



ACADEMIC  
PRESS

Available online at [www.sciencedirect.com](http://www.sciencedirect.com)

SCIENCE @ DIRECT®

Computer Speech and Language 17 (2003) 311–328

COMPUTER  
SPEECH AND  
LANGUAGE

[www.elsevier.com/locate/csl](http://www.elsevier.com/locate/csl)

# Acoustic model clustering based on syllable structure

Izhak Shafran, Mari Ostendorf \*

*Department of Electrical Engineering, University of Washington, Box 352500, Seattle, WA 98195-2500, USA*

Received 27 February 2001; received in revised form 4 September 2002; accepted 23 September 2002

---

## Abstract

Current speech recognition systems perform poorly on conversational speech as compared to read speech, arguably due to the large acoustic variability inherent in conversational speech. Our hypothesis is that there are systematic effects in local context, associated with syllabic structure, that are not being captured in the current acoustic models. Such variation may be modeled using a broader definition of context than in traditional systems which restrict context to be the neighboring phonemes. In this paper, we study the use of word- and syllable-level context conditioning in recognizing conversational speech. We describe a method to extend standard tree-based clustering to incorporate a large number of features, and we report results on the Switchboard task which indicate that syllable structure outperforms pentaphones and incurs less computational cost. It has been hypothesized that previous work in using syllable models for recognition of English was limited because of ignoring the phenomenon of resyllabification (change of syllable structure at word boundaries), but our analysis shows that accounting for resyllabification does not impact recognition performance.

© 2003 Elsevier Science Ltd. All rights reserved.

---

## 1. Introduction

Recognizing conversational speech has proved to be more challenging than read speech for automatic speech recognition (ASR) systems. For the best systems reporting results on the 1999 DARPA Broadcast News benchmark tests, word error rates on the spontaneous speech portion of the test set (14–16%) were nearly double those on the baseline condition of planned recordings (8–9%) (Pallett, Fiscuss, Garofolo, & Martin, 1999). Those sites that also participated in 2000 DARPA Conversational Speech benchmark tests, performed at error rates of roughly 30%.

---

\* Corresponding author. Fax: 1-206-543-3842.

*E-mail addresses:* [zak@u.washington.edu](mailto:zak@u.washington.edu) (I. Shafran), [mo@ee.washington.edu](mailto:mo@ee.washington.edu) (M. Ostendorf).

Though error rates continue to drop, the gap between conversational and news speech remains, with error rates for Broadcast news being half of what is reported for the Switchboard corpus (Le, 2002). The degradation in performance may be due to many factors such as channel effects, variability in speaking rate and dialect of speakers, less careful pronunciation, loosely structured language, and the presence of disfluencies. In 1996, (Weintraub, Taussig, Hunicke-Smith, & Snodgrass, 1996) demonstrated that a large part of the degradation is related to acoustic variation associated with speaking style. In this study, spontaneous speech was recorded and then the transcript of the speech was both read and acted by the same speakers to control for effects related to speaker, channel and language. While spontaneous speech was recognized with a word error rate of 52.6%, the acted and read versions were recognized at 37.4% and 28.8% error rates, respectively. The degradation with increasingly casual speaking style was observed across telephone-band and wide-band speech and under matched training and test conditions (Saraclar, Nock, & Khudanpur, 2000).

In many ASR systems, the acoustic variation of words are modeled at two levels – the pronunciation model which maps word sequences to phonemes, and the acoustic model which maps phoneme sequences to multivariate acoustic models. Work with simulated data which was produced using the acoustic models of speech, have pointed to pronunciation variability as a key problem in recognizing conversational speech (McAllister, Gillick, Scattone, & Newman, 1998). However, the work on pronunciation modeling in terms of phoneme-level substitutions, deletions and insertions has so far only yielded small performance gains (Byrne et al., 1997; Riley et al., 1999). In other work, (Saraclar et al., 2000) showed that modeling pronunciation at the state level and allowing the components of a Gaussian mixture model to be shared across alternate pronunciations is more beneficial than modeling pronunciation at phoneme level. Experiments by Hain and Woodland (2000) demonstrate an advantage in using phonetic context to directly influence the model sequence. Both these studies support the notion that there is a need to represent variation of a more gradient nature where higher-level context (beyond triphones) influences the acoustic models of the phonemes as well as the pronunciation of a word.

Conventionally, phone-level acoustic variation has been captured by conditioning the acoustic models for a phoneme on the context of neighboring phonemes in the hypothesized sequence. Typically, in large vocabulary ASR, phonemes with immediate neighbors (triphones) and possibly two neighbors (pentaphones) are used. Conditioning only on phonemic context does not capture the acoustic variation of conversational speech fully. Already, position of phoneme in word has been found to be useful in acoustic modeling. This conforms to observations about word-position effects in linguistic studies of different consonants with electropalatography (EPG) (Keating, Wright, & Zhang, 1999). The linguopalatal (tongue-palate) contact, which affects the strength and duration of the sound produced, differs significantly for word-initial vs. word-final consonants.

Our hypothesis is that, in English, syllable structure is also useful in modeling the variation not accounted for by phoneme context. Consider the phoneme “t” (in the context “iy t er”) in “beater”, “beat Ernest” and “return”. Even though it is the same triphone, the articulation of phone “t” in the three contexts is distinctly different – in the first it is flapped, in the second it is an unreleased closure and in the third it is a closure plus a release. These differences are closely related to syllable structure, and correspond to ambisyllabic, syllable-final, and syllable-initial contexts, respectively. The use of syllable structure is motivated in part by results from psychoacoustic studies, which argue for the syllable as a unit of perception, e.g., Massaro (1972); see

also Wu (1998) for a survey. Further support comes from corpus analyses. In a series of studies, Greenberg and Fosler found systematic variation with respect to the syllable constituent, namely onset, nucleus and coda (Greenberg, 1998; Fosler-Lussier, Greenberg, & Morgan, 1999; Greenberg & Fosler-Lussier, 2000). They used a subset of Switchboard corpus for conversational speech and a standard corpus called TIMIT for read speech, annotated by linguists at word, phone and syllable level. The perceived phones were compared with the lexical expansion of the spoken words. Both corpora show that the onset of a syllable maintains its canonical identity at most times (85–91%) regardless of the speaking style, and more so in the presence of consonant clusters. In general, the nucleus is prone to substitution by a wide range of vowels. The coda is less often realized in canonical form in conversational speech (63%) than in read speech (81%). The coda is prone to deletion, but the absence of a coda does not impact canonical realization of nucleus. These results together support the notion of syllable-initial strengthening, which has been observed as a more gradient phenomena in EPG studies that also suggest that the amount of strengthening may be equal to that in word-initial position (Keating et al., 1999). An analysis of errors made by state-of-the-art systems on recognizing conversational speech (Greenberg & Chang, 2000) suggests that accurate recognition of syllable onsets is more important for word recognition than syllable codas. While categorical phonetic changes can be accommodated by a larger phone inventory and a good pronunciation model, as in the TIMIT labeling conventions, phone substitutions and deletions fail to capture more gradient aspects of variation such as strength of a stop release. Thus, it is not surprising that state-based pronunciation models outperform phone-based models. In this work, we look at a complementary approach to state-based pronunciation modeling, which is acoustic model context conditioning on high-level contexts, specifically syllable and word structure.

