

A Low-Cost Phonetic Transcription Method

Pablo Fetter

Udo Haiber

Peter Regel-Brietzmann

Daimler-Benz AG, Research and Technology, Wilhelm-Runge-Str. 11,
D-89081 Ulm, Germany
e-mail: fetter@dbag.ulm.daimlerbenz.com

ABSTRACT

In this paper our goal is to find the phonetic transcription of spoken utterances. We present a method which uses information extracted directly from the word-based search to compute the most likely phoneme sequence. Utterances are transcribed during recognition, so that the phonetic representation of the input is available after the search. Using this method, the computational cost of the word-based search remains almost unaltered, and the phonetic transcription is obtained almost for free.

1. INTRODUCTION

Phonetic transcriptions of spoken speech are necessary in a number of speech recognition issues. They can be used to incorporate new words into the system vocabulary, or to compute confidence measures, to name just two possibilities.

Most of the previous research on transcribing new words [1, 2, 6] and on using transcriptions to compute confidence [8] have the following characteristics in common:

- Two separate search processes are used (one for the word-based search and one for the phoneme-based search), and consequently,
- the computational cost of the phonetic transcription is high.

In the current study, we present a phonetic-transcription method which uses information extracted directly from the word-based search to compute the most likely phoneme sequence. Thus, a second search is avoided, and little additional computation is required for the phoneme transcription.

We will next describe our transcription approach, and its implementation in the Daimler-Benz large-vocabulary continuous-speech recognition system. We will then present some experimental results and general conclusions.

2. THE “CHEAP TRANSCRIPTION ALGORITHM”

The cheap transcription algorithm aims at obtaining an accurate phonetic transcription of the spoken utterance. In

the framework of a word recognizer, this method is “cheap” because it is embedded in the word-based search.

In the Daimler-Benz speech recognition system [4], the lexicon is implemented as a tree for efficiency reasons. This tree is compiled off-line and can be directly accessed and traversed during decoding. In the tree, every node is a (pointer to an) HMM state. During recognition, all active paths are expanded in every frame so that the static lexicon structure becomes a highly dynamic, growing tree (see Figure 1 for an example). Here, the search for the best path is performed by means of the Viterbi algorithm.

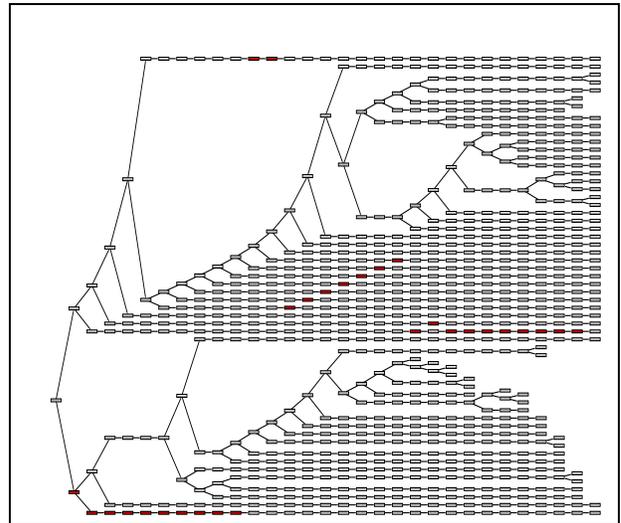


Figure 1: The search space of the first 30 frames (0.3 sec) in a Verbmobil sentence. Every column in the graph represents a frame and every node an HMM state. The nodes marked in black correspond to the best path in the current frame.

Keeping this picture in mind, it is intuitively clear that the phonetic information encoded in the HMMs is accessible during the lexical search. As mentioned before, every node in the lexicon tree represents one state of an HMM. And every HMM is again one context-dependent subword unit. So for phonetic transcription we store the best lexicon node in every

best transcription. This outcome is somehow intuitive: The best path per frame has a degree of inertia, that compensates the lack of constraints mentioned earlier. On the other hand, the best instantaneous phoneme (BMF) is not subject to any stabilizing force, and can therefore vary widely from frame to frame. Since the inferiority of the latter parameter was confirmed in further experiments, we abandon it at this point.

Best path per frame (BPF)	FER = 37.70 %
Best model per frame (BMF)	FER = 45.76 %

Table 1: Frame error rates (FER) for the best path per frame (BPF), and best model per frame (BMF).

