain rule in all places where this context occurs.

3. EXPERIMENTS

There have been several attempts to incorporate knowledge about pronunciation into standard methods for ASR. Most of them (with some exceptions, e.g. [2]) didn't yield any improvements. The argument was that the advantage of a better modelling on the lexical level is eaten up by the fact that the search space and/or the dictionary ambivalence increases. However, most of the literature did not take into account reliable statistics (because they were simply not available) and used acoustic models that were trained using canonical or most likely pronunciations. Our hypothesis is that an increase in recognition performance can only be achieved if the following two conditions are satisfied:

1. A reliable statistical model for pronunciation (which very likely has to be adapted to the specific task).
2. Acoustical models that match the modelling on the lexical level in terms of discriminative power.

We conducted several experiments with a standard HTK recogniser ([12]) for the 1996 VERBMOBIL evaluation task. In this paper we will report about experiments using the direct statistics approach and some preliminary results using recognition rule sets.

As a reference system a standard recogniser of HTK 2.0 with the following properties was designed for the experiment:

The speech signal is mean subtracted, emphasized and preprocessed into 12 MFCCs + log Energy, Delta, Delta-delta every 10 msec. Training and test sets are defined in the 1996 VERBMOBIL evaluation task ([7], 'Kuer', test

corpus: 6555 words). The canonical dictionary contains 840 different entries. The language model is a simple bigram calculated exclusively from the training set. The acoustic models are monophone left-to-right HMMs with 3-5 states containing a variable number of mixtures without tying. We use 46 models from the extended German SAM-PA including one model for silence and one model for non-speech events.

We trained and tested the recogniser with the same amount of data in two different fashions:

- *Baseline System*

Standard bootstrapping to manually labelled data (1h40) and iterative embedded re-estimation (segmental-k-means) using 30h of speech until the performance on the independent test set converged. The re-estimation process used a canonical pronunciation dictionary with one pronunciation per lexical entry. The system was tested with the same canonical dictionary.

- *MAUS System*

This system was bootstrapped to one third of the training corpus (approx. 10h of speech) using the MAUS segmentation and then iteratively re-estimated (30h of speech) using not the canonical dictionary but the transcripts of the MAUS analysis (note that the segmental information of the MAUS analysis is NOT used for the re-estimation).

The system was tested with the probabilistic pronunciation model described in section 2.1. using the pruning parameters $N = 20$ and $M = 0\%$.

Figure 3 shows the performance of both systems during the training process. Note that the MAUS system starts with a much higher performance because it was bootstrapped to 10h of MAUS data (compared to 1h40min of manually labelled data for the baseline system). After training, the MAUS system converges on a significantly higher performance level of 66.35% compared to 63.44% of the baseline system.

To verify our hypothesis we did two 'cross-check' experiments, where we used

- the statistical pronunciation model together with the acoustic models of the baseline system
- the canonical pronunciation dictionary together with the MAUS trained acoustic models.

As we expected the performance in these 'cross-check' experiments was not improved (in fact the performance dropped significantly in case A).

We also conducted several variations of the experiment where we used only the pronunciation variant with the highest a-posteriori probability or normalised the a-posteriori probabilities $P(V|W)$ within one lexical entry W to $\max(P(V|W))$, but none of these experiments showed a significant improvement to the baseline system.

The latter was surprising because we expected that words with very many variants will be 'punished' during recognition because of the very small $P(V|W)$. A possible explanation for this effect might be that words with very many observed variants are more frequent than words with few variants and therefore being favoured by the bigram language model. The fact that both, language model and pronunciation probabilities, have to be scaled by approximately the same factor (10.0) is another clue for this.

