



Available at
www.ElsevierComputerScience.com
POWERED BY SCIENCE @ DIRECT®

SPEECH
COMMUNICATION

Speech Communication xxx (2003) xxx–xxx

www.elsevier.com/locate/specom

2 The use of subword linguistic modeling for multiple tasks 3 in speech recognition ☆

4 Stephanie Seneff *

5 *Laboratory for Computer Science, Massachusetts Institute of Technology, Spoken Languages System Group,*
6 *200 Technology Square, Room 643, Cambridge, MA 02139, USA*

7 Abstract

8 Over the past several years, I have been conducting research on subword modeling in speech recognition. The re-
9 search is most specifically aimed at the difficult task of identifying and characterizing unknown words, although the
10 proposed framework also has utility in other recognition tasks such as phonological and prosodic modeling. The
11 approach exploits the linguistic substructure of words by describing graphemic, phonemic, phonological, syllabic, and
12 morphemic constraints through a set of context-free rules, and supporting the resulting parse trees with a corpus-
13 trained probability model. A derived finite state transducer representation forms a natural means for integrating the
14 trained model into a recognizer search. This paper describes several research projects I have been engaged in, together
15 with my students and associates, aimed at exploring ways in which recognition tasks can benefit from such formal
16 modeling of word substructure. These include phonological modeling, hierarchical duration modeling, sound-to-letter
17 and letter-to-sound mapping, and automatic acquisition of unknown words in a speech understanding system. Results
18 of several experiments in these areas are summarized here.

19 © 2003 Published by Elsevier B.V.

21 1. Introduction

22 1.1. Background

23 Speech is first and foremost a *communicative*
24 signal. It is a complex encoding of linguistic mes-
25 sages for the purpose of conveying information

among humans who share the code. Speech sci- 26
entists have been studying various aspects of the 27
speech code for many decades, and engineers have 28
been involved in designing computer systems that 29
attain a certain degree of competence in under- 30
standing and rendering the code. 31

At the core of human communication is the 32
notion of “words” as the fundamental units. 33
Above the word level, it is apparent that words 34
group into phrases, and phrases group into higher 35
level units such as clauses and sentences. In addi- 36
tion to studies of how words are organized into 37
meaning, studies of the *substructure* of words in 38
multiple languages have revealed a number of 39
organizational principles (Scalise, 1986). The exact 40
specification of that substructure still eludes us, 41

☆ This research was supported by DARPA under contracts N66001-94-C-6040 and N66001-96-C-8526, monitored through Naval Command, Control and Ocean Surveillance Center, and contract NBCH1020002, monitored through the Department of the Interior, National Business Center, Acquisition Services Division, Fort Huachuca, AZ, and by the National Science Foundation under grant no. IRI-96187321.

* Tel.: +1-617-253-0451; fax: +1-617-258-8642.

E-mail address: seneff@sls.lcs.mit.edu (S. Seneff).

42 however, particularly for languages such as Eng- 82
 43 lish with a rich borrowing from other languages. 83
 44 The inventory of *phonemes* for any particular 84
 45 language can usually be enumerated quite specifi- 85
 46 cally. We are now also reasonably confident that 86
 47 the syllable exists as an intermediate layer between
 48 words and phonemes, although most speech rec-
 49 ognition systems make little or no use of this syl-
 50 lable layer. There is also the possibility of breaking
 51 words down into *meaning* units (i.e., morphemes),
 52 which may not necessarily align precisely with
 53 syllable units based strictly on phonology and
 54 sonority. The difficulty of defining exactly how the
 55 phonemes of a word might group themselves into
 56 natural subunits has been a major hurdle to the
 57 design of systems that utilize this intermediate
 58 structure.

59 1.2. *Historical context*

60 My interest in developing hierarchical struc-
 61 tures to represent word substructure was inspired
 62 by research conducted in the 1970s and 1980s by a
 63 team of researchers and students at MIT, dating
 64 back originally to Chomsky and Halle's theory of
 65 generative phonology (Chomsky and Halle, 1968).
 66 A doctoral thesis by Kahn led to the formalization
 67 of phonological phenomena with respect to sylla-
 68 ble structure (Kahn, 1976). In the 1980s, a team of
 69 researchers led by Allen developed a sophisticated
 70 letter-to-sound generation system called MITalk,
 71 based on a decomposition of words into meaning
 72 units called morphs (Allen et al., 1987). At the
 73 same time, Randolph developed formal rules to
 74 parse words into syllables, with the aim of for-
 75 mally encoding a distinctive-feature formalism
 76 (Randolph, 1989). Zue was also codifying his
 77 acoustic phonetic knowledge into formal, ordered,
 78 context-sensitive rules that could be utilized to
 79 expand lexical pronunciations for a speech recog-
 80 nition task (Zue, 1983). Church's doctoral thesis
 81 (1983) proposed applying context-free rules¹ to

parse syllables, in order to capture phonological
 effects such as flapping and palatalization, arguing
 that conditions for phonological phenomena could
 be encoded effectively in category names, thus
 avoiding explicit context dependencies.

1.3. *Motivation*

87
 88 While researchers have made significant pro- 88
 89 gress in speech recognition in the last two decades 89
 90 (Bacchiani and Ostendorf, 1998; Cook et al., 1998; 90
 91 Nguyen et al., 1995; Woodland and Young, 1993), 91
 92 there are still many remaining problems which 92
 93 could be addressed through the use of formal 93
 94 representation of the substructure of words. 94
 95 Nearly all speech recognition systems today are 95
 96 based on a simple model in which words are rep- 96
 97 resented explicitly in a lexicon encoding their 97
 98 phonemic realizations, and class *n*-gram language 98
 99 models provide linguistic constraint. One short- 99
 100 coming of such a representation is that unknown 100
 101 words are not formally represented, and therefore 101
 102 will be recognized as an acoustically similar sub- 102
 103 stitution of a known word, often adversely affect- 103
 104 ing the recognition of neighboring words as well 104
 105 (Hetherington, 1994; Hetherington and Zue, 105
 106 1993). Another limitation is the lack of a syllable- 106
 107 based framework for characterizing phonological 107
 108 rules, as well as the difficulty to capture such rules 108
 109 in an appropriate probability formulation.² The 109
 110 durational aspects of phonemes depend on their 110
 111 position within the syllable and the word, but this 111
 112 information is usually not available to a recog- 112
 113 nizer. Finally, the task of modeling an association 113
 114 between letters and their pronunciations is likely to 114
 115 benefit from knowledge of the linguistic context of 115
 116 each letter. 116

1.4. *Overview*

117
 118 Through my earlier and ongoing work in 118
 119 parsing words into meaning, via the TINA natural 119
 120 understanding system (Seneff, 1992), I have come 120
 121 to believe that a similar approach can be used 121

¹ A context free rule is a rule that rewrites a symbol generally into a sequence of zero or more symbols. A context-sensitive rule attaches conditions under which the symbol is permitted to be rewritten.

² Although one could argue that the use of triphone modeling does provide some *implicit* information about phonetic context.