One way to model syllable structure is to use syllable-sized units rather than phones. For small vocabulary tasks, a few researchers have successfully used the syllable as a unit for acoustic modeling (Jones, Downey, & Mason, 1997; Hamaker, Ganapathiraju, Picone, & Godfrey, 1998). Others (Lleida, Marino, Nadeu, & Salavedra, 1991; Marino, Nogueiras, & Bonafonte, 1997) split the syllable into demi-syllable units and used them for acoustic modeling. However, these two approaches lack the ability to effectively model syllables that are rarely or not seen in the training data. To overcome this deficiency, triphones were used in addition to frequent syllables in (Doddington, Corrada, & Wheatley, 1997). This approach in conjunction with improvement in temporal structure of the acoustic model gave a 2% absolute reduction in error rate (from 49% to 47%) on a conversational speech recognition task. A major part of this improvement came from modeling syllables in monosyllabic words separately from other instances of the same syllable. A disproportionate number of errors were found to be in words modeled by triphones rather than syllables, and may be due to poor sharing of parameters between the triphones and syllables (i.e., the triphones within infrequent syllables did not share parameters with those in frequent syllables).

The focus of our work is on learning contextual variation directly in the acoustic model using both word- and syllable-level information, since they seemed promising in both pronunciation models and previous acoustic studies mentioned above. In contrast to modeling the syllable explicitly as a unit, we use a tree-based clustering mechanism to allow sharing of parameters across all contexts for robust estimation. To tackle the problems that arise in using a large number of contextual features, we have extended the decision tree based clustering to use multiple stages of clustering.

The paper is organized as follows. In Section 2, we present a brief review of tree-based acoustic model clustering, followed by our use of syllable features in them, and then outline issues related to resyllabification. The details of the multiple stage clustering approach is presented in Section 3. Experimental results on the use of syllable structure, a resyllabification model, and multistage clustering are reported in Section 4, using the Switchboard corpus of conversational speech. Finally, Section 5 concludes and discusses future work.

## 2. Clustering with syllable features

Since we use decision trees extensively in this work, a brief review of tree-based clustering is provided, which is followed by a discussion of related work and description of our use of syllable features.

### 2.1. Tree-based clustering

For large vocabulary ASR systems, decision tree distribution clustering is used to map the large number of possible triphone (or pentaphone) contexts into a smaller set of distributions that can be robustly estimated (Young & Woodland, 1994). This technique is particularly attractive for parameter tying, as it allows mapping of any sub-word unit that is not seen in the training data to a cluster made up of acoustically similar units. Typically, a fixed HMM topology of 3–5 five states is associated with each phoneme, and separate trees are used for each state of a phoneme.

In training, all the context-specific observations associated with a particular phone state and observed in the training data are pooled at the root node of the tree. A set of predefined questions, typically about phonetic context (e.g., “Is the left phoneme a vowel?”), is used to define candidate binary partitions of a node in the tree. Assuming that all the data in a partition share a common Gaussian, the question corresponding to the partition that maximizes the likelihood of the data in a node is chosen as a candidate for the next split. From amongst these candidates, the node with the best likelihood gain from using 2 versus 1 Gaussian is split. The best partitions of the new clusters resulting from this split are added to the list of the candidate splits, and thus the tree is grown until some stopping criterion is met (e.g., limit on the number of leaves or terminal nodes).

More specifically, evaluating goodness of a particular split involves computing a generalized log likelihood ratio:

$$\log \left[ \frac{(\max_{\mu_L, \Sigma_L} p(\mathcal{X}_L | \mu_L, \Sigma_L)) (\max_{\mu_R, \Sigma_R} p(\mathcal{X}_R | \mu_R, \Sigma_R))}{\max_{\mu_P, \Sigma_P} p(\mathcal{X}_P | \mu_P, \Sigma_P)} \right],$$

where  $\mu_i$  and  $\Sigma_i$  are means and covariances, respectively; L, R and P subscripts indicate the left, right and parent nodes; and  $\mathcal{X}_i$  indicates a data subset, where  $\mathcal{X}_P = \mathcal{X}_L \cup \mathcal{X}_R$ . Efficient implementation of this ratio for the different candidate splits involves using sufficient statistics computed for each observed context combination (Kannan, Ostendorf, & Rohlicek, 1994), so tree design complexity depends on the number of observed contexts  $N$ . The computational complexity of growing a tree is  $O(QN)$ , where  $Q$  is the number of candidate questions (the sum of the number of partitions of a conditioning variable over all variables). For variables of cardinality  $M$ , the number of partitions can be as large as  $2^M$ , but typically the question set is restricted in some way

to reduce this factor to be linear in  $M$ . Note that the variables included in the context influence both the size of  $Q$  and  $N$ : each new conditioning factor adds linearly to  $Q$  and increases the number of possible triphone contexts by a factor of  $M^3$ , though typically the number of observed contexts is smaller. Hence, the cost of tree design increases polynomially with the cardinality of each new conditioning variable.

The tree is typically designed using a Gaussian state distribution assumption, for simplicity, then later more complex Gaussian mixture distributions are estimated to model the data in the leaves using the Estimation-Maximization (EM) algorithm (Rabiner & Juang, 1993). In building a word model for decoding, a particular context-specific phoneme is dropped down the tree and is guided by the linguistic questions at the branches. As illustrated in Fig. 1, the distribution associated with the leaf that it lands in is associated with a state in the context-dependent phone model.

In most ASR systems, decision tree questions are based on hand-specified phonetic classes (e.g., grouped by manner and/or place of articulation) on the neighboring phonemes. By incorporating a symbolic description of phonemes in the lexicon such as stress, position of the phone in the word, and position of the phone in the syllable, it is possible to capture new phenomena with decision tree clustering, such as a tendency to reduce unstressed vowels and to more strongly release a stop consonant in word onset position. The phonemes in the training data are marked with these symbols. Then, the tagged models are estimated and clustered, just as for triphones, except that the decision tree must choose among questions about these tags as well as those defined in terms of phonetic context. Clustering with information other than phonetic neighbors is sometimes referred to as *tagged clustering*.