3.2. Phoneme-Based Experiments

In order to obtain the actual phonetic transcription, the frame-based chains are smoothed and then collapsed. Collapsing consists of simply reducing successions of identical tokens to one token. For smoothing, a sliding weighted window was used: To decide on the identity of a phoneme, all symbols within an interval given by the window width are considered. The weighted frequency of every symbol inside the window is computed. The weight is determined by an exponential function of the distance of the symbol in question to the window center (symbols further apart contribute less). The symbol with the highest weighted frequency replaces the phoneme in the center of the window.

After the phoneme chains were smoothed and collapsed, we measured the Phoneme Error Rate (PER), which is calculated as usual:

$$\text{PER} = \frac{\text{Count(Ins)} + \text{Count(Sub)} + \text{Count(Del)}}{\text{Count(Ref)}} \quad (2)$$

where $\text{Count}(\dots)$ represents the number of insertions (Ins), substitutions (Sub), and deletions (Del) among the hypothesized phonemes, when aligned to the reference phonemes (Ref) in the test set³. A phoneme error rate slightly over 33% was obtained.

Next, we experimented with BPF chains containing a variable number of alternative phonemes per frame n , with $n \geq 1$. Using a frame-based phoneme language model, the most likely phoneme sequence ($n = 1$) can be extracted similarly to the way the best word sentence is extracted from a word graph [4]. This sequence is different to the one belonging to the best path in every frame (BPF), because the context is also considered for the decision on the identity of the current phoneme.

This chain was then smoothed and collapsed, and finally aligned to the spoken phoneme references. Figure 3 shows the results. A very significant improvement of the phoneme

³In order to calculate *frame* error rates, no alignment is necessary because the strings to be compared have the same length. For *phoneme* error rates, the alignment with the smallest Levenstein distance is used.

accuracy was obtained: The phoneme error rate was reduced from $\approx 33\%$ to $\approx 29\%$. It is interesting to note that as the number of alternatives presented to the language model grows, the importance of smoothing decreases: In this case, the language model does the smoothing.

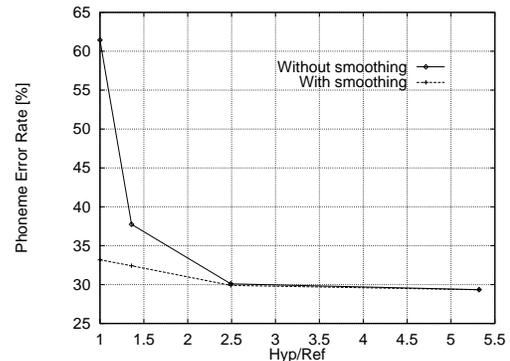


Figure 3: Phoneme error rates of the most likely phoneme sequence. This phoneme sequence is extracted from a graph structure containing a number of alternatives per spoken frame (*Hyp/Ref*).

3.3. Baseline System

For comparison purposes, the Daimler-Benz speech recognizer was also run in phonetic recognition mode. In this mode, phonemes are treated as if they were words. These “words” can be either context-independent (CI) or context-dependent (CD) HMM recognition units. In the former case, the lexicon consists of the 36 phonetic symbols used throughout this chapter, plus two models for silence and noise. In the latter case, we used approx. 1500 different phonemes in context. In addition, we also performed experiments with phoneme bigram backoff language models, which were used for constraining the search. The context-independent bigram language model was trained with the half of the Kiel Corpus of Spontaneous Speech not used for testing⁴. For the context-dependent bigram language model the whole VerbMobil corpus (excluding development and test sets) was used. Table 2 summarizes the characteristics of the language models. It is interesting to note that the language model of context dependent units (CD) is much more robust, even though its vocabulary is much larger. This is due to the strong contextual constraints imposed by context dependent units.

The results obtained for different system configurations are shown in the upper part of Table 3. As expected, the performance of the recognizer with context-dependent units (CD) is significantly better than with context-independent units (CI). Further, Table 3 clearly shows the importance of using a phoneme bigram language model. The language model

⁴Tests with other training sets were performed as well. The bigram trained with the actually spoken phoneme strings from the Kiel Corpus of Spontaneous Speech yielded the lowest test set perplexity and was therefore chosen for recognition experiments.

LM	No. of turns	No. of tokens	Vocab. size	PP
CI	184	20 K	38	13
CD	9.6 K	1 M	1.5 K	4

Table 2: Characteristics of the language models used for the baseline phoneme recognizers.

is specially advantageous for large vocabularies, as in the context-dependent case. The best performance obtained is $\approx 28\%$ phoneme error rate using context-dependent units and a bigram language model. This outcome is comparable to the ones reported by other Verbmobil partners [7].