122 effectively below the word level, yielding a parsimonious and trainable hierarchical representation of word substructure (Seneff, 1998). It is my belief that such representations may have significant advantages over a flatter structure, in that they should be capable of generalizing knowledge across similar contexts. I have subsequently investigated various ways in which such linguistic substructure can be utilized in speech recognition tasks. These investigations are predicated on a common theme that involves parsing words into their underlying linguistic constituents via a formal grammar, expressed through context-free rewrite rules. The resulting structural information is then augmented with a probability framework, where probabilities are determined by tabulating counts in parse trees obtained by parsing a large corpus of representative speech materials. A final optional step is to reformulate the trained parse trees as a finite state transducer (Hetherington, 2001), typically with inputs and outputs associated with the terminals and preterminals of the parse tree, respectively. This step then enables a straightforward mechanism for incorporating the linguistic models directly into a recognizer search.

147 The ideas discussed above have been formalized into a framework, called *ANGIE*, and several different topics of research have been investigated by members of the Spoken Language Systems group within this framework. These include phonetic recognition (Chung, 2001; Chung and Seneff, 1998; Lau, 1998; Lau and Seneff, 1998), hierarchical duration modeling (Chung, 1997; Chung and Seneff, 1997), sound-to-letter and letter-to-sound generation systems (Meng, 1995; Meng et al., 1996; Chung et al., in press; Seneff et al., 1996), unknown word detection and modeling (Chung, 2001; Mou et al., 2001; Parmar, 1997), and phonological modeling (Seneff and Wang, 2002). This paper will provide motivation for the approach we have taken, and will describe instances of all of the above applications within a common thread. While some of the investigations are on-going, it seems appropriate at this time to provide a detailed accounting of this research, partly in the hope that others might be inspired to pursue similar avenues of research.

2. *ANGIES* subword linguistic representation

169

170 Most of the work in speech recognition to date has been focused on the task of correctly producing the sequence of words that were spoken. The notion of characterizing any information beyond the word sequences is usually not treated as part of the explicit goal, although some amount of phonological and semantic knowledge is generally viewed as a necessary adjunct to success. Usually, each word is represented in the lexicon as a sequence of phonemes, and in some systems a phonological rule framework permits the expansion of lexical entries to explicitly account for phonological effects like flapping or devoicing (Cohen, 1989; Gauvain et al., 1993; Glass and Hazen, 1998; Weintraub et al., 1989). Typically the rules are precompiled into the lexicon, yielding an expanded lexicon of alternate pronunciations.

187 In order to address the issue of out-of-vocabulary (OOV) words, some recognition systems have included a generic model for unknown words as part of the recognizer's phonetic model. The approach typically adopts a generalized probabilistic subword model as a pronunciation model, such as a phone bigram, for the "word" OOV (Bazzi and Glass, xxxx). The word OOV then competes with known words, and the goal is that it would score better than known words for spoken out-of-vocabulary words, preventing the system from erroneously substituting a vocabulary entry with a similar pronunciation. An interesting example of a more sophisticated use of this technique is the work by Onishi et al. (2001), which developed distinct subword models for two different *classes* of unknown words: city-name and surname. In an evaluation experiment, they were able to achieve perfect disambiguation of the unknown word class whenever an unknown was detected, and with a slight improvement in overall recognition performance, when compared with a baseline that had no unknown word model.

210 In the *ANGIE* framework, we are interested in building a single subword linguistic model that can be effective for both the known and the unknown words. Thus, the purpose for building hierarchical structure below the word level is multifold. One main goal is to predict phone sequences of the

216 language without explicit ties to a particular
217 vocabulary. A bottom-up parsing procedure has
218 the important property that it supports significant
219 structure sharing among both in-vocabulary and
220 OOV words that begin with the same phone se-
221 quence. If words are further decomposed into
222 syllables, which then form the basic recognition
223 unit, even greater sharing is possible, since words
224 such as “retention” and “contention” can share
225 everything except their prefix in common syllable
226 nodes.

227 Exactly what linguistic knowledge should be
228 encoded in the ANGIE parse trees is open for de-
229 bate. In all of the experiments we have conducted,
230 syllable structure plays a critical role. At this
231 point, we have developed several grammars with
232 distinctive symbol sets at the preterminal and ter-
233 minal layers of the parse tree, but we have con-
234 sistently distinguished stressed and unstressed
235 syllables, which are further decomposed into on-
236 set, nucleus, and coda, according to standard syl-
237 lable theory (Selkirk, 1982; Kahn, 1976). This
238 choice of decomposition reflects in part our belief
239 that the position of a consonant within its syllable
240 plays a critical role in its phonetic expression. For
241 example, a /t/ in a syllable onset /st/ cluster is
242 unaspirated, whereas a /t/ would normally be
243 aspirated in onset position. Greenberg (1999) has
244 shown, through studies on a large corpus of hand
245 transcribed Switchboard data (Godfrey et al.,
246 1992), that 28% of consonants in *coda* position
247 were deleted, a rate that is substantially higher
248 than the rate for onset position.

249 The substructure that is captured in ANGIE’s
250 grammar rules includes morphology, stress, sylla-
251 ble structure, and phonological variants. Prob-
252 abilities are trained automatically from a parsed
253 corpus. We have used the approach of seeding on
254 phonetic transcriptions provided by automatic
255 alignment of training data using our segment-
256 based SUMMIT speech recognizer (Glass et al.,
257 1996; Glass and Hazen, 1998), which expands
258 idealized phonemic baseforms into phonetic
259 alternatives via formal phonological rules.

260 The shared probability model is important for
261 generalizing phenomena over similar contexts.
262 Rare words can benefit from observations of
263 common words that have the same local phonetic

264 environment. And words that are completely un- 264
265 known to the recognizer can be generated with a 265
266 non-zero probability by following the parse tree 266
267 fragments of words with localized equivalent pat- 267
268 terns. For example, “queen” can be decomposed 268
269 into the onset of “quick” and the rhyme of “seen.” 269

270 Parse trees in ANGIE are further characterized 270
271 by structural units that encode positional roles of 271
272 the syllables in a word. Thus, unstressed syllables 272
273 are identified as “prefix” if they begin a word, and 273
274 as “suffix” if they are a terminal unit carrying 274
275 syntactic information, such as “-ing” for present 275
276 participle or “-ness” marking a nominalization. 276
277 This additional information is beneficial not only 277
278 for further constraints but also because the posi- 278
279 tion of a syllable within a word impacts other as- 279
280 pects, such as the prosodic characteristics. For 280
281 example, in prepausal lengthening a final un- 281
282 stressed syllable is affected much more strongly 282
283 than an initial one (Chung and Seneff, 1997). 283

284 In addition to the parse framework, a set of 284
285 explicit subword lexical units can offer further 285
286 constraint in phonetic recognition tasks. One 286
287 possibility is to augment a phonetic recognizer 287
288 with a lexicon containing the inventory of all the 288
289 unique syllables present in a corresponding word 289
290 lexicon for the task. A syllable *n*-gram will then 290
291 provide additional language model support to 291
292 improve the quality of the phone or phoneme 292
293 graph being proposed as outputs of the recognizer. 293
294 In addition to simple syllable units, we have also 294
295 investigated the use of more detailed units which 295
296 we call “morphs,”³ essentially syllables marked 296
297 for both their spelling and positional information 297
298 within the word. Another unit above the syllable is 298
299 the metrical foot (Hayes, 1995), which consists of a 299
300 stressed syllable and zero or more adjacent un- 300
301 stressed syllables. This unit of recognition, which is 301
302 intermediate between phonemes and words, pro- 302
303 vides a convenient compromise in yielding fairly 303
304 strong constraint while still supporting substantial 304
305 coverage of novel words and disfluencies in con- 305
306 versational speech. 306

³ This follows roughly the definition given in (Allen et al., 1987, p. 24), which is a representation of morphological units such as prefix and root that is also tied to the word’s spelling.

