

# PRONUNCIATION MODELING USING A FINITE-STATE TRANSDUCER REPRESENTATION

*Timothy J. Hazen, I. Lee Hetherington, Han Shu, and Karen Livescu*

Spoken Language Systems Group  
MIT Laboratory for Computer Science  
Cambridge, MA, USA

## ABSTRACT

The MIT SUMMIT speech recognition system models pronunciation using a phonemic baseform dictionary along with rewrite rules for modeling phonological variation and multi-word reductions. Each pronunciation component is encoded within a finite-state transducer (FST) representation whose transition weights can be probabilistically trained using a modified EM algorithm for finite-state networks. This paper explains the modeling approach we use and the details of its realization. We demonstrate the benefits and weaknesses of the approach both conceptually and empirically using the recognizer for our JUPITER weather information system. Our experiments demonstrate that the use of phonological rewrite rules within our system reduces word error rates by between 4% and 8% over different test sets when compared against a system using no phonological rewrite rules.

## 1. INTRODUCTION

Pronunciation variation has been identified as a major cause of errors for a variety of automatic speech recognition tasks [8]. In particular, pronunciation variation can be quite severe in spontaneous, conversational speech. To address this problem, this paper presents a pronunciation modeling approach that has been under development at MIT for more than a decade. Our approach systematically models pronunciation variants using information from a variety of levels in the linguistic hierarchy. Pronunciation variation can be influenced by the higher level linguistic features of a word (e.g., morphology, part of speech, tense, etc.) [12], the lexical stress and syllable structure of a word [5], and the specific phonemic content of a word sequence [11, 16]. When all of the knowledge in the linguistic hierarchy is brought to bear upon the problem, it becomes easier to devise a consistent, generalized model that accurately describes the allowable pronunciation variants for particular words. This paper presents the pronunciation modeling approach that has been implemented and evaluated within the SUMMIT speech recognition system developed at MIT.

Pronunciation variation in today's speech recognition technology is typically encoded using some combination of a lexical pronunciation dictionary, a set of phonological rewrite rules, and a collection of context-dependent acoustic models. The component which models a particular type of pronunciation variation can be different from recognizer to recognizer. Some recognizers rely almost entirely on their context-dependent acoustic models to cap-

ture phonological effects, while other systems explicitly model phonological variation with a set of phonological rewrite rules. Some systems ignore phonological rules entirely and simply account for alternate pronunciations directly in the pronunciation dictionary. In this paper we use the SUMMIT recognizer to examine the advantages and disadvantages of accounting for general phonological variation explicitly with phonological rules using distinct allophonic models versus implicitly within context-dependent models. We also describe a pronunciation variation modeling approach which uses a cascade of finite-state transducers, each of which models different variations resulting from different underlying causes.

## 2. OVERVIEW

### 2.1. Segment-Based Recognition

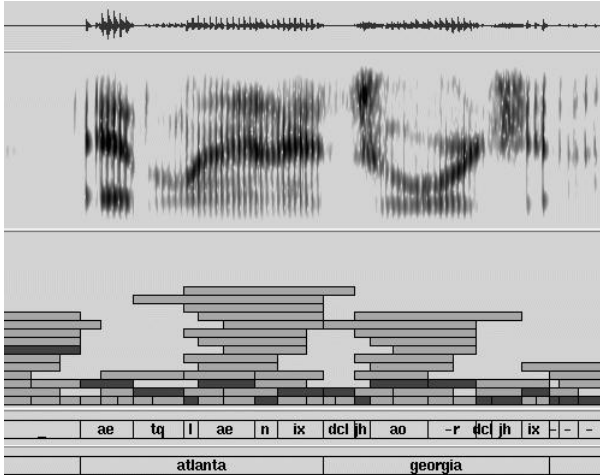
The experiments presented in this paper use the SUMMIT speech recognition system. SUMMIT uses a segment-based approach for acoustic modeling [3]. This approach differs from the standard hidden Markov modeling (HMM) approach in that the acoustic-phonetic models are compared against pre-hypothesized variable-length segments instead of fixed-length frames. While HMM systems allow multiple frames to be absorbed by a single phoneme model via self-loops on the HMM states, our segment-based approach assumes a one-to-one mapping of hypothesized segments to phonetic events. This approach allows the multiple frames of a segment to be modeled jointly, removing the frame independence assumption used in the standard HMM. Details of SUMMIT's acoustic modeling technique can be found in [15].