System	+RTF	PER [%]
Baseline: CI, no LM	≈ 0.4	48.03
Baseline: CI, bigram	≈ 0.5	42.17
Baseline: CD, bigram	≈ 6	27.79
BPF	≈ 0	33.19
BPF, frame trigram	≈ 0.2	29.33

Table 3: Overhead real time factor (+RTF) and phoneme error rate (PER) for different transcription methods. The baseline systems are phoneme recognizers with different configurations. The cheap transcription algorithm is evaluated with and without a language model.

3.4. Comparison of Performance

Finally, we performed a comparative analysis of the computational resources needed by the different transcription methods presented above. For all methods tested, we measured the *elapsed time* needed to complete a recognition run. Then, we normalized them by the duration of the test set to obtain the *real time factor (RTF)*. We report here the overhead RTF (+RTF), which represents the additional time required to generate a phonetic transcription. All experiments were conducted on a Digital AlphaServer 8200 5/300 (SPECfp95 11.7, as reported in <http://open.specbench.org>). There, our word recognizer runs in 1.3 times real time (RTF ≈ 1.3).

Table 3 summarizes our measurements. The cheap transcription algorithm allows the input utterance to be transcribed with no (BPF) or very small computational overhead (BPF with frame trigram). In both cases, the phonetic transcriptions are more accurate and less expensive than those obtained using context-independent models with a phoneme recognizer. The lowest error rates are still obtained with a phoneme recognizer using context-dependent units and a phoneme bigram language model. Yet this performance gain is very costly (note that +RTF ≈ 6) due to the high confusability of the vocabulary, which makes pruning techniques ineffective. For comparison, we also tuned the recognizer parameters to allow more pruning. We could accelerate the context-dependent baseline recognizer from +RTF ≈ 6 to +RTF ≈ 3 with a minor performance degradation (PER = 30.97%). Still, the transcriptions provided by the cheap transcription algorithm are a better compromise

between performance and computational costs.

4. CONCLUSIONS

We presented a method we dubbed the “cheap transcription algorithm”, which can be used to transcribe utterances during recognition, so that the phonetic representation of the input is available as a by-product of the (word-based) search. Using this method, the computational cost of the word-based search remains almost unaltered. As a bonus, the method does not require any changes in the system lexicon, since the transcription is extracted directly from the models competing in the search. Further, the phonetic transcriptions obtained using the cheap transcription algorithm are almost as accurate as the ones obtained by a phoneme recognizer at a much higher cost.

ACKNOWLEDGMENTS

The authors gratefully thank Gunnar Evermann and Henrik Heine (University of Hamburg) for sharing their ideas and contributing to this work with interesting discussions. We also thank Matthias Pätzhold and Adrian Simpson (IPdS, University of Kiel) for their support with the Kiel Corpus of Spontaneous Speech.

REFERENCES

1. Ayman Asadi, Richard Schwartz, and John Makhoul. Automatic modeling for adding new words to a large-vocabulary continuous speech recognition system. In *Proc. ICASSP'91*, pages 305–308, 1991.
2. Katunobu Itou, Satoru Hayamizu, and Hozumi Tanaka. Detection of unknown words and automatic estimation of their transcriptions in continuous speech recognition. In *Proc. ICSLP'92*, Banff, Alberta, Canada, October 1992.
3. Klaus J. Kohler. Labelled Data Bank of Spoken Standard German: The Kiel Corpus of Read/Spontaneous Speech. In *Proc. ICSLP'96*, Philadelphia, USA, 1996.
4. Thomas Kuhn, Pablo Fetter, Alfred Kaltenmeier, and Peter Regel-Brietzmann. DP-Based Wordgraph Pruning. In *Proc. ICASSP'96*, volume 2, page 861, Atlanta, USA, 1996.
5. Jörg Reinecke. Evaluierung der signalnahen Spracherkennung im Verbundprojekt Verbmobil (Herbst 1996). Verbmobil Memo 113, Technische Universität Braunschweig, Germany, 1996. available at <http://sbvsrv.ifn.ing.tu-bs.de/eval.html>, in German.
6. Bernd Suhm and Monika Woszcyna. Detection and transcription of new words. In *Proc. EUROSPEECH'93*, pages 2179–2182, 1993.
7. Henrik Heine Uwe Jost and Gunnar Evermann. What is wrong with the lexicon—an attempt to model pronunciations probabilistically. In *Proc. EUROSPEECH'97*, Rhodos, Greece, 1997.
8. Sheryl R. Young. Detecting Misrecognitions and Out-Of-Vocabulary Words. In *Proc. ICASSP'94*, pages II-21–II-24, Adelaide, Australia, 1994.