307 In *ANGIE*, we currently represent our lexicon in
 308 two tiers—words are entered as sequences of
 309 morphs, and morphs are in turn entered as se-
 310 quences of phonemes. We currently distinguish for
 311 English five different possible morph positions:
 312 prefix, stressed root, unstressed root, “dsuf” and
 313 “isuf”.⁴ Context-free rules encode positional
 314 constraints for the morph units—for example,
 315 unstressed root always follows immediately after
 316 stressed root, and isuf’s are always terminal.

317 As mentioned previously, it is often not obvious
 318 where to place syllable boundaries in English
 319 words. There are many cases of ambisyllabicity, as
 320 in the word “connect” where it is not clear whether
 321 the intermediate consonant belongs with the pre-
 322 ceding or following syllable. Placement of the
 323 boundary can also be influenced by the underlying
 324 morphology—when there is a clear inflectional
 325 ending our policy has been *not* to shift the terminal
 326 consonant of the root into onset position, even
 327 though this would be in accord with a maximal-
 328 onset rule. Hence “dancing” becomes “danc-ing”
 329 rather than “dan-cing”. Often we introduce a
 330 double consonant phonemically as a means of
 331 implementing explicit ambisyllabicity, which re-
 332 duces via a gemination rule to a single phonetic
 333 realization. Hence, “connect” becomes “con-nect”
 334 with two /n/ phonemes at the phonemic layer
 335 reducing to one at the phonetic layer. This makes
 336 the boundary between the word-internal syllables
 337 behave analogously to boundaries between word
 338 sequences, as in “on next” or “seven nine.” Such
 339 lexicalized geminations are nearly always associ-
 340 ated with a spelling that includes a doubleton let-
 341 ter, such as the “nn” in “connect.”

342 2.1. Example parse tree

343 One of the main goals of *ANGIE*’s modeling is
 344 to provide letter-to-sound and sound-to-letter
 345 mappings, and, particularly for this purpose, we
 346 have found it beneficial to provide a pair of
 347 grammars with a shared superstructure but two

348 distinct sets of rules mapping preterminals to termi-
 349 nals: one expecting phonetic units as the termi-
 350 nals and the other expecting graphemics. The
 351 preterminal layer contains the phonemic sequence
 352 exactly matched to the entries in the morphs of the
 353 two-tiered lexicon. The terminals are either the
 354 letters of the spelling of the word or the phones of
 355 the particular spoken realization. Thus letter-to-
 356 sound and phonological rules are licensed on the
 357 preterminal-to-terminal mappings. The upper lay-
 358 ers capture syllabification, morphology, and stress.

359 Example parse trees in *ANGIE* for the word
 360 “commission” are given in Figs. 1 (letter termi-
 361 nals) and 2 (phone terminals). The lexical repre-
 362 sentation of the word consists of a prefix (com-) a
 363 stressed root (mis+) and an inflectional suffix
 364 (= sion). Phonemically, there are both a final /m/
 365 for the prefix and an onset /m/ for the root. These
 366 geminate in the phonetic realization into a single
 367 [m].⁵ Similarly, the “mis+” unit ends phonemi-
 368 cally with an /s/. The /s/ is palatalized to a [sh] at
 369 the phonetic level, with the onset /sh/ of the
 370 “= sion” marked as deleted. Fig. 3 illustrates how
 371 sharing of subword units can be achieved, using
 372 the examples “mis+” and “= sion.”

373 2.2. Lexicon creation

374 *ANGIE* relies heavily on the availability of a
 375 specifically prepared two-tiered lexicon, in which
 376 words are represented in terms of their underlying
 377 morphs. We first obtained, through careful hand-
 378 editing, a seed lexicon of some 10,000 words, de-
 379 rived from the common words of the Brown cor-
 380 pus (Kucera and Francis, 1967) augmented with
 381 words from some of our conversational domains
 382 such as ATIS (Zue et al., 1991) and Jupiter (Glass
 383 and Hazen, 1998). We have since converted all the
 384 common words of Pronlex⁶ into *ANGIE*’s lexical
 385 format (Parmar, 1997). We have utilized a semi-
 386 automatic process which first parses the letters of

⁵ [-m] is a code for “deleted in the context of preceding [m]”.

⁶ A pronunciation lexicon for the words in the Comlex lexicon, produced and distributed by the Proteus Project at New York University, under the auspices of the Linguistic Data Consortium (see <http://www ldc.upenn.edu>).

⁴ “dsuf” roughly corresponds to “derivational suffix,” and “isuf” to “inflectional suffix,” but we sometimes violate strict conventions for pragmatic reasons.

word						
pre		sroot			isuf	
uonset	nuc	onset	nuc_lax+	coda	uonset	nuc
k!	em	m!	ih+	s	sh!	en
c	o	m	m2	i	s2	i o n
com-		mis+			=sion	

Fig. 1. ANGIE parse tree for the word “commission,” with letters as the terminals. An aligned sequence of morphs is shown below the parse tree. Note: “!” denotes onset position and “+” marks stress. The second letter in a doubleton is specially tagged for additional constraint (m2, s2).

word						
pre		sroot			isuf	
uonset	nuc	onset	nuc_lax+	coda	uonset	nuc
k!	em	m!	ih+	s	sh!	en
kcl	k	ax	m	-m	ih	sh -sh ax n
com-		mis+			=sion	

Fig. 2. ANGIE parse tree for the word “commission,” with phones as the terminals. An aligned sequence of morphs is shown below the parse tree. The highlighted entries illustrate units involved in the trigram language model as applied to the bottom-up prediction of the preterminal layer.

Word Lexicon

commission com- mis+ =sion
mister mis+ ter
mansion man+ =sion

Morph Lexicon

com- k! em
man+ m! ae+ n
mis+ m! ih+ s
sion sh! en
ter t! er

Fig. 3. Selected entries from a word and morph lexicon for ANGIE.

387 each word into a set of hypothesized phonemic
388 alternatives, and then parses the phonetic units as
389 provided by Pronlex into phonemes, constrained
390 by the choices produced by the letter-parsing step.
391 Of course the automatic procedures are not error-
392 free, so extensive hand correction is required to
393 perfect the lexicon.