Figure 1 shows the recognizer's graphical display containing a segment graph (with the recognizer's best path highlighted) along with the corresponding phonetic transcription. It is important to note that SUMMIT pre-generates a segment network based on measures of local acoustic change before the search begins. The smallest hypothesized segments can be as short as a single 10 millisecond frame, but segments are typically longer in regions where the acoustic signal is relatively stationary.

The segment-based approach presents several modeling issues which are essentially not present in frame-based HMM systems. For example, in our segment-based approach plosives must be explicitly modeled as two distinct phonetic events, a closure and a release. In HMM recognizers the closure and burst regions can be implicitly learned by multi-state phoneme models. However, in a segment-based approach they must be explicitly separated into different phonetic models because the segmentation algorithm will observe two distinct acoustic regions and may not hypothesize a

---

This research was supported by DARPA under contract N66001-99-1-8904, monitored through Naval Command, Control and Ocean Surveillance Center.



**Fig. 1.** The output of a graphical interface displaying a sample waveform, its spectrogram, the hypothesized SUMMIT segment network with the best path segment highlighted, the time-aligned phonetic transcription of the best path, and the time-aligned word transcription of the best path.

single segment spanning both the closure and the burst regions.

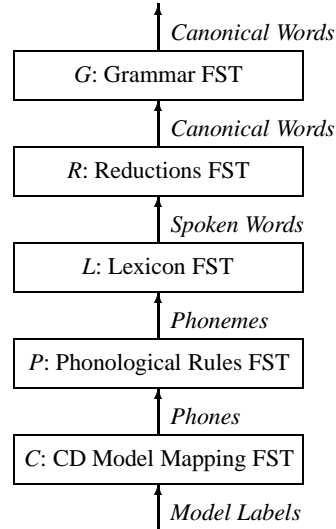
Another issue faced by our segment-based approach is its difficulty in absorbing deleted or unrealized phonemic events required in its search path. An HMM need only absorb as little as one poorly scoring frame when a phonemic event in its search path is not realized, while SUMMIT must potentially absorb a whole multi-frame segment. As a result, accurate phonetic modeling that accounts for potentially deleted phonemic events is more crucial for segment-based approaches than for HMM approaches. It is our belief that accurate phonetic segmentation and classification is important for distinguishing between acoustically confusable words.

## 2.2. FST-Based Search

The SUMMIT recognizer utilizes a finite-state transducer (FST) representation for the lexical and language modeling components. The FST representation allows the various hierarchical components of the recognizer's search space to be represented within a single parsimonious network through the use of generic FST operations such as composition, determinization and minimization [9]. The full search network used by SUMMIT is illustrated in Figure 2. The figure shows the five primary hierarchical components of the search space: the language model ( $G$ ), a set of word-level rewrite rules for reductions and contractions ( $R$ ), the lexical pronunciation dictionary ( $L$ ), the phonological rules ( $P$ ), and the context-dependent model mapping ( $C$ ). Each of these components can be independently created and represented as an FST. By composing the FSTs such that the output labels of the lower-level components become the inputs for the higher-level components, a single FST network is created which encodes the constraints of all five individual components. The full network can be represented mathematically with the following expression:

$$N = C \circ P \circ L \circ R \circ G$$

This paper focuses on the reductions FST  $R$ , the lexicon FST  $L$  and the phonological rules FST  $P$ .



**Fig. 2.** The set of distinct FST components which are composed to form the full FST search network within the SUMMIT recognizer.

## 2.3. Levels of Pronunciation Variation

In our pronunciation modeling approach we distinguish between four different levels of pronunciation variation: (1) variations that depend on word-level features of lexical items (such as part-of-speech, case, tense, etc.), (2) variations that are particular to specific lexical entries, (3) variations that depend on the stress and syllable position of phonemes, and (4) variations that depend only on local phonemic or phonetic context. In the following paragraphs we provide examples of these variants specifically for English.