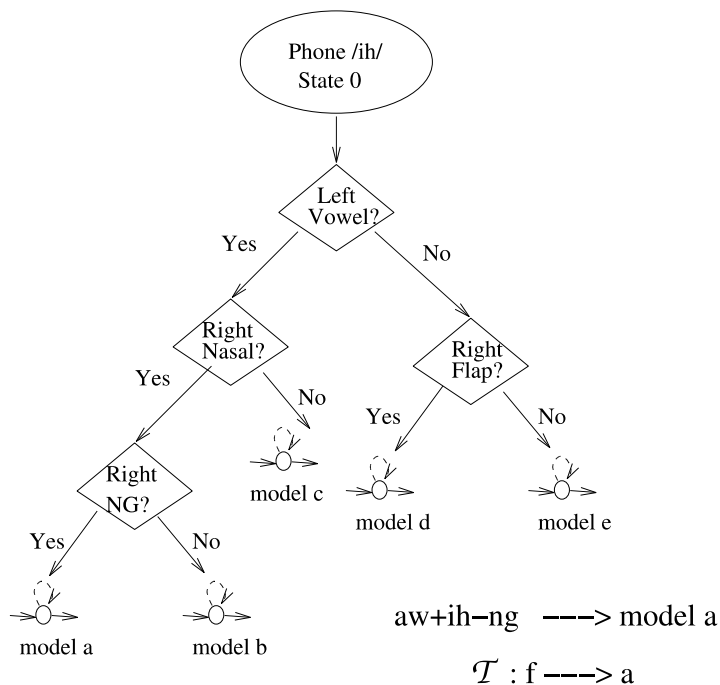


Fig. 1. Decision tree maps context-specific phonemes to acoustic models.

## 2.2. *Related work on tagged clustering*

The idea of using word- and syllable-level features in a decision tree framework for conversational speech is supported by a study conducted at a summer workshop in JHU (Ostendorf, Byrne, Bacchian, & Finke, 1996). Clustering a subset of standard training data for conversational speech with triphones that were coded with these features, it was found that questions regarding them were asked early, i.e., near the top of the tree. This may lead to finding better equivalence classes, and thus improve acoustic models. To make this study possible, the computational cost was reduced by ignoring the triphones that span word boundaries, i.e., evaluating performance with a word-internal triphone system. Thus, the usefulness of these features was not demonstrated conclusively.

The use of word position (initial, medial, final) as a context-conditioning feature has been shown to be useful in several studies, for both conversational speech (Finke, 1997; Gunawardana, 1998) and read speech (Reichl & Chou, 1999), and is used in many research systems. The use of syllable position (initial, medial, final) alone has not so far proved to be useful (Paul, 1997; Gunawardana, 1998), though Paul reports a small gain when syllable position is used in combination with lexical stress tags. Paul's results for lexical stress are also mixed, with gains depending on the dictionary used. The mixed results on read speech could be due to a variety of reasons, including inadequate levels for coding features or the sensitivity to the alignment of the training data used. Hence, we chose to re-evaluate the use of syllable vs. word position in clustering for conversational speech.

Tagged clustering studies have also looked at other features. Word type (function vs. content) was found to be useful in experiments recognizing read speech (Lee & O'Shaughnessy, 1997), and a preliminary study with prosody shows the potential for improving acoustic models for conversational speech (Shafran, Ostendorf, & Wright, 2001). In our work, we will restrict our experimental study to word and syllable features, but the development of the multi-stage clustering approach makes possible the use of a greater number of features in general, which might include these and other features.

## 2.3. *Use of richer syllable features*

In this work, we used a rich set of symbols to represent syllable structure, which includes consonant cluster and ambisyllabicity. The lexicon and syllable coding system used in this work was developed at the 1996 JHU workshop, and later extended for new words. The lexical expansion of words are coded at the phoneme level as illustrated in Table 1. Note that the position of the phone in the syllable distinguishes between onset consonants which are and are not in clusters, and marks consonants as onset even if they are not syllable initial, unlike previous work on syllable position. Also, unlike previous work, stress is marked using three levels (taken from Pronlex and as in most dictionaries): primary, secondary and unstressed. Interpreting primary stress as the default position of the strong syllable of a word and secondary stress as a potential position for a strong syllable (i.e., can receive a pitch accent), we labeled monosyllabic content and function words as having primary and secondary stress, respectively, so that all syllables in the dictionary were marked with one of the three levels. Here, stress is not an indicator of acoustic prominence, but rather the potential for a syllable to be relatively stronger or weaker.

Table 1  
Coding of word- and syllable-features in the dictionary

Phone position in word	Syllable position in word	Phone position in syllable	Stress
First	First	Onset initial	Stress-less
Middle	Middle	Onset other	Primary
Last	Last	Nucleus	Secondary
Only	Only	Coda only	
		Coda initial	
		Coda other	
		Ambisyllabic	

The state-level segmentation of the training data (from a triphone system) for each phoneme was encoded with the word- and syllable-level features from the corresponding lexical expansion of the word in the lexicon. Initial acoustic models for each context-specific phoneme were estimated from these alignments. These models were then clustered using decision tree training, and questions about the word- and syllable-level features were used in addition to the standard contextual questions about neighboring phonemes. Subsequently, the estimate of the distribution associated with each cluster was refined using EM iterations.

Syllabification may vary systematically at word boundaries, depending on the neighboring word. For example, “just a” may be resyllabified as “[jh ah s][t ax]” instead of “[jh ah s t][ax]”, “choirs and” as “[k w ay r][z ax n]” instead of “[k w ay r z][ax n]” and “it’s a” as “[ih t][s ax]” instead of “[ih t s][ax]”, where the latter forms are obtained by stringing together syllables of each word. A complete description of the process of resyllabification in English is relatively complex. However, the process of resyllabification can be explained to a large extent by an empirical rule – the Sonority Dispersion Rule (Clements, 1990; Kenstowicz, 1994). The simplest syllabification is the one with maximal and most evenly distributed rise in sonority at the beginning and the minimal drop in sonority at the end. Based on this principle, the Sonority Dispersion Rule moves the syllable boundary amongst the consonants to minimize the slope of sonority in the nucleus-coda demi-syllable of the pre-boundary word, and maximize it in the onset-nucleus demi-syllable of the following word. Sonority ranks can be assigned to groups of phonemes based on their phonetic properties; here we followed the convention in (Clements, 1990).

Since a small number of rules captures a large fraction of the cases of resyllabification, it is possible to incorporate resyllabification into the acoustic model of a word by allowing alternate “pronunciations”. The specific method used in our work is as follows:

- Candidates for resyllabification include only open-vowel syllables that are word initial and do not follow a pause or a vowel, based on forced alignments in training and N-best hypotheses in testing. (Our analysis of ICSI transcriptions shows that 73% of the resyllabifications in speech occurs in open-vowel word-initial syllables.)
- If the syllable is preceded by a single consonant, mark that consonant as optionally ambisyllabic.
- If the syllable is preceded by more than one consonant, apply the Sonority Dispersion Rule to obtain the alternate syllable boundary.