394 We hope to use the resulting morph lexicon as a
395 basis for a generic morph-based recognizer for
396 general English. A phonological model can then be

397 trained on any large corpus of spoken utterances.
398 There would still be some possibility of unseen
399 morphs in new material, but these would likely be
400 covered generatively by the rule base. Such a lex-
401 icon is also useful for training a reversible letter-to-
402 sound system. Ultimately, we would like to aug-
403 ment it with additional information such as part-
404 of-speech, and perhaps add a feature propagation
405 mechanism to ANGIE’s framework to utilize such
406 features, similar to the one developed for the TINA
407 natural language understanding system (Seneff,
408 1992).

2.3. Probability model

409 In ANGIE, a parse tree is obtained for each
410 word by expanding the rules of a carefully con-
411 structed context-free grammar. The grammar is
412 intentionally arranged such that every parse tree
413 lays out as a regular two-dimensional grid, as
414 shown in Fig. 4. Each layer is associated with a
415 particular aspect of subword structure: migrating
416 from morphemics to syllabics to phonemics to
417 phonetics at the deepest layer. Although the rules
418 are context free, context dependencies are captured
419 through a superimposed probability model. The
420 particular choice for the probability model was
421 motivated by the need for a balance between suf-
422 ficient context constraint and potential sparse data
423 problems from a finite observation space. We were
424 also motivated to configure the probability model
425 such that it would be causal, with strong locality,
426 for practical reasons having to do with the nearly
427 universal left-to-right search path in recognition
428 tasks, as well as the convenience of providing arc
429

sentence								
word								
sroot		uroot			sroot2			
nuc_lax+	coda	uonset	nuc	onset	lnuc+	lcoda		
ih+	n	t!	r	ow	d!	uw+	s	
ih	n	-n	rx	-rx	dcl	d	uw	s
in+		tro			duce+			

Fig. 4. ANGIE parse tree for the word “introduce,” showing phonological rules expressed in preterminal-to-terminal mappings. The morph sequence is shown below the terminal phones.

430 probabilities for a finite state transducer (FST)
431 representation (Hetherington, 2001).

432 Given these considerations, the probability
433 formulation we have developed for ANGIE can be
434 written as follows:

$$P(C_i|C_{i-1}) = P(a_{i,0}|C_{i-1}) \prod_{j=1}^{N-1} P(a_{i,j}|a_{i,j-1}, a_{i-1,j}) \quad (1)$$

436 where C_i is the i th column in the parse tree and
437 $C_i = \{a_{i,j}, 0 \leq j \leq N\}$, and $a_{i,j}$ is the label at the j th
438 row of the i th column in the two-dimensional
439 parse grid.⁷ In words, each phone is predicted
440 based on the entire preceding column, and the
441 column probability is built bottom-up based on a
442 trigram model, considering both the child and the
443 left sibling in the grid. These probabilities are ac-
444 quired by tabulating counts in a large corpus of
445 parsed sentences, mapping words to their corre-
446 sponding phonetic realizations. This process will
447 become clearer when we give an explicit example in
448 the next section.

449 ANGIE’s language model, while restricted to
450 phone-to-phone transitions, is very powerful, and
451 captures generic linguistic knowledge of English
452 while a partial word is under construction. We
453 have determined empirically that, within the ATIS
454 flight information domain, ANGIE is able to
455 achieve a significantly lower perplexity on unseen
456 data than a phone trigram similarly trained (Lau
457 and Seneff, 1997). Once a word is completed,
458 higher level language models can be incorporated
459 as well (e.g., syllable/word n -grams).

460 2.4. Phonological rule expression

461 ANGIE’s ability to encode and generalize pho-
462 nological rules is best illustrated through an
463 example. Consider the parse tree shown in Fig. 4
464 for the word “introduce” pronounced casually as
465 “innerduce.” The two special phones $[-n]$ and
466 $[-rx]$ are “deletion” phones, meaning that they
467 occupy no temporal space and have no acoustic
468 model. The deletion category is tied to the pre-
469 ceding phone’s identity. The grammar developer

would specify that $/t/$ can be realized as “ $[-n]$ ”,
meaning “ $/t/$ in onset position can be deleted after
 $[n]$.” The probability model captures the important
context conditions—falling stress and following
schwa. The deletion of the $/ow/$ is predicated on
the realization of the preceding $/r/$ as a retroflexed
schwa ($[rx]$).

Fig. 5 illustrates the context conditions that are
learned, with regard to this t -deletion rule. The
column above the $[n]$ encodes coda position in a
stressed syllable. It predicts a deletion after $[n]$
with no awareness of which phoneme actually
follows. The trigram column-building step decides
which phoneme was deleted. Other possibilities
would be $/t/$, $/d/$, $/d!/$, and $/n!/$. The training pro-
cedure would collapse together the $/t/$ deletion here
with other similar environments, such as “in-
tegrate,” “cantaloupe,” “entertain,” “Santa
Clause,” “hunter,” and “pantyhose.” The column
above the $[-n]$ would learn through training that it
is rarely followed by anything other than $[ax]$, $[rx]$,
and $[ix]$. The system would thus learn from
examples that the right context must be a schwa,
but it could be front, back or retroflexed. This
“fact” was not informed by any rule, but rather
discovered from observation of training data.

2.5. Spellnemes

In the original grammars we developed for
ANGIE, we adopted the point of view that there
would be two parallel grammars with identical

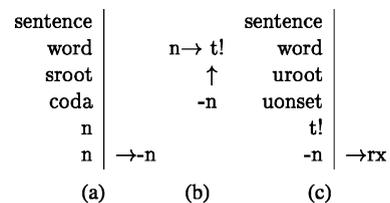


Fig. 5. Schematic of probability model in ANGIE, and its accounting of the context conditions for t -deletion in words such as “introduce.” In (a) and (c) are the two column contexts for the predictions of the phones “ $-n$,” symbolizing deletion after $/n/$, and “ rx ,” a retroflexed schwa. (b) Illustrates the bottom-up trigram prediction of the deleted phone’s parent phoneme, which in the example is “ $t!$,” a $/t/$ in onset position (see text for discussion).

⁷ j indexing begins at the bottom of the column.

500 parse tree superstructure, but terminating in pho-
501 netic units on the one hand and in graphemic units
502 (spellings) on the other hand. The phonetic ter-
503 minals would characterize possible phonological
504 rules, such as alveolar flapping, vowel reduction,
505 anticipatory palatalization, etc. (Zue, 1983). The
506 graphemic units would encode letter-to-sound
507 rules of English, with the terminals consisting of
508 the letters of the English alphabet, sometimes
509 grouped to form natural paired units such as “th”
510 or “ng.”