Type (1) variants include contractions (*what's*, *can't*, etc.), reductions (*gonna*, *wanna*, etc.), part-of-speech variants (as in the noun and verb versions of *record*), and tense variants (as in the past and present tense versions of *read*). In most speech recognition systems, these types of variants are handled in very superficial manners. Reductions and contractions are typically entered into the pronunciation lexicon as distinct entries independent of the entries of their constituent words. All alternate pronunciations due to part of speech or tense are typically entered into the pronunciation lexicon within a single entry without regard to their underlying syntactic properties. In our system reductions and contractions are handled by the reductions FST ( $R$ ), while all other type (1) variants are encoded as alternate pronunciations within lexical entries in the lexicon FST ( $L$ ). In future work we may investigate methods for explicitly delineating pronunciation variations caused by the part-of-speech, case or tense of a word.

Type (2) variants are simply word-dependent pronunciation variants which are not the result of any linguistic features of that word. A simple example of a word with a type (2) variant is *either*, which has two different phonemic pronunciations as shown here:

either: ( iy | ay ) th er

These variants are typically encoded manually by lexicographers. In our system these variants are all handled as alternate pronunciations in the lexicon FST ( $L$ ).

Variants of type (3) in English are typically related to the realization of stop (or plosive) consonants. The set of possible allophones of a stop consonant in English is heavily dependent on its

position within a syllable and the stress associated with the syllables preceding and following the stop. For example, a stop in the suffix or coda position of a syllable can be unreleased, while stops in the prefix position of a stressed syllable must be released. An example is shown here using the word *laptop*:

laptop: l ae pd t aa pd

In this example, the label /pd/ is used to represent a /p/ within a syllable suffix or coda whose burst can be *deleted*. The /t/ in this example is in the onset position of the syllable and therefore must have a burst release. Type (3) variants are encoded using syllable-position-dependent phonemic labels directly in the lexicon FST ( $L$ ). The details of the creation of the pronunciation lexicon using these special labels are presented in Section 3.2.

Variants of type (4) can be entirely determined by local phonemic or phonetic context and are independent of any higher-level knowledge of lexical features, lexical stress, or syllabification. Examples of these effects are vowel fronting, place assimilation of stops and fricatives, gemination of nasals and fricatives, and the insertion of epenthetic silences. To account for type (4) variants we have developed our own FST mechanism for applying context-dependent phonological rules. The details of the syntax and application of the rules are described in [6]. Examples of these rules will be presented in Section 3.3. In relation to Figure 2, type (4) variants are generated by the phonological rules FST ( $P$ ).

#### 2.4. Modeling Variation with Context-Dependent Models

When devising an approach for capturing phonological variation there is flexibility in the specific model in which certain types of phonological variation are captured. In particular, certain forms of phonological variation can easily be modeled either explicitly with phonological rules using symbolically distinct allophonic variants, or implicitly using context-dependent acoustic models which capture the acoustic variation from different allophones within their probability density functions [7]. One example is the place assimilation effect, which allows the phoneme /d/ to be realized phonetically as the palatal affricate [jh] when followed by the phoneme /y/ (as in the word sequence *did you*). The effect could be modeled symbolically with a phonological rewrite rule allowing the phoneme /d/ to be optionally realized as [jh]. Alternately, it can be captured in a context-dependent acoustic model which implicitly learns the [jh] realization within the density function for the context-dependent model for the phoneme /d/ in the right context of the phoneme /y/.

Modeling effects such as place assimilation within the context-dependent acoustic model has several advantages. First, this type of model simplifies the search by utilizing fewer alternate pronunciation paths in the search space. The likelihoods of the alternate allophones are encoded directly into the observation density function of the acoustic models. Additionally, no hard decision about which allophone is used is ever made during either training or actual recognition.