To study the effect of resyllabification, we performed a series of experiments, using the rules mentioned above, as further described in Section 4.

### 3. Multi-stage clustering

There are a few limitations in using standard decision tree design techniques for clustering phonemes when they are coded with a large number of features. The number of elementary coded phonemes (or, context combinations) increases drastically with the number of features, increasing the memory requirements for storing sufficient statistics of all possible coded phones. Each code of cardinality  $M$  increases the number of unique context combinations by a factor of  $M^3$  when the code applies to left, center and right contexts (in practice, they are constrained by the diversity of the data). In addition, the large number of partitions that need to be tested to use these features – potentially  $3 \times 2^{M-1}$  questions to test at each stage – raises the computational cost of clustering. Furthermore, phonemes in infrequent contexts, which constitute a large fraction of the phonemes, are estimated poorly and the partitions learned may not represent general trends in speech. For example, the percent of observed contexts with fewer than five frames are 3% of 32k triphones, 42% of 2100k pentaphones, and 6% of 186k syllable-coded phone-states.

These problems have restricted previous work on tagged clustering. For example, in (Ostendorf et al., 1996), experimental costs were reduced by restricting the use of syllable boundary and stress only to word-internal triphones, discarding triphones that spanned word boundaries. In other work, cross-word context is used but only with simple tag sets such as word position (begin, middle, end). We address these problems by introducing a new tree design technique based on multi-stage clustering.

Our approach to reduce the storage and computational costs for clustering is based on dividing the task into multiple stages. The decision tree can be viewed as a function,  $\mathcal{F}$ , that maps a feature vector,  $\mathbf{f}$ , consisting of contextual information to an index,  $a$ , a particular state of an acoustic model, thus  $\mathcal{F} : \mathbf{f} \rightarrow a$ . As illustrated in Fig. 2, for two-stage clustering, we group the contextual

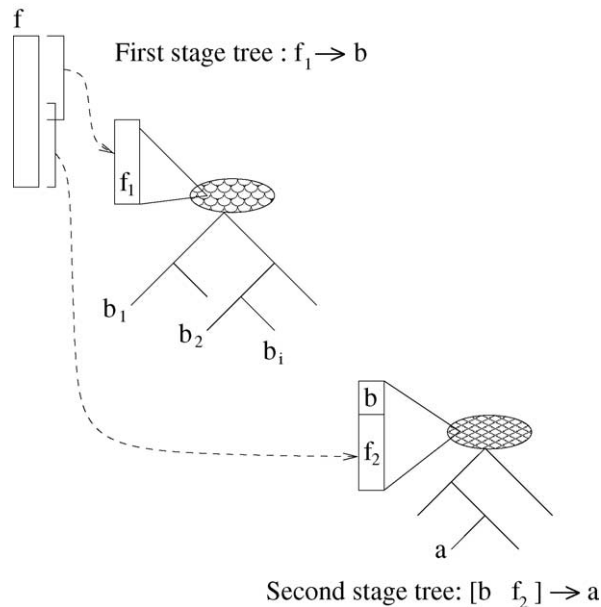


Fig. 2. Multi-stage clustering illustrated here with two stages.



information into two feature vectors  $\mathbf{f}_1$  and  $\mathbf{f}_2$ , optionally allowing some common components between them. In the first stage, the training data is annotated only with the values of vector  $\mathbf{f}_1$ . Using the annotated data, we grow a decision tree,  $\mathcal{T}_1$ , which maps the different values of  $\mathbf{f}_1$  to the index of its leaves  $B$ , thus  $\mathcal{T}_1 : \mathbf{f}_1 \rightarrow b$  where  $b \in B$ . In the second stage, the training data is annotated with  $\mathbf{f}_2$  along with the value of  $b$  which is obtained by dropping its context  $\mathbf{f}_1$  down the tree  $\mathcal{T}_1$ . Using the newly annotated training data, a new decision tree,  $\mathcal{T}_2$  is grown that maps  $[b \mathbf{f}_2]$  to the index of acoustic models as represented by the leaves of  $\mathcal{T}_2$ , thus  $\mathcal{T}_2 : [b \mathbf{f}_2] \rightarrow a$ .

In current decision tree clustering for speech recognition, questions about features are defined by hand and are linguistically motivated. This is straightforward for the features in  $\mathbf{f}_1$  and  $\mathbf{f}_2$ , but not for the index  $b$ . Allowing all possible partitions of  $B$  is impractical since there are  $2^{|B|}$  binary partitions, and to use the features in the first stage adequately, the size of  $\mathcal{T}_1$  (i.e., the number of leaves,  $|B|$ ) needs to be large. To solve this problem we define binary questions that test whether a node  $b$  belongs to a subtree of the first tree  $\mathcal{T}_1$  or not. Such questions are equivalent to compound questions which are obtained by performing an “and” operation on the set of binary questions about the features in  $\mathbf{f}_1$  that lead to node  $b$ . Defining questions on subtrees permits the decision tree to test a large number of partitions, and is more efficient than allowing all partitions. For example, asking only 9 questions leads to  $M = 10$  leaves, in which case the naive partition allows  $2^9 = 512$  questions on the associated variable. Instead, we define questions over subtrees of  $\mathcal{T}_1$ , so that the number of questions is  $O(M)$ . Restricting questions to subtrees of size 3 on the 10-leaf tree would lead to a set of 4–8 questions, depending on the balance of the tree. In our experiments with larger trees, we restrict questions to subtrees of at least five leaves.

Once the second stage tree  $\mathcal{T}_2$  is grown, the questions on subtrees in  $\mathcal{T}_1$  are replaced with the equivalent compound questions to obtain a single tree. For example, if  $\mathcal{T}_2$  chooses the question, “ $b \in \text{subtree}(j)$ ?”, then it is replaced with the sequence of questions from  $\mathcal{T}_1$  that lead to the root node of subtree  $j$ . Note that, in principle, there is no limit to the number of stages, but this work considers only two.

The multi-stage clustering techniques helps ameliorate the problem of sparse data by reducing the number of coded units for which sufficient statistics need to be estimated, since only a subset of the features are used at each clustering stage. The root node at every stage has all the data available to it or, in this case, all the data associated with the particular state and the phone. The number of elementary units that need to be clustered in stage  $i$  depends on the features  $\mathbf{f}_i$  used in that stage and, if  $i > 1$ , the number of leaves of the preceding tree  $\mathcal{T}_{i-1}$ . Both of these factors can be controlled to reduce the effects of data fragmentation, essentially by trading off the potential for more directly modeling interaction between features (with a large dimension  $\mathbf{f}_i$ ) with the robustness (and computational) advantages of a low dimension feature set  $\mathbf{f}_i$ . Note that robust estimation of statistics of elementary units also benefits from the general principle of increasing system complexity incrementally. In particular, we use phone alignments from our best triphone system, rather than bootstrapping from monophone models, as shown to be important in (Gunawardana, 1998).