511 In more recent work, however, (Chung,
512 2000a,b; Chung, 2001) we have begun to explore
513 potential benefits of introducing a new unit type
514 which we refer to as a letter-phoneme, or “spell-
515 neme.” These are units which tie together letters or
516 letter sequences with their corresponding pronun-
517 ciation in a single symbol string. For example, the
518 spellneme unit, “a_x+” symbolizes a short /a/,
519 pronounced as in the word “cat,” whereas the
520 “a_l+” symbol string refers to a long /a/ as in
521 “cape.” Other symbols distinguish “soft” and
522 “hard,” for consonants such as ‘c’ (compare
523 “cent” with “car”), etc. In general, we expect a
524 sequence of spellnemes to encode both the spelling
525 and the pronunciation of the word.

526 An advantage of these spellneme units is that
527 they provide greater constraint, in that they are
528 more specific than either letters or phonemes. We
529 expect this to translate into better performance,
530 not only on the sound-to-letter task, but also on
531 other tasks such as phonetic recognition, due to
532 the richer language model, which has in fact been
533 borne out by experiments (Chung and Seneff,
534 1998). There is also the notational convenience
535 that both the spelling and pronunciation can be
536 derived through simple string manipulations. The
537 “commission” example, with spellnemes as pre-
538 terminals and phones as terminals, is illustrated in
539 Fig. 6.

540 **3. Incorporating ANGIE into recognition tasks**

541 The ANGIE parsing framework provides pow-
542 erful mechanisms for learning important aspects of
543 subword structure. However, parsing is in general
544 a computationally expensive process, and so it

word									
pre			sroot				isuf		
uonset	nuc		onset	nuc_lax+	coda	uonset	nuc		
c_k!	om		m!	i_x+	s	si!	on		
kcl	k	ax	m	-m	ih	sh	-sh	ax	n
com-			mis+				=sion		

Fig. 6. ANGIE parse tree for the word “commission,” with spellnemes as the preterminals (highlighted), and phones as the terminal units. This is to be compared with the corresponding parse tree in Fig. 2.

545 becomes important to consider ways of encapsu- 545
546 lating the results of ANGIE parse training into a 546
547 speech recognition task while still preserving close 547
548 to real-time performance. Fortunately, the SUM- 548
549 MIT segment-based recognizer has been formu- 549
550 lated in terms of a FST framework (Hetherington, 550
551 2001), which provides the opportunity to encode 551
552 complex linguistic knowledge and embed it in the 552
553 core recognizer search engine. The dominant ap- 553
554 proach we have taken is to reconfigure ANGIE 554
555 parse trees as a “column bigram” (Chung, 2000b), 555
556 where a “column” is a unique path from a termi- 556
557 nal node up to the root node. An assignment of the 557
558 probability associated with a transition from one 558
559 column to the next can be computed directly from 559
560 the ANGIE parse framework. In most of our 560
561 experiments, we have preserved only a portion of 561
562 the information in the ANGIE parse tree via the 562
563 input and output symbols associated with the FST. 563
564 These simplifications were thought to be necessary 564
565 both because of the complexity inherent in an FST 565
566 representing the full ANGIE parse space, and be- 566
567 cause the full parse space would likely lead to an 567
568 impractical FST size (but see (Mou et al., 2001) for 568
569 a more general FST solution). 569

570 For many applications, the full linguistic model 570
571 obtained from ANGIE is correlated with other 571
572 linguistic knowledge encoded in higher level lan- 572
573 guage models introduced either in the same stage 573
574 or in later stages of a multi-stage recognizer. An 574
575 intended goal is to influence the search through 575
576 language model constraints before the lexical en- 576
577 tries have been retrieved. Once a proposed word is 577
578 known, and at the point where the word *n*-gram 578
579 score is being introduced, the linguistic component 579
580 of the ANGIE contribution to the word can be 580
581 subtracted out, leaving behind only the phono- 581

582 logical component. In other configurations, a more
583 precise formulation using Bayes' formula to re-
584 move from each column-column score all contri-
585 butions except the phoneme-to-phone probability
586 assignment, where the intent is to use ANGIE
587 strictly as a phonological rule model. The latter
588 technique is particularly effective when the speech
589 corpus used for phonology training is not well
590 matched to the application domain, as is often the
591 case when a new domain is being launched. In such
592 cases, the upper layers of ANGIE's parse tree
593 would obtain probability training from an inap-
594 propriate linguistic model, which can be elimi-
595 nated by the normalization procedure. The
596 likelihoods of the phoneme-to-phone mappings
597 should be independent of the domain, and thus
598 can be used to capture phonological rule proba-
599 bilities generically (Seneff and Wang, 2002).

600 In the remainder of this section we will touch on
601 several applications where the ANGIE linguistic
602 hierarchy has been found to be useful for
603 improving speech recognition performance. We
604 begin with the task of duration modeling in rec-
605 ognition, evaluated in the ATIS flight information
606 domain (Chung and Seneff, 1997). This is followed
607 by an experiment in acquiring probabilities on
608 phonological rule productions, conducted in the
609 Mercury flight reservation domain (Seneff and
610 Polifroni, 2000), but making use of a large training
611 corpus from the Jupiter weather domain (Seneff
612 and Wang, 2002). Next, we address the issue of
613 detecting and accounting for unknown words,
614 through experiments conducted in the Jupiter do-
615 main (Chung, 2000a,b). We conclude with a dis-
616 cussion of our research in the highly related topic
617 of automatic new-word acquisition (Chung and
618 Seneff, 2002). Due to space restrictions, each of the
619 topics is described only briefly. The interested
620 reader is referred to the literature for details of the
621 experiments.

622 3.1. Duration modeling

623 ANGIE's parse trees can provide access to
624 intermediate structures within words, which can be
625 useful for characterizing prosodic information.
626 Thus far we have only attempted to characterize
627 prosody through *timing* measures. However, we

628 have found that significant improvements in both
629 phonetic recognition and word spotting can be
630 gained through the use of relative duration models
631 relating parents to children at all layers of an
632 ANGIE parse tree (Chung, 1997; Chung and Se-
633 neff, 1997). The approach adopted involved normal-
634 izing the duration of each constituent in the
635 parse tree with respect to its particular *children*,
636 and then measuring the portion it occupies of its
637 *parent's* total duration. The procedure propagates
638 to the top of the tree to yield a word-by-word
639 speaking rate parameter, which can then be folded
640 back into the phonemic layer to tighten the dis-
641 tributions on absolute phoneme duration. This too
642 leads to improved overall recognition.

643 To quantitatively assess the effectiveness of the
644 hierarchical duration model, an experiment was
645 conducted on phonetic recognition in the ATIS
646 flight information domain, where the sophisticated
647 duration models were benchmarked against a
648 phonetic recognizer configuration which made use
649 of the raw phone duration as a feature in the
650 phone-based acoustic models. The enhanced sys-
651 tem augmented the standard system with two sets
652 of Gaussian models, as suggested above: relative
653 duration models, across the entire parse tree, and
654 models for absolute phoneme duration, normal-
655 ized with respect to the estimated word-by-word
656 speaking rate parameter. It was found that the
657 performance gains attributable to the hierarchical
658 duration models were stronger when the linguistic
659 models included a richer knowledge base, with the
660 best gains yielding improvements in phonetic error
661 rate from 29.7% to 27.4%.⁸

662 The word spotting experiments (Lau, 1998)
663 were also conducted in the ATIS domain, where
664 the task was to detect all city names in the user
665 utterances, treated as keywords. Results were re-
666 ported in terms of a "figure of merit" (FOM),
667 derived by integrating over a receiver operator
668 characteristic (ROC) curve, which gives detection
669 rate as a function of false alarm rate. The addition
670 of the hierarchical duration model to the scores for
671 the keywords yields performance gains on the

⁸ For the details of these experiments, please see Chung (1997).