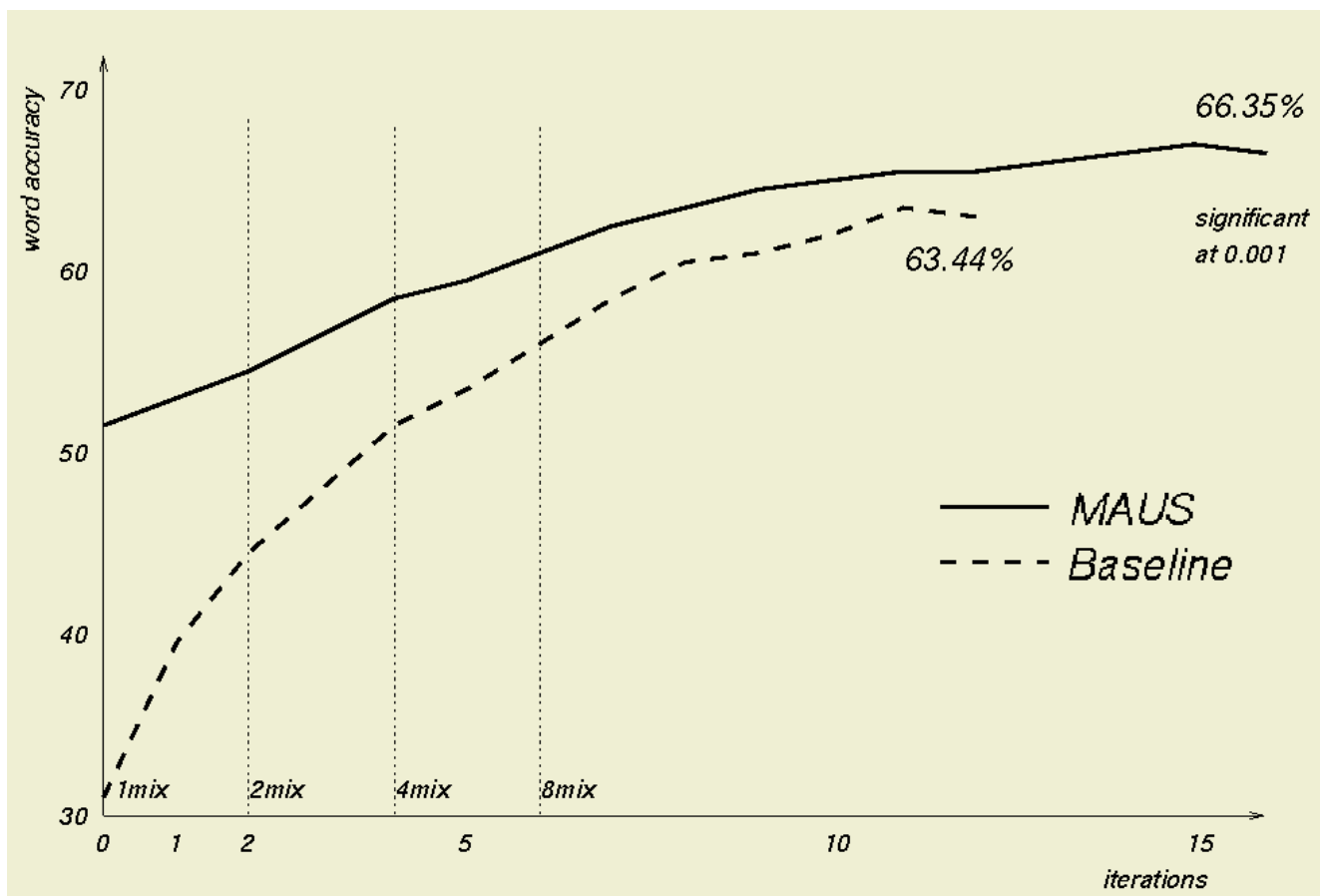


Figure 3. Performance of baseline system compared to the system trained with MAUS data and probabilistic pronunciation model

At the time of press we've only conducted some preliminary experiments using the 'generalized statistics' approach described in section 2.2.. We build a recognition rule set from the VERBMOBIL training set and applied this rule-set (together with the appropriate probabilities) to all words of the canonical dictionary. However, non of these experiments showed an improvement to the baseline system. At the moment our explanation is that the rule set with its limited context of $+/- 1$ is not able to generalize to other phoneme strings correctly. The obvious solution would be to extend the context to values higher than 1, but then we run into typical data scarcity problems.