The storage and computational cost of the multi-stage clustering depends on various factors. The number of sufficient statistics that need to be clustered in the two stages is determined by the number of components used in  $\mathbf{f}_1$  and  $\mathbf{f}_2$ , and the size of  $\mathcal{T}_1$ , as mentioned above. The number of sufficient statistics is limited by the diversity of the data, and also depends on how uniformly the training data is divided into the clusters (how balanced the tree is). If a maximal tree is grown for  $\mathcal{T}_1$ , then the multistage clustering will be computationally more expensive than clustering a single

tree. The size of  $\mathcal{T}_1$  may be set using a constraint on total number of leaves, or saturation of likelihood increase. To reduce the number of partitions that are tested in the second stage, we selected only a subset of subtrees (the largest) from the first stage during clustering at the top of the second stage, where the largest proportion of the computation occurs in training decision trees.

## 4. Experiments, results and observations

### 4.1. Experimental framework

#### 4.1.1. Speech corpus

We used the Switchboard and Callhome corpora, which together provide a collection of about 140 h of spontaneous telephone conversations between pairs of callers in American English (Godfrey, Holliman, & McDaniel, 1992). All the experiments were tested on a subset of the 1998 NIST Hub-5 development test set (NIST, 1998), consisting of about 12.5k words in approximately an hour of speech from 14 conversations. In addition, tests were also performed on 2000 NIST Hub-5 evaluation test set (NIST, 2000), consisting of about 42k words in approximately three and a half hours of speech from 40 conversations. All results use the standard criterion of word error rate (WER), where  $WER = (C - I)/R$ , the number of words correct is  $C$ , inserted is  $I$  and the total number of words in the reference transcript is  $R$ .

#### 4.1.2. Recognition system

The speech data is preprocessed to produce a 14-dimensional vocal-tract-length-normalized MFCC vector sequence augmented with its first-order derivatives, at a rate of 100 vectors per second. This serves as acoustic input to the recognition system (Zavaliagos, McDonough, Miller, & El-Jaroudi, 1998). Two types of acoustic models were used in our tests. In our initial experiments, a standard left-to-right HMM topology with 5 states and skips is used to model the acoustic units, with a single full covariance Gaussian as the observation distribution for each state. Allophones of each phoneme and state are clustered to produce 10,000 means and 2000 subtree-shared covariances. By avoiding expensive mixture splitting, this allows acoustic models to be trained quickly, at the cost of a small performance degradation. In the experiments on the evaluation test set, three-state HMMs (with no skips) and 8000 clustered states were modeled by more complex 16-mixture Gaussian distributions with diagonal covariances. These models were trained in the feature space containing derivatives as well as accelerations of the MFCC. In each table of experiments described below, the number of parameters in the acoustic model were kept fixed.

To reduce experimentation time, we restricted our experiments to a lattice re-scoring decoding paradigm and did not adapt the models to the speaker being tested. The word-level lattices for the development set were derived from the 100 best hypotheses provided by BBN. It has an oracle WER of about 20% and 1-best WER of 42.6%.<sup>1</sup> These lattices contained language model scores

---

<sup>1</sup> The 1-best WER for the development set is actually better than the results reported here, because the N-best hypotheses were generated with a significantly different system. While it is often useful in rescoring to combine scores from different acoustic models, since our focus was on understanding the behavior of the syllable features, the BBN and the AT&T acoustic model scores were not used in the results reported here.

from a part-of-speech smoothed trigram trained with Broadcast news data as well as the Switchboard and Callhome data (Iyer & Ostendorf, 1997). The lattices for evaluation set was generated using gender-independent models of AT&T switchboard system, and has an oracle WER of about 10% and 1-best WER of 35.4% (Ljolje, Hindle, Riley, & Sproat, 2000). For rescoreing them, language model scores were obtained from AT&T 6-gram language model.

#### 4.1.3. *Lexicon*

For this study, we use the syllabified dictionary and the coding system that was developed at the 1996 summer workshop at John Hopkins University (JHU-WS-Lexicon, 1996). A brief review of the lexicon is given here; further details can be found in (Ostendorf et al., 1996). The stress markings for the multisyllabic words in this lexicon are taken from Pronlex dictionary; monosyllabic words were marked with either primary or secondary stress depending on whether or not the word was a content word, as described in Section 2.3. Syllabification for this lexicon was performed automatically using Fisher's implementation of Kahn's principles for English syllabification (Fisher, 1996). The syllabification is performed by assigning as many consonants to syllable onsets as possible (maximal onset rule) where permitted onsets were predefined. Among the possible syllabifications of a word, the most casual variant was selected to represent the nature of conversational speech. In this process, the morpheme boundaries were not taken into consideration. However, the use of casual variants of syllabification mitigates the associated syllabification errors, since many of the consonants at the boundary were labeled as ambisyllabic, rather than with the wrong syllable. To syllabify foreign words, an augmented list of permitted onsets was applied on those words that initially failed to parse.

#### 4.2. *Testing syllable features*

We developed gender-dependent systems using information about word and syllable structure, as mentioned in Section 2.3. During acoustic model training, the decision trees were allowed to ask questions about syllable and word information of the center and the immediate neighboring phonemes, in addition to the questions about triphone context. The recognition performance of these systems were compared with triphone and pentaphone systems with same number of model parameters. In all the systems described below, the clustered triphones were trained using a single stage of clustering from the same base triphone alignment and then re-estimated with a few passes of EM training. The results are summarized in Table 2 for the development set and in Fig. 3 for the evaluation set.

While the difference in performance of the triphone and pentaphone systems on the development set is not statistically significant,<sup>2</sup> the system based on word- and syllable-features is significantly better than triphone system, according to NIST statistical significance tests (at the level of  $p = 0.001$  for matched pair sentence segment, Wilcoxon signed rank and McNemar). Contrary to other reported results, we find consistent gain from using features in addition to word position, particularly at higher beams (with a significance level of  $p = 0.002$  for McNemar test). Due to its

---

<sup>2</sup> Other systems showing improved performance with pentaphones appear to have increased numbers of parameters in the pentaphone system, whereas here the number is constrained to be roughly the same.