Pushing the modeling of allophonic variation into the context-dependent acoustic model does have potential drawbacks as well. In particular, context-dependent acoustic models may not accurately represent the true set of allophonic variants in cases where stress and syllable-boundary information is required for predicting the allowable set of allophones. For example, consider the two word sequences “*the speech*” and “*this peach*”. Both of these word sequences can be realized with the same phonetic sequence:

th ix s pcl p iy tcl ch

In this particular example, there are two acoustically distinct allophonic variants of /p/; the /p/ in “*the speech*” is unaspirated while the /p/ in “*this peach*” is aspirated. The exact variant of /p/ is determined by the location of the fricative /s/ in the syllable structure. In “*the speech*” the /s/ forms a syllable-initial consonant cluster with the /p/ thereby causing the /p/ to be unaspirated. In “*this peach*” the /s/ belongs to the preceding syllable thereby causing the /p/ to be aspirated. A standard context-dependent acoustic model will model these variants inexactly, allowing the /p/ to be either aspirated or unaspirated in either case. In essence, pushing the modeling of phonological variation into the context-dependent acoustic models runs the risk of creating models which *over-generate* the set of allowable realizations for specific phonemic sequences.

### 3. PRONUNCIATION MODELING IN SUMMIT

#### 3.1. Deriving the Reduction FST

To handle reductions and contractions, a reduction FST is created which encodes rewrite rules that map contractions and other multi-word reductions to their underlying canonical form. Some examples of these rewrite rules are as follows:

gonna → going to  
 how’s → how is  
 i’d → i would | i had  
 lemme → let me

In some cases, such as the contraction *i’d*, a contracted form could represent more than one canonical form. All contractions and reductions which are inputs to the reduction FST ( $R$ ) are re-written such that the input to the grammar FST ( $G$ ) contains only canonical words thereby allowing/constraining the grammar to operate on the intended sequence of canonical words, irrespective of their surface realization. In the JUPITER weather information domain, the reduction FST ( $R$ ) contains 120 different contracted or reduced forms of word sequences.

#### 3.2. Deriving the Lexicon FST

The lexicon FST represents the phonemic pronunciations of the words in the system’s vocabulary (including contractions and reductions). This FST is created primarily by extracting pronunciations from a syllabified version of the PronLex dictionary, which expresses the pronunciations with a set of 41 phonemic labels [10]. A set of rewrite rules is used to generate special phonemic stop labels, which capture information about the allowable phonetic realizations of each stop based on stress and syllable position information. For example, stops in an onset position of a syllable retain their standard phonemic label (/b/, /d/, /k/, etc.) while stops in the suffix or coda of a syllable are converted to labels indicating that their closure can be unreleased with the burst being *deleted* (/bd/, /dd/, /kd/, etc.). In total, the set of 6 standard stop labels are converted into a set of 20 different stop labels for the purpose of encoding the allowable allophones for each stop. The remainder of the phonemic labels are essentially the same as the PronLex phonemic labels. To provide an example about the typical number of alternate pronunciations in  $L$ , roughly 17% of the entries in our JUPITER weather information lexicon contain more than one pronunciation.

### 3.3. Deriving the Phonological FST

To encode the possible pronunciation variants caused by phonological effects, we have developed a syntax for specifying phonological rules and a mechanism for converting these rules into an FST representation. In this approach phonological rules are expressed as a set of context-dependent rewrite rules. All of the phonological rules in our system have been manually derived based upon acoustic-phonetic knowledge, and upon actual observation of phonological effects present within the spectrograms of the data collected by our systems. The full set of phonological rules contains nearly 200 context-dependent rewrite rules. A full description of the expressive capabilities of the phonological rule syntax and the parsing mechanism for converting the rules into an FST can be found in [6].

To demonstrate some of the expressive capabilities of our phonological rule syntax, we now provide some examples of the phonological rules used in our system. Two example phonological rules for the phoneme /s/ are:

$$\{l m n ng\} s \{l m n w\} \rightarrow [epi] s [epi] ;$$
$$\{\} s \{y\} \rightarrow (s | sh) ;$$

The first rule expresses the allowed phonetic realizations of the phoneme /s/ when the preceding phoneme is an /l/, /m/, /n/, or /ng/ and the following phoneme is an /l/, /m/, /n/, or /w/. In these phonemic contexts, the phoneme /s/ can have an epenthetic silence optionally inserted before and/or after its phonetic realization of [s]. In the second rule the phoneme /s/ can be realized as either the phone [s] or the phone [sh] when followed by the phoneme /y/ (i.e., the /s/ can be palatalized).