672 FOM from 89.3% to 91.6%. Details of this
673 experiment can be found in (Chung, 1997; Lau,
674 1998).

675 We believe that this direction of research has
676 many as yet unexplored branches, both in terms of
677 incorporating hierarchies above the word level and
678 in incorporating other prosodic measures such as
679 fundamental frequency and energy.

680 3.2. Acquiring phonological rule probabilities

681 This section describes a set of experiments we
682 have conducted, aimed at acquiring a probability
683 model to support phonological rules describing the
684 mappings from the phonemic baseforms of a lex-
685 icon to the actual phonetic realizations in sponta-
686 neous speech. In this case, we are only interested in
687 the component of ANGIE's probability model that
688 predicts the terminal phone unit of each sub-
689 sequent column. It is not enough to simply discard
690 the predictions of the chain of parents moving up
691 the right hand column, but rather they must be
692 considered as contributing to the conditioning
693 context for the terminal phone.

694 The procedures we have adopted appear com-
695 plex, but are relatively straightforward to execute,
696 given a SUMMIT recognizer with an associated
697 phonological rule set, an ANGIE grammar with
698 ANGIE's phonemes in the preterminal layer and
699 SUMMIT's phones in the terminal layer, a lexicon
700 of words in the domain, with baseforms available
701 in both SUMMIT's phonemic units and in ANGIE's
702 phonemic units, and a large corpus of utterances
703 for training.

704 The approach we have taken, then, is to start
705 with a SUMMIT recognizer, complete with its
706 standard set of phonological rules, which, when
707 applied to the SUMMIT baseforms, yields a finite
708 state transducer specifying all the phonetic vari-
709 ants possible for each word in the lexicon, but in
710 the process, losing the mapping from phones to
711 phonemes (Hetherington, 2001) (this transducer
712 inputs phones and outputs words). We then insert
713 a column bigram FST mapping phones to ANGIE
714 phonemes, along with an ANGIE baseforms FST
715 to map from ANGIE phonemes to words. The
716 column bigram FST will contain probabilities
717 computed by training ANGIE on an observation

space obtained by parsing the phonetically aligned 718
corpus. The process can be iterated. 719

The ANGIE model intentionally captures both 720
phonological and linguistic aspects of the lan- 721
guage, such as the frequency of different syllable 722
onset patterns. However, for the purpose of 723
modeling the likelihood of the phonological vari- 724
ants, the linguistic contribution to the probability 725
model needs to be removed. Specifically, our 726
phonological model is designed to predict each 727
subsequent phone, using the entire previous col- 728
umn and the column above the new phone as the 729
context. This can be achieved by essentially 730
inverting the probability model of the right col- 731
umn such that the predictor focuses totally on the 732
prediction of $a_{i,0}$, the *phonetic* realization associ- 733
ated with the right column. In practice, this means 734
summing over all observed instances of $a_{i,0}$ fol- 735
lowing C_{i-1} to compute a total probability for each 736
particular set of $\{a_{i,j}, j > 0\}$, i.e., each unique up- 737
per column. This sum then becomes the denomi- 738
nator in a normalization step. Thus, the 739
probability of the right column's phone is mod- 740
elled as the probability of the phone and the upper 741
column, normalized by the total probability of the 742
upper column, given the left column: 743

$$\begin{aligned} &P(a_{i,0}|C_{i-1}, \{a_{i,j}, j > 0\}) \\ &= \frac{P(a_{i,0}, \{a_{i,j}, j > 0\}|C_{i-1})}{P(\{a_{i,j}, j > 0\}|C_{i-1})} \\ &= \frac{P(C_i|C_{i-1})}{\sum_{a_{i,0}} P(a_{i,0}, \{a_{i,j}, j > 0\}|C_{i-1})} \\ &= \frac{P(C_i|C_{i-1})}{\sum_{a_{i,0}} P(C_i|C_{i-1})} \end{aligned} \quad (2)$$

$$\begin{aligned} &P(a_{i,0}|C_{i-1}, \{a_{i,j}, j > 0\}) \\ &= \frac{P(a_{i,0}|C_{i-1}) \prod_{j=1}^{N-1} P(a_{i,j}|a_{i,j-1}, a_{i-1,j})}{\sum_{a_{i,0}} P(a_{i,0}|C_{i-1}) \prod_{j=1}^{N-1} P(a_{i,j}|a_{i,j-1}, a_{i-1,j})} \end{aligned} \quad (3)$$

To acquire the probability model for the col- 746
umn bigram, the corpus is first processed through 747
forced alignment using standard methods avail- 748
able in SUMMIT, to yield a phonetic transcription 749
associated with each utterance. The ANGIE 750
grammar is then trained on parse trees associated 751
with the corpus. Next, the corpus is reparsed, but 752

753 this time using the trained grammar, and with the
754 intent of producing a column bigram mapping
755 phones to ANGIE phonemes, removing the lin-
756 guistic predictions through the procedures de-
757 scribed above. An attractive aspect of this
758 approach is that the corpus for training does not
759 have to be restricted to the intended application
760 domain, since the language model component of
761 the column bigram probability space has been
762 completely removed. A summary of the steps in
763 this procedure is given in Fig. 7.

764 To demonstrate the viability of this approach,
765 we have trained the system on a corpus consisting
766 of a mixed set of 80,700 utterances⁹ from the
767 Jupiter weather domain and 13,800 utterances
768 from the Mercury flight reservation domain. The
769 trained model was then tested on an independent
770 test set of 848 utterances exclusively from the
771 Mercury domain. Results are summarized in Table
772 1. We were able to realize a significant reduction in
773 word error rate, compared with the SUMMIT
774 baseline system, when training on a trigram lan-
775 guage model. Perhaps more significantly, when we
776 evaluated on *understanding* error rates, the per-
777 formance improvement was even greater: concept
778 error rate dropped from 11.9% in the baseline
779 system to 10.4% with the phonological probability
780 modeling, suggesting that the probabilities are
781 differentially helping the content words. For fur-
782 ther details concerning these experiments, please
783 see Seneff and Wang (2002).

784 3.3. Modeling unknown words

785 One of the most significant applications for
786 ANGIE in subword modeling is both the detection
787 and the characterization of new words. Our ap-
788 proach to this problem is predicated on the notion
789 that the known words can serve as a model for the
790 unknown words: by decomposing words into their
791 linguistic constituents, novel combinations of these
792 constituents can yield representations for the un-
793 known words. To fully characterize new words,
794 one needs both their phonemic and their graphe-
795 mic representations. Thus, if a subword hierarchy

796 can capture both of these aspects, then it has utility
797 to provide both constraint and valuable linguistic
798 knowledge.