Table 2  
WER of systems using different features in clustering (development test set)

System	WER (%)
Triphone	44.56
Pentaphone	44.37
Triphone + word-position	44.31
Triphone + word-position + syllable-features	44.05

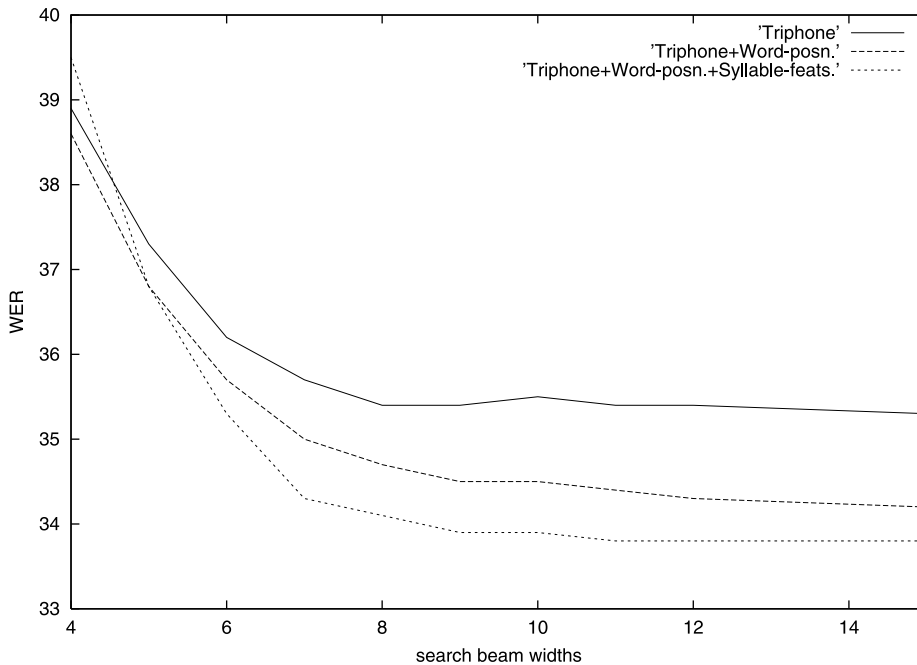


Fig. 3. WER of systems using different features in clustering (evaluation test set), at different search beam widths.

complexity, we did not rescore the bushier lattices of the evaluation set with the pentaphone system.

The computational cost for decoding and training were also significantly lower for the system with word- and syllable-features. In training, the pentaphone system required testing 350 potential partitions for clustering up to 2.5 M acoustic units, while the system with syllable features required testing only 200 potential partitions for 700k acoustic units. In decoding, the system with syllable features was 20% faster. In addition, unlike pentaphones, which incur extra computational expense and software flexibility to span the two forward contexts, the coded-triphones only look ahead as much as a single phone and could be incorporated in a standard first pass triphone decoder.

#### 4.3. Observations on the use of syllable features

To study how the syllable features were used, we analyzed the questions that were chosen in the decision tree for clustering acoustic units. Two measures were used to characterize the utility of

Table 3

Word- and syllable-feature usage in decision trees trained on the Switchboard and Callhome corpora (about 140 h of speech)

Feature	Questions in the tree (%)	Data affected (%)	Degrees of freedom
Triphone	66	70	$2 \times 90 = 180$
Phone in word	10	37	$3 \times 6 = 18$
Syllable in word	11	25	$3 \times 6 = 18$
Phone in syllable	3	31	$14 + 3 + 4 = 21$
Stress	11	32	$3 \times 2 = 6$

features – the number of questions asked about a feature in the tree and the percentage of total data affected or split by the feature. To make a fair comparison across the columns, the number of partitions tested (or degrees of freedom) for each feature also needs to be taken into account (see Table 3).

Even though questions about phone identity allow three times as many degrees of freedom (in the sense of number of candidate binary partitions) as all other features combined, the decision trees chose the latter at least one in three times. Among all the features, the fewest questions are asked about position of the phone in the syllable. However, it affected disproportionately a larger amount of data. This feature is likely to be used higher up in the tree, suggesting that questions about it generalize over a large fraction of the data. The lexical stress and position of the syllable in the word is least useful, as it affects the least amount of data. It was also observed that the questions about the position of the center and the right phone in the syllable is significantly more important (4–6 times the fraction of data affected) than that of the left phone. Whether the center phone was in a monosyllabic word was among the top questions about the position of syllable in the word, as one might expect from the gains observed in modeling monosyllabic word (Dodding et al., 1997), but the amount of data affected was not high so it did not stand out as a particularly important feature. In general, the pattern of usage of the features across gender is similar.

Interestingly, even though the pentaphone models had a higher likelihood on the training data, the syllable system had a better likelihood on the test data. Analogous to using likelihood on independent data as a model selection criterion, we argue that the higher likelihood indicates that syllable system generalizes better, as confirmed by improvements in word error rate.

#### 4.4. Impact of resyllabification

First, we studied the effect of resyllabification on the training data. The method described in Section 2.3 was applied on a single phoneme path to generate alternate syllabifications across word boundaries for all vowel-initial words in the training set. This produced about 13% coded phonemes in the alternate resyllabified paths, and constituted only 50 new types of coded phonemes (about 3% of the total number of uniquely coded phones in the lexicon) distributed across 15 phonemes. Next, using the acoustic models developed in Section 2.3, we let the decoder choose the best path from the lattice of possible paths. A few examples of resyllabification that were chosen are listed in Table 4.

The number of alternate coded phonemes that the decoder chose was only about 7% of those hypothesized. This was about 1% of the total labels in the training data and involved many different

Table 4

Examples of resyllabifications automatically chosen in the training data

Partial word sequence	Resyllabifications
bridge across	[b r ih [jh] ax]. . .
did it	[d ih [d] ih t]
work and	[w er [k] ax n]
choirs and	[k w ay r][z ax n]. . .
drugs out there	[d r ah g][z aw t]. . .
lots of times I'd	. . . [t ay m][z ay d]
its an	[ih t][s ax]. . .
minutes up	[m ih][n ih t][s ah p]
takes a	[t ey k][s ax]

phones. In comparison, we observed about 4.5% vowel-initial resyllabifications in the ICSI transcripts. Thus, we may have captured about one in four resyllabifications. Examination of the chosen resyllabifications indicate that a large number of them involved phoneme “s” and “z” or were ambisyllabic. Since the total data affected was low, we did not expect resyllabification to affect our acoustic model significantly, and continued using the same models without bootstrapping.