To provide another example, the following rule accounts for the optional deletion of /t/ in a syllable suffix position when it is preceded by an /f/ or /s/ (as in the words *west* and *crafts*):

$$\{f s\} td \{\} \rightarrow [tcl [t]] ;$$

In this example the /t/ can contain a closure and a release, can be realized with an unreleased closure, or can be completely deleted.

To provide one more example, the following rule can be used to optionally insert a transitional [y] unit following an /iy/ when the /iy/ is following by another vowel or semivowel:

$$\{\} iy \{VOWEL r l w hh\} \rightarrow iy [y] ;$$

While this specific type of phonological effect is typically handled within the context-dependent acoustic models of a recognizer, this type of rule can be effective for providing additional detail to time-aligned phonetic segmentations. This can be especially helpful when utilizing automatically derived time-alignments for corpus-based concatenative synthesis.

### 3.4. Training the Pronunciation FSTs

To incorporate knowledge about the likelihoods of the alternate pronunciations encoded within the various component FSTs, we have implemented an EM training algorithm for training arbitrary FST networks [1, 2]. The full details of this algorithm are presented in [14]. When using the training algorithm is important to note that the size of the trained FSTs can be larger than the untrained FSTs. This is a result of a change in the FST topology which requires a determinization step during the creation of a trained FST in order to properly to account for the probability space (as explained in [14]).

The training algorithm can be used to train the individual component FSTs independently or jointly. When training the components independently (i.e.,  $tr(P) \circ tr(L) \circ tr(R)$ ) the likelihoods of specific phonological rules can be generalized across all words sharing these rules. When training the components jointly (i.e.,  $tr(P \circ L \circ R)$ ) the phonological rule probabilities are not shared across words and the likelihood of a particular realization of a phonological rule becomes dependent on the word in which it is applied. In previous experiments we found that joint training dramatically increased the size of the final static FST without improving the recognizer's accuracy [14].

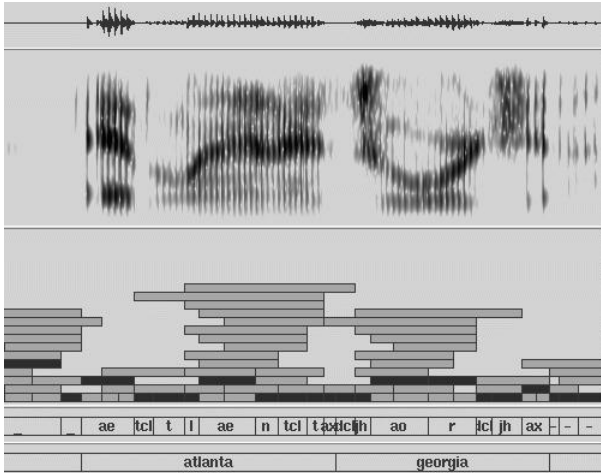
## 4. EXPERIMENTS & RESULTS

To investigate the effectiveness of using phonological rules, we evaluated three different sets of rules. These rule sets can be described as follows:

- Basic phoneme rules: This set of rules generates a one-to-one mapping of phonemes to phones. This is essentially the same as applying no rules except for the fact that we split stop and affricate phonemes into two phonetic segments to represent the closure and release portions of the phones with different models.
- Insertion and deletion rules: This set of rules augments the basic set with a collection of rules for inserting or deleting phonetic segments in certain contexts. This primarily includes the deletion of stop bursts or entire stop consonants, the reduction of stops to flaps, the insertion of epenthetic silences near strong fricatives, and the replacement of schwa-nasal or schwa-liquid combinations with syllabic nasal or syllabic liquid units.
- Full rule set: This set augments the insertion and deletion rules with a large set of rules for allophonic variation. This includes the introduction of new allophonic labels for stops and semivowels as well as rules for place assimilation and gemination.