4. CONCLUSION

The MAUS system can be used effectively to fully automatically label and segment read and spontaneous speech corpora into broad phonetic alphabets. The MAUS system enables us for the first time to derive statistical models on different processing levels (acoustic, phonetic, lexical) on the basis of very large databases.

Using a very simple statistical model based on reliable labelling data showed a significant improvement to a standard ASR baseline system. We regard this as a first hint to our hypothesis that

1. statistical modelling of pronunciation for ASR is feasible.

2. it's effective to use a simple model based on reliable statistical data.

The MAUS principle is not language dependent (however, the required resources are!). Therefore we strongly encourage colleagues in other countries to adopt the MAUS principle for their specific languages and produce similar resources as are currently produced at the Bavarian Archive for Speech Signals (BAS, [5]) for the German language.

REFERENCES

- [1] A. Kipp, M.-B. Wesenick (1996): Estimating the Quality of Phonetic Transcriptions and Segmentations of Speech Signals; in: Proceedings of the ICSLP 1996. Philadelphia, pp. 129-132.
- [2] T. Sloboda, A. Waibel (1996): Dictionary Learning for Spontaneous Speech; in: Proceedings of the ICSLP Philadelphia, pp. 2328.
- [3] Technical Report SAM, ESPRIT Project 2589, 1990. or www.phon.ucl.ac.uk/home/sampa/home.htm
- [4] M.B. Wesenick, F. Schiel (1994): Applying Speech Verification to a Large Data Base of German to obtain a Statistical Survey about Rules of Pronunciation, Proceedings of ICSLP 1994, pp. 279 - 282, Yokohama.
- [5] H.G. Tillmann, Chr. Draxler, K. Kotten, F. Schiel (1995): The Phonetic Goals of the new Bavarian

Archive for Speech Signals, Proceedings of the ICPHS 1995, pp. 4:550-553, Stockholm Sweden.

- [6] A. Kipp, M.-B. Wesenick, F. Schiel (1996): Automatic Detection and Segmentation of Pronunciation Variants in German Speech Corpora; in: Proceedings of the IC-SLP 1996. Philadelphia, pp. 106-109, Oct 1996.
- [7] J. Reinecke: Evaluierung der signalnahen Spracherkennung im Verbundprojekt VERBMOBIL (Herbst 1996), Verbomobil Memo 113, TU Braunschweig, Nov 1996.
- [8] M.-B. Wesenick (1996): Automatic Generation of German Pronunciation Variants; in: Proceedings of the ICSLP, Philadelphia, pp. 125-128, Oct 1996.
- [9] A. Kipp, M.-B. Wesenick, F. Schiel (1997): Pronunciation Modeling Applied to Automatic Segmentation of Spontaneous Speech; Proceedings of the Eurospeech 1997. Rhode, Greece, pp. 1023-1026.
- [10] M. D. Riley: A Statistical Model for Generating Pronunciation Networks, ICASSP 1991, pp. 737-740.
- [11] L. Lamel, G. Adda: On Designing Pronunciation Lexicons for Large Vocabulary Continuous Speech Recognition, ICSLP 1996 pp. 6-9.
- [12] S. Young et al. (1995/1996): The HTK Book; Cambridge University.
- [13] Partitur Format, www.phonetik.uni-muenchen.de/Bas/BasFormatsdeu.html
or
F. Schiel, S. Burger, A. Geumann, K. Weilhammer (1997): The Partitur Format at BAS. In: FIPKM 98, Institute of Phonetics, University of Munich, to appear.
- [14] PHONOLEX, www.phonetik.uni-muenchen.de/Bas/BasPHONOLEXeng.html
or
F. Schiel (1997): The Bavarian Archive for Speech Signals; in: FIPKM 1997, Institute of Phonetics, University of Munich, to appear.
- [15] IPDS (ed.): The Kiel Corpus of Spontaneous Speech; CDROM 1 + 2, University of Kiel, 1995.