The impact of resyllabification on test data was evaluated using N-best re-scoring. After expanding each of the 100 hypotheses into a lattice with alternate resyllabification paths, we let the decoder choose the best path using the acoustic models developed in Section 4.2. The word error rate did not show any improvements. This may be explained simply by the small number of tokens affected by resyllabification; however, there are other possible explanations for this result. For example, resyllabification tends to occur in high frequency word pairs, where it may be that the language model score adequately compensates for any loss in acoustic match. Further, it is likely that the combination of right phonetic context and word position allows the clustering trees to encode resyllabification implicitly,<sup>3</sup> in which case explicit modeling of resyllabification is unnecessary. It may also be the case that resyllabifications where a coda becomes ambisyllabic are missed by our system because of the acoustic similarity. In acoustic clustering, among the data affected by questions about syllable position, only 5% was affected by questions about coda and 11% about ambisyllabic. In cases where an explicit resyllabification provides an improved characterization and the models are initially biased toward the un-resyllabified hypotheses, training acoustic models iteratively may improve the results, but the number of such cases may not warrant the added system complexity.

#### 4.5. Verification of multi-stage clustering

To evaluate the effectiveness of multi-stage clustering we trained gender-specific pentaphone systems using standard single-stage clustering and two-stage clustering. The systems were trained from a base triphone alignment with one pass of Viterbi training and a few passes of EM. In the first stage of the second system, we clustered the data into 1000 clusters for each of the five states,

<sup>3</sup> In a large number (>90%) of resyllabifications in the ICSI transcriptions, the syllable boundary moves across the word by one phone (e.g., [th ih ng z] [aa n] → [th ih ng] [z aa n]).

Table 5

Word error rates of systems trained with one vs. two stages of clustering

System	
(a) Pentaphone: 1 stage	44.37%
(b) Pentaphone: 2 stage	44.39%

using the second phone neighbors as features. In the second stage, we clustered the data using the leaf indices of  $\mathcal{T}_1$  along with the triphone context to obtain the final models. This specific order of the split features ( $\mathbf{f}_1 : \pm 2$  neighbors;  $\mathbf{f}_2 : \pm 1$  neighbors) was chosen after comparing the likelihoods of resulting models from both orders on an independent data set. While much more sophisticated mechanisms could be envisioned for choosing feature subsets, we did not expect this to lead to a significant gain in performance on our task.

The two-stage pentaphone system for both genders performed as well as the one-stage systems. Thus the result shows that incorporating features in multiple stages is a viable method for using a large number of features in acoustic modeling (see Table 5).

The memory used in clustering is directly proportional to the number of unique contexts to be clustered. In the single-stage pentaphone system, we had about 2M unique contexts, whereas in two stage system, we had only about 78K in the first stage and less than 0.5M in the second stage, thus reducing the memory requirement at any time by a factor of 4. This could be reduced further by shrinking the size of the first tree. The computational cost of two-stage clustering in this case is half that of single-stage clustering.

## 5. Conclusions

We have shown a small but consistent improvement in using syllable structure in addition to word position in a large vocabulary recognition task. This, is in contrast to other reported results, and may be due to our use of high quality state alignments and a more detailed syllable coding system. The results also suggest that the syllable features generalize better than long span (pentaphone) phonetic context. Perhaps more importantly, the system using syllable features has lower training and decoding computational costs than a pentaphone system of equivalent size. In addition, our studies show that alternate paths of resyllabification predicted by general linguistic rules do not provide any improvement in recognition performance, though it is possible that gains may be obtained with further iterations of training.

To take full advantage of syllable features, we conjecture that temporal variation must also be modeled. For example, the fixed-state topology could be replaced with a context-specific topology. This may be carried out within a decision tree framework such as (Eide, 1999).

We have also developed a multi-stage clustering system that enables the use of a large number of features in clustering. Multi-stage clustering addresses the general issue of unreliable estimates of infrequent context as well as the higher computational cost incurred in clustering them. Speech recognition experiments show that the approach performs as well as a single stage of clustering with significantly reduced computational costs. In the work described here, the feature groups were chosen heuristically based on linguistic intuitions. The approach could be extended to incorporate data-driven techniques for organizing features into hierarchies for use in the different stages.

The features explored in this work included phonetic context and syllable structure, but they could easily be expanded to include other features such as speaking rate (quantized into finite levels), hyperarticulation, word type (function vs. content word), prosodic constituent structure, or other factors that have been shown to have some effect on ASR performance. An advantage of the clustering approach over explicit estimation of different models is that data is not divided unnecessarily.

## Acknowledgements

We thank Prof. Richard Wright for advice on resyllabification rules, Dr. Michiel Bacchiani and Dr. Michael D. Riley for making their software tools available, and Dr. Owen Kimball at BBN and Dr. Andre Ljolje at AT&T for providing their first stage recognition outputs. This work was supported by the National Science Foundation, Grant No. ISI-9618926. The views and conclusions contained in this document are those of the authors and should not be interpreted as reflecting the official policies of the funding agency.

## References

- Byrne, B., Finke, M., Khudanpur, S., McDonough, J., Nock, H., Riley, M., Saraclar, M., Wooters, C., Zavaliagkos, G., 1997. Pronunciation modeling for conversational speech recognition: a status report from WS97. In: IEEE Workshop on Automatic Speech Recognition and Understanding Proceedings (ASRU), Santa Barbara, CA, pp. 26–33.
- Clements, G.N., 1990. The role of the sonority cycle in core syllabification. In: Kingston, J., Beckman, M.E. (Eds.), *Papers in Laboratory Phonology I: Between the Grammar and Physics of Speech*. Cambridge University Press, Cambridge, MA, pp. 283–333.
- Doddington, G., Corrada, A., Wheatley, B., et al., 1997. Syllable based speech processing. In: *Proceedings of the 1997 Summer Workshop on Speech Recognition*, <http://www.clsp.jhu.edu/ws97/syllable>.
- Eide, E., 1999. Automatic modeling of pronunciation variation. In: *Proceedings of the European Conference on Speech Communication and Technology*, vol. 1, pp. 451–454.
- Finke, M., Rogina, I., 1997. Wide context acoustic modeling in read vs. spontaneous speech. In: *Proceedings of the International Conference on Acoustic, Speech and Signal Processing*, Munich, 1997. vol. III, pp. 1743–1746.
- Fisher, W., 1996. A C implementation of Daniel Kahn's theory of English syllable structure, <ftp://jaguar.ncsl.nist.gov/pub/tsylb2-1.1.tar.Z>.
- Fosler-Lussier, E., Greenberg, S., Morgan, N., 1999. Incorporating contextual phonetics into automatic speech recognition. In: *Proceedings of the ESCA Workshop on Modeling Pronunciation Variation for Automatic Speech Recognition*, pp. 611–614.
- Greenberg, S., 1998. Speaking in shorthand – a syllable-centric perspective for understanding pronunciation variation. In: *Proceedings of the Workshop Modeling Pronunciation Variation in Automatic Speech Recognition*, pp. 47–56.
- Greenberg, S., Chang, S., 2000. Linguistic dissection of Switchboard-corpus automatic speech recognition systems. In: *Proceedings of the ISCA Workshop on Automatic Speech Recognition: Challenges for the New Millennium*.
- Greenberg, S., Fosler-Lussier, E., 2000. The uninvited guest: information's role in guiding the production of spontaneous speech. In: *Proceedings of the Crest Workshop on Models of Speech Production: Motor Planning and Articulatory Modeling*, Kloster Seon, Germany.
- Godfrey, J., Holliman, E., McDaniel, J., 1992. Switchboard: telephone speech corpus for research and development. In: *Proceedings of the International Conference on Acoustic, Speech and Signal Processing*, vol. 1, pp. 517–520.
- Gunawardana, A., 1998. Personal communication.