By creating these three distinct sets of phonological rules we can first examine the effectiveness of introducing rules that account for phonetic insertions and deletions against the basic set of rules which do not allow substitutions and deletions. Figure 3 shows the phonetic alignment obtained by the SUMMIT recognizer using only the basic set of phonological rules on the same utterance presented earlier in Figure 1. An examination of the phonetic alignment in Figure 3 presents anecdotal evidence that the recognizer is not able to model the true sequence of phonetic events with the minimal set of phonological rules. This is particularly obvious in the word *atlanta* where the recognizer was forced to insert [t] releases for both /t/ phonemes despite the fact the speaker actually used the glottal stop allophone for the first /t/ and completely deleted the second /t/. Despite the poor phonetic transcription, the recognizer was still able to recognize this utterance correctly.

By adding rules to cover allophonic variation independent of rules which cover phonetic insertions and deletions, we can investigate the effectiveness of modeling allophonic variation implicitly using context-dependent acoustic models versus explicitly using context-dependent phonetic rewrite rules. Anecdotal evidence of the effectiveness of utilizing explicit rewrite rules to capture allophonic variation can be seen in the example in Figure 1 (on the first page of this paper). By examining the phonetic transcription in this figure, it can be observed that the recognizer successfully



**Fig. 3.** The output of a graphical interface displaying a sample waveform, its spectrogram, the hypothesized SUMMIT segment network with the best path segment sequence highlighted, the time-aligned phonetic transcription of the best path, and the time-aligned word transcription of the best path.

Phonological Rule Set	Word Error Rate (%)	
	Full Test Set	Clean Test Set
Basic Set	19.2	11.9
Ins./Del. Set	18.4	10.9
Full Rule Set	19.0	11.6

**Table 1.** Performance of JUPITER recognizer on the full test set and on the clean test set using three different sets of phonological rules and untrained FSTs.

identified the use of the glottal stop variant of /t/ at the beginning of *atlanta* and the use of fronted schwas at the end of both *atlanta* and *georgia*.

Our experiments were conducted using the SUMMIT recognizer trained specifically for the JUPITER weather information system, a conversational interface for retrieving weather reports and information for over 500 cities around the world [4, 19]. This recognizer has a vocabulary of 1915 words (excluding contracted or reduced forms) and includes 5 noise models for modeling non-speech artifacts and 3 models for filled pauses. The system was tested on a randomly selected set of 1888 utterances from calls made to JUPITER’s toll-free telephone line (we call this the *full test set*). Results are also reported for a 1293 utterance subset of the test data containing only in-vocabulary utterances with no non-speech artifacts (we call this the *clean test set*).

Table 1 contains the results of our experiments when using untrained versions of the component FSTs. As can be observed in the table, incorporating phonological rules for handling insertions and deletions of phonetic events resulted in a relative word error rate reduction of 8% (from 11.9% to 10.9%) on the clean test set. Over the full test set the error rate reduction was a more modest 4% (from 19.2% to 18.4%). These results demonstrate that standard context-dependent models by themselves are not sufficient for modeling contextual effects that cause the number of realized phonetic events to be different from the underlying canonical form.

Table 1 also shows that the additional rules added to create the full rule set actually degrade performance. These additional

Training Condition	Word Error Rate (%)	
	Ins./Del. Set	Full Rule Set
$P \circ L \circ R$	18.4	19.0
$tr(P) \circ L \circ R$	18.4	18.6
$P \circ L \circ tr(R)$	18.2	18.8

**Table 2.** Performance of JUPITER recognizer on full test set when training the phonological FST ( $P$ ) and the reductions FST ( $R$ ).

Phonological Rule Set	CD Acoustic Models		Full Static FST	
	Models	Gaussians	States	Arcs
Basic Set	1173	38349	39380	213550
Ins./Del. Set	1388	41677	45212	279340
Full Set	1630	45976	54641	386500

**Table 3.** Effect of phonological rules on size of context-dependent acoustic models and untrained static FST search network.

rules explicitly model allophonic variations which do not alter the number of phonetic events (such as palatalization, vowel fronting, etc.). This suggests that the context-dependent acoustic models are sufficient for modeling allophonic variation caused by phonetic context, and that the added complexity required to explicitly model these effects hinders the recognizer’s performance.