799 In our early work on sound-to-letter and letter-
800 to-sound tasks (Seneff et al., 1996) we formulated a
801 grammar whose terminals were the letters of the
802 spelled form, with the preterminals encoding
803 phonemic information. We conducted experiments
804 using the Brown corpus and obtained competitive
805 results on the letter-to-sound task (91.5% phoneme
806 accuracy on an unseen test set). For the sound-to-
807 letter task, we modified the search such that the
808 preterminal phonemes were provided as inputs and
809 the parsing process then predicted the most likely
810 letter sequence corresponding to these phonemic
811 specifications. This strategy gave a reasonable
812 performance (89.2% letter accuracy on an inde-
813 pendent test set) but it assumed a perfect phonemic
814 transcription as the input sequence, and it still falls
815 far short of the performance level necessary for
816 new word enrollment.

817 In later experiments conducted by Chung
818 (2000a), we attempted the much more ambitious
819 *waveform-to-letters* task. These experiments were
820 conducted within the Jupiter weather domain
821 (Glass and Hazen, 1998), and we selected as a test
822 set a set of utterances that contained unknown city
823 names. The task therefore involved first identifying
824 the presence of the portion of the speech waveform
825 associated with the unknown city, subsequently
826 proposing a possible spelling for that city.

827 We attacked this problem through a two-stage
828 procedure, where the first stage utilized subword
829 structure mainly as a language model in support of
830 phonetic recognition, and the second stage in-
831 volved parsing the resulting phone graph into a
832 sequence of proposed known and unknown words.
833 For both stages, we utilized a grammar that con-
834 tained SUMMIT phonetic units as the terminals
835 and *spellnemes* as the preterminals. Thus we
836 mapped phones directly to units that encode both
837 the phonemic and the graphemic information. This
838 approach is distinguishable from the approaches
839 addressed by Bazzi and Glass (xxxx) and Onishi
840 et al. (2001) in that a detailed sound-to-letter sys-
841 tem is embedded in the linguistic model charac-
842 terizing the unknown words. We anticipate that
843 linguistic constraint achieved as a side-effect of the

⁹ Since we have available a much larger corpus from Jupiter.

1. Obtain corpus of phonetic alignments using SUMMIT.
2. Train ANGIE grammar on phonetic alignments.
3. Obtain column bigram FST by parsing corpus using trained grammar and inverting right-column production probabilities.
4. Compose diphone FST from SUMMIT with ANGIE column bigram to yield FST mapping phonetic models to ANGIE phonemes.
5. Compose with ANGIE's lexicon.
6. Add standard word *n*-grams as language models.

Fig. 7. Steps in training phonological rule probabilities using ANGIE and SUMMIT.

Table 1

Speech recognition (WER) and understanding (CER) performance for telephone quality speech collected within the Mercury flight reservation domain

No. of Utterances	WER (%)		CER (%)	
	Baseline	+ Angie PM	Baseline	+ Angie PM
848	17.3	16.3	11.9	10.4

The system which utilized an ANGIE pronunciation model (+ Angie PM) is contrasted with a baseline system that utilized the same set of phonological rules but had no probability model for the alternate pronunciations.

844 statistical sound-to-letter mappings will provide a
845 richer linguistic model for the unknown words,
846 with the additional benefit (and goal) of providing
847 a full characterization of the unknown words in
848 terms of both their spellings and pronunciations.

849 Since all of the linguistic information in the first
850 stage was ultimately discarded, we were not re-
851 quired to represent this information perfectly. In
852 fact, for certain words, we felt that the training
853 process would generalize better if we discarded
854 rare forms, allowing their more common homo-
855 morphs to stand in for them. A good example to
856 clarify this point is the word "champagne," whose
857 unique spelling is very difficult to predict from
858 observations of other words in English. An Eng-
859 lish sound-to-letter system not explicitly trained on
860 this word would likely produce something like
861 "shampain." We decided to formalize such "mis-
862 takes" by introducing these odd spellings inten-
863 tionally in order to reduce the perplexity of the
864 task and better generalize the models. We realized
865 further that even the boundaries of the words were

not necessary to preserve in the first stage of our 866
system. We therefore decided to license a realign- 867
ment of word boundaries by reorganizing the syl- 868
lables into foot-like units, each of which contained 869
a single stressed syllable and zero or more un- 870
stressed syllables on either side. Furthermore, we 871
developed an iterative procedure which realigned 872
these foot-like units with each iteration. Each 873
realignment would support the same phonetic se- 874
quence as the original but with a reduced per- 875
plexity. The result of all of this training was a 876
significant net reduction in the size of the FST, 877
along with a reduction in the perplexity of the task, 878
both of which are positive outcomes. The foot-like 879
pseudo-words were supported by a standard class 880
trigram, to yield additional constraint in the 881
search. 882

The second stage of this system parsed the 883
phone graph into ANGIE parse trees, this time 884
using a grammar that had been trained on the 885
Jupiter word lexicon, and allowing unknown 886
words to compete with known words in the search. 887
Table 2 reports some recognition and under- 888

Table 2

Recognition performance in terms of word error rate (WER) and concept error rate (CER) in the Jupiter weather domain, for utterances containing out-of-vocabulary city names

	WER (%)	CER (%)
Baseline	24.6	67.0
Two-stage ANGIE	15.6	31.3
Three-stage ANGIE/ TINA	17.4	21.8

The percentages indicate error rates for *all* the words. Unknown words were counted as correct if they were identified as such. The baseline system had no capability to handle unknown words.

889 standing error rates for three different systems,
 890 where words are considered as correct if they are
 891 unknown to the recognizer and correctly identified
 892 as unknown cities. Every utterance in this selected
 893 test set contained an unknown city, although the
 894 systems were unaware that this was the case. All
 895 systems used the same set of context-dependent
 896 acoustic models. The baseline system was a stan-
 897 dard SUMMIT recognizer configuration utilizing a
 898 word trigram language model, but with no capa-
 899 bility to deal with unknown words. The two-stage
 900 system utilized the foot-based ANGLE grammar
 901 reconfigured as an FST in the first stage, sup-
 902 ported by a foot-trigram language model. The
 903 output of this stage was a phone graph that was
 904 subsequently searched in a second stage, parsing
 905 with an ANGLE grammar that mapped to the rec-
 906 ognizer vocabulary, supporting novel unknown
 907 generations as unknown words. An optional third
 908 stage parsed the word graph proposed by the
 909 second stage using our TINA natural language
 910 system (Seneff, 1992), which sometimes favored a
 911 solution that was suboptimal in the second stage
 912 scores. The two-stage system yielded a 36.6%
 913 reduction in word error rate (WER), and a 53.3%
 914 reduction in concept error rate (CER), compared
 915 to the baseline system. The reason that CER is not
 916 100% for the baseline is that other concepts, such
 917 as dates, state names, and topic of inquiry (e.g.,
 918 “will it rain?”) are also counted. With the addition
 919 of NL support, the concept error rate improved
 920 further to 22%, a net reduction of 67.5%, although
 921 this was accompanied by an increase in recogni-
 922 tion error.