- Hain, T., Woodland, P.C., 2000. Modelling sub-phone insertions and deletions in continuous speech recognition. In: Proceedings of the International Conference on Speech and Language Processing, vol. 4, pp. 172–176.
- Hamaker, J., Ganapathiraju, A., Picone, J., Godfrey, J., 1998. Advances in alphadigit recognition using syllables. In: Proceedings of the International Conference on Acoustics, Speech and Signal Processing, vol. 1, pp. 421–424.
- Iyer, R., Ostendorf, M., 1997. Transforming out-of-domain estimates to improve in-domain language models. In: Proceedings of the European Conference on Speech Communication and Technology, vol. 4, pp. 1975–1978.
- JHU-WS-Lexicon, 1996. [http://www.clsp.jhu.edu/ws96/ws96\\_workshop.html](http://www.clsp.jhu.edu/ws96/ws96_workshop.html).
- Jones, R.J., Downey, S., Mason, J.J., 1997. Continuous speech recognition using syllables. In: Proceedings of the European Conference on Speech Communication and Technology, vol. 3, pp. 1171–1174.
- Kannan, A., Ostendorf, M., Rohlicek, J.R., 1994. Maximum likelihood clustering of Gaussians for speech recognition. *IEEE Transactions on Speech and Audio Processing* 2 (3), 453–455.
- Keating, P., Wright, R., Zhang, J., 1999. Word-level asymmetries in consonant articulation. In: *UCLA Working Papers in Phonetics* No. 97.
- Kenstowicz, M., 1994. The syllable and syllabification. In: *Phonology in Generative Grammar*. Blackwell Publishers, UK, pp. 250–309.
- Le, A., 2002. The 2002 NIST RT Evaluation Speech-to-Text Results. In: Proceedings of the Rich Transcription 2002 Workshop, May 2002.
- Lee, C.Z., O'Shaughnessy, D., 1997. Clustering beyond phoneme contexts for speech recognition. In: Proceedings of the European Conference on Speech Communication and Technology, vol. 1, pp. 19–22.
- Ljolje, A., Hindle, D.M., Riley, M.D., Sproat, R.W., 2000. The AT&T LVCSR-2000 system. In: Proceedings of the 2000 Speech Transcription Workshop.
- Lleida, E., Marino, J.B., Nadeu, C., Salavedra, J., 1991. Demisyllable-based HMM spotting for continuous speech recognition. In: Proceedings of the International Conference on Acoustics, Speech and Signal Processing, vol. 1, pp. 709–712.
- Marino, J.B., Nogueiras, A., Bonafonte, A., 1997. The demiphone: an efficient sub-word unit for continuous speech recognition. In: Proceedings of the European Conference on Speech Communication and Technology, vol. 3, pp. 1215–1218.
- Massaro, D.W., 1972. Perceptual images, processing time, and perceptual units in auditory perception. *Psychological Review* 79, 124–145.
- McAllister, D., Gillick, L., Scattoni, F., Newman, M., 1998. Fabricating conversational speech data with acoustic models: a program to examine model-data mismatch. In: *International Conference on Spoken Language Processing*, vol. 5, pp. 1847–1850.
- NIST-HUB5-1998, [http://www.nist.gov/speech/tests/ctr/hub5e\\_98/current-plan.htm](http://www.nist.gov/speech/tests/ctr/hub5e_98/current-plan.htm).
- NIST-HUB5-2000, [http://www.nist.gov/speech/tests/ctr/h5\\_2000](http://www.nist.gov/speech/tests/ctr/h5_2000).
- Ostendorf, M., Byrne, B., Bacchian, M., Finke, M., et al., 1996. Modeling systematic variation in pronunciation via language-dependent hidden speaking mode. In: Proceedings of the 1996 Summer Workshop on Speech Recognition, <http://www.clsp.jhu.edu/ws96/zilla/hidden-mode/group.html>.
- Pallett, D., Fiscuss, J., Garofolo, J., Martin, A., et al., 1999. 1998 broadcast news benchmark test results: English and non-English word error rate performance measures. In: Proceedings of the DARPA Broadcast News Workshop, pp. 5–12.
- Paul, D.B., 1997. Extensions to phone-state decision tree clustering: single tree and tagged clustering. In: *International Conference on Acoustics Speech and Signal Processing*, vol. 2, pp. 1487–1490.
- Rabiner, L., Juang, B.H., 1993. *Fundamentals of speech recognition*, Prentice Hall (Signal processing series).
- Reichl, W., Chou, W., 1999. A unified approach of incorporating general features in decision tree based acoustic modeling. In: *International Conference on Acoustics Speech and Signal Processing*, vol. 2, pp. 573–576.
- Riley, M., Byrne, B., Finke, M., Khudanpur, S., Ljolje, A., McDonough, J., Nock, H., Saraclar, M., Wooters, C., Zavalagkos, G., 1999. Stochastic pronunciation modeling from hand-labeled phonetic corpora. *Speech Communication* 29, 209–224.
- Saraclar, M., Nock, H., Khudanpur, S., 2000. Pronunciation modeling by sharing Gaussian densities across phonetic models. *Computer Speech and Language* 14 (2), 137–160.

- Shafran, I., Ostendorf, M., Wright, R., 2001. Prosody and phonetic variability: lessons learned from acoustic model clustering. In: Proceedings of ISCA Tutorial and Research Workshop on Prosody in Speech Recognition and Understanding, pp. 127–131.
- Weintraub, M., Taussig, K., Hunicke-Smith, K., Snodgrass, A., 1996. Effect of speaking style on LVCSR performance. In: Int Conf. Speech and Language Processing Supplement, pp. 16–19.
- Wu, S.-L., 1998. Incorporating information from syllable-length time scales into automatic speech recognition, Ph.D. Thesis, UC Berkeley.
- Young, S.J., Woodland, P.C., 1994. State clustering in HMM-based continuous speech recognition. *Computer Speech and Language* 8 (4), 369–384.
- Zavaliagos, G., McDonough, M., Miller, D., El-Jaroudi, A., et al., 1998. The BBN Byblos 1997 large vocabulary conversational speech recognition system. In: Proceedings of the International Conference on Acoustic, Speech and Signal Processing, vol. 2, pp. 905–908.