Table 2 shows the results on the full test set when various component FSTs are trained. By examining the first and second lines of Table 2, we see that training the phonological FST ( $P$ ) improves the performance of the system using the full rule set (from 19.0% to 18.6%). This is a similar improvement to past results we have obtained [14]. Unfortunately, training the  $P$  FST for the insertion/deletion rule set did not improve performance. Even when using an untrained  $P$ , the insertion/deletion rule set achieves a lower error rate than the full rule set using a trained  $P$ .

A comparison of the first and third lines of Table 2 shows that training the reductions FST ( $R$ ) provides modest improvements to both systems. We also attempted to train the lexical FST ( $L$ ) but did not achieve any performance improvement for either system from this training. We are also unable to report results for any system that combines a trained  $P$  with a trained  $R$  because the memory requirements for computing the composition of the individual component FSTs were prohibitively large. In past results using a slightly different pronunciation approach, where reductions were encoded directly within  $L$ , we were able to build a system which used both a trained  $P$  and a trained  $L$  within the final static FST to achieve a modest performance improvement [14]. We are currently investigating approximation methods to help reduce the size of the trained FSTs (and hence the memory requirements for building the final static FST).

To further demonstrate the effect that adding phonological rules has on the recognizer’s complexity, Table 3 shows the size of the recognizer for each of the three different rule sets in terms of the number of context-dependent acoustic models, the total number of Gaussian components in the acoustic model set, and the number of states and arcs in the pre-compiled untrained FST network. The number of acoustic models is determined for each rule set automatically based on phonetic context decision tree clustering. The number of Gaussians per context-dependent model is determined heuristically based on the number of training samples available for each model. The table shows a dramatic increase in the number of parameters required for the acoustic model and the complexity of the search space as additional phonological rules are added to the system.

## 5. PRONUNCIATION VARIATION FOR SYNTHESIS

Although this paper has focused on speech recognition, we have also utilized the same pronunciation framework in our group's concatenative speech synthesis system ENVOICE [17, 18]. When applying the framework for synthesis, the FST network is given a sequence of words and is searched in the reverse direction (i.e., in generative mode) to find an appropriate sequence of waveform segments from a speech corpus to concatenate. In generative mode the phonological rules can also be weighted in order to provide preferences for specific types of phonological variation. For example, the rules can be weighted to prefer reduced words, flaps and unreleased or deleted plosives in order to generate casual, highly-reduced speech. To generate well articulated speech the rules can be weighted to prefer unreduced words and fully articulated plosives.

## 6. SUMMARY

This paper has presented the phonological modeling approach developed at MIT for use in the segment-based SUMMIT speech recognition system. We have evaluated the approach in the context of the JUPITER weather information domain, a publicly-available conversational system for providing weather information. Results show that the explicit modeling of phonological effects that cause the deletion or insertion of phonetic events reduced word error rates by 8% on our clean, in-vocabulary test set. Our results also demonstrated that phonological effects which cause allophonic variation without altering the number of phonetic events can be modeled implicitly with context-dependent models to achieve better accuracy and less search space complexity than a system which models these effects explicitly within phonological rewrite rules.

Anecdotal visual examinations of the phonetic transcriptions generated using a full set of phonological rules also demonstrate a dramatic improvement in phonetic segmentation and classification accuracy during forced path recognition over a system using no phonological rules. This may not be of great consequence for word recognition, but it is vitally important for corpus-based concatenative synthesizers that rely on accurate automatically-derived time-aligned phonetic transcriptions in order to generate natural-sounding synthesized waveforms.

## 7. FUTURE WORK

While our work in this paper has been evaluated on spontaneous speech collected within a conversational system, we have found that human-human conversations tend to have even greater phonological variation than the human-machine data we have collected. Thus, we hope to evaluate our phonological modeling techniques on human-human corpora such as Switchboard or SPINE. We believe accurate modeling of phonological variation will have even greater benefits for these tasks.