923 3.4. Automatic acquisition of new words

924 In addition to proposing unknown words, the
 925 system described above is also capable of propos-
 926 ing spellings for these words. Some examples of
 927 proposed spellings are given in Fig. 8. A recogni-
 928 tion evaluation of the proposed spellings in terms
 929 of letter substitutions, insertions, and deletions,
 930 was computed for the unknown words that were
 931 correctly tagged as such. The result was a 57.8%
 932 letter error rate, which, while quite high, is still
 933 substantially better than chance performance. The
 934 significant result is that we have formulated a

alameda	→	alumida
hanover	→	anover
hatteras	→	sateras
madagascar	→	madigasgar
mapleton	→	mapelton
mountainview	→	mounonvue
youngstown	→	janston

Fig. 8. Some examples of unknown cities and their proposed spellings, produced by the ANGLE two-stage recognizer. Spellings were extracted from the letter-phonemes at the preterminal layer of the ANGLE parse tree.

935 procedure for modeling unknown words by a
 936 technique of generalizing from the known words,
 937 and have been able to locate the unknown words
 938 in a user utterance and propose a set of alternate
 939 spellings and pronunciations for these words.

940 A further experiment aimed at improving the
 941 sound-to-letter performance, conducted by
 942 Gabovich (2002), utilized the PhoneBook (Dupont
 943 et al., 1990) isolated word corpus as the acoustic
 944 data. PhoneBook is a set of approximately 92,000
 945 isolated word utterances spoken over the tele-
 946 phone by a large number of native speakers of
 947 American English. It utilizes a vocabulary of
 948 about 8000 words, and the data have been as-
 949 signed to speaker-disjoint and vocabulary-disjoint
 950 training and test sets.

951 Our interest was in understanding how reliably
 952 we can spell the unseen words from PhoneBook if
 953 we choose to represent them only by the column
 954 bigram FST obtained by training an ANGLE
 955 grammar on the corpus. We compared perfor-
 956 mance on the training set with that on the test set,
 957 to measure how well the models are able to gen-
 958 eralize to unseen data. We are aware that sparse
 959 data problems will cause a certain percentage of
 960 the words to fail, even if given a perfect phonetic
 961 transcription, which would force these words to
 962 choose a suboptimal solution. But a more serious
 963 source of error is the difficult recognition task of
 964 producing a phonetic transcription from a wave-
 965 form without explicit knowledge of a lexicon.

966 We began with a lexicon of words in Phone-
 967 Book represented in terms of morph sequences
 968 whose pronunciations were in turn represented as

969 phonemic units.¹⁰ We developed a procedure to
 970 convert the phonemic baseforms into *spellnemic*
 971 baseforms by utilizing a letter-to-phoneme *ANGIE*
 972 grammar and then inferring the correct spellneme
 973 by associating the terminal letter sequence with the
 974 corresponding preterminal phoneme. We could
 975 then derive an *ANGIE* grammar mapping *pho-*
 976 *nemes* to spellnemes semi-automatically. Further-
 977 more, we took advantage of the research on
 978 phone-to-phoneme modeling to create a statistical
 979 representation of the phonological rules account-
 980 ing for the variations in pronunciation of the
 981 individual words. We trained the phone-to-*pho-*
 982 *neme* grammar on aligned phonetic transcriptions
 983 for the training corpus, and derived a corre-
 984 sponding column bigram FST. In parallel, the
 985 phoneme-to-spellneme grammar was trained on
 986 phonemic representations automatically obtained
 987 by parsing the letters of the training corpus using
 988 a letter-to-phoneme grammar, verified at the morph
 989 level against the lexicon. The language model was
 990 then just a composition of the two resulting col-
 991 umn bigram FSTs. Weights were optimized on an
 992 independent development set.

993 This system was tested on PhoneBook data, with
 994 the main goal of observing how well training would
 995 generalize to words that the system had never ob-
 996 served in training. The task is more difficult than
 997 phonetic recognition, in that a sound-to-letter
 998 system is embedded in the overall task. For exam-
 999 ple, a recognition of “fragmental” as “fraggmittle”
 1000 has only a single phonetic error (missing /n/), but
 1001 gets a 50% letter error rate.

1002 Results are summarized in Table 3. Overall
 1003 letter error rate (LER) increased from 34.1% to
 1004 41.0% when comparing the training set with the
 1005 test set, which we feel reflects fairly good general-
 1006 ization capabilities. We were also interested in
 1007 assessing how well the system would perform on
 1008 the test set if the phonetic transcription were per-
 1009 fect. Notice that this is different from and more
 1010 difficult than the phoneme-to-letter task discussed
 1011 in Section 3.3, since it is mapping from *phones* to
 1012 letters. This experiment will measure the capabili-

Table 3

Recognition performance (letter error rate) on the training and test sets for the task of automatically proposing a spelling of an unknown spoken word, for the PhoneBook telephone-quality isolated word corpus

	Sub (%)	Del (%)	Ins (%)	LER (%)
Training Corpus	16.6	7.5	10.0	34.1
Test Corpus	21.2	9.3	10.5	41.0

See text for details.

ties of the sound-to-letter system independent of 1013
 the phonetic recognition subtask. We obtained a 1014
 LER of 12.7% using as inputs the forced phonetic 1015
 alignments, for the subset (70.4%) of the test set 1016
 that had any solution at all through the FST space. 1017
 Nearly 30% failed to parse, clearly indicating that 1018
 we need to add a back-off mechanism to support 1019
 them. However, overall performance on the set 1020
 that parsed versus the set that failed drops by less 1021
 than 6%. 1022

3.4.1. Integrating pronunciation and spelling information 1023 1024

In the context of an interactive dialogue system, 1025
 there are further options available to help with the 1026
 task of unknown word acquisition. Having de- 1027
 tected that there may be an unknown city, the 1028
 system can solicit from the user a spoken spelling 1029
 form for the word. An *ANGIE* grammar with let- 1030
 ters as terminals can be used to parse a letter graph 1031
 produced by the *SUMMIT* recognizer. This inde- 1032
 pendent source can be matched against the pro- 1033
 posed solutions from the word pronunciation in 1034
 order to select something that is consistent with 1035
 both sources. A final resource that is available with 1036
 telephone input is a keypad entry of the unknown 1037
 word. This provides the interesting constraint that 1038
 each key disambiguates into one of three possible 1039
 letters. This can be formulated as a strict language 1040
 model and provide further constraint to the 1041
 problem. 1042

We have been pursuing the above ideas in a 1043
 joint research project with Chung.¹¹ For the 1044
 experiments described below, we have created an 1045

¹⁰ This lexicon had been prepared by Livescu as part of her research on duration modeling.

¹¹ Now at CNRI in Washington, DC