While our paper has focused on modeling phonological variation within a sequence of independent FST layers, our group is also pursuing an approach which integrates the multiple layers within a single probabilistic hierarchical tree structure. This approach, called ANGIE, has the potential advantage of learning generalizations across the layers of the hierarchy which are currently modeled independently in our FST approach [13].

## 8. ACKNOWLEDGEMENTS

The authors would like to acknowledge the efforts of both Jim Glass, who developed the initial versions of the JUPITER recognizer and lexicon used in this paper, and Jon Yi who syllabified the PronLex dictionary. Jim and Jon are also the primary developers of the ENVOICE synthesizer discussed in this paper.

## 9. REFERENCES

- [1] A. Dempster, N. Laird and D. Rubin, "Maximum likelihood from incomplete data via the EM algorithm," *Journal of the Royal Statistical Society, Series B*, vol. 39, pp. 1-38, June 1977.
- [2] J. Eisner, "Parameter estimation for probabilistic finite-state transducers," *Proc. of the Annual Meeting of the Association for Computational Linguistics*, Philadelphia, PA, July 2002.
- [3] J. Glass, J. Chang, and M. McCandless, "A probabilistic framework for feature-based speech recognition," in *Proc. ICSLP*, Philadelphia, PA, October 1996.
- [4] J. Glass, T. J. Hazen, and I. L. Hetherington, "Real-time telephone-based speech recognition in the JUPITER domain," *Proc. ICASSP*, Phoenix, March, 1999.
- [5] S. Greenberg, "Speaking in shorthand - A syllable-centric perspective for understanding pronunciation variation," *Speech Communication*, vol. 29, pp. 159-176, November 1999.
- [6] I. L. Hetherington, "An efficient implementation of phonological rules using finite-state transducers," *Proc. EUROSPEECH*, Aalborg, Denmark, September 2001.
- [7] D. Jurafsky, *et al*, "What kind of pronunciation variation is hard for triphones to model?," *Proc. ICASSP*, Salt Lake City, UT, May 2001.
- [8] D. McAllester, L. Gillick, F. Scattone, and M. Newman, "Fabricating conversational speech data with acoustic models: A program to examine model-data mismatch," in *Proc. ICSLP*, Sydney, Australia, December 1998.
- [9] F. Pereira and M. Riley, "Speech recognition by composition of weighted finite automata," in *Finite-State Language Processing* (E. Roche and Y. Schabes, eds.), pp. 431-453, Cambridge, MA, MIT Press, 1997.
- [10] PronLex, COMLEX English Pronunciation Dictionary, available from <http://www ldc.upenn.edu>.
- [11] M. Riley, *et al*, "Stochastic pronunciation modelling from hand-labelled phonetic corpora," *Speech Communication*, vol. 29, pp. 209-224, November 1999.
- [12] S. Seneff, "The use of linguistic hierarchies in speech understanding," Keynote address at *ICSLP*, Sydney, Australia, November 1998.
- [13] S. Seneff and C. Wang, "Modelling phonological rules through linguistic hierarchies," in these proceedings.
- [14] H. Shu and I. L. Hetherington, "EM training of finite-state transducers and its application to pronunciation modeling," *Proc. ICSLP*, Denver, CO, September 2002.
- [15] N. Ström, I. L. Hetherington, T. J. Hazen, E. Sandness and J. Glass, "Acoustic modeling improvements in a segment-based speech recognizer," *Proc. IEEE ASRU Workshop*, Keystone, CO, December 1999.
- [16] G. Tajchman, E. Fosler and D. Jurafsky, "Building multiple pronunciation models for novel words using exploratory computational phonology," *Proc. EUROSPEECH*, Madrid, Spain, September 1995.
- [17] J. Yi, J. Glass and L. Hetherington, "A flexible, scalable finite-state transducer architecture for corpus-based concatenative speech synthesis," *Proc. ICSLP*, Beijing, China, October 2000.
- [18] J. Yi and J. Glass, "Information-theoretic criteria for unit selection synthesis," *Proc. ICSLP*, Denver, Colorado, September 2002.
- [19] V. Zue, *et al*, "JUPITER: A telephone-based conversational interface for weather information," *IEEE Trans. on Speech and Audio Processing*, vol. 8, no. 1, January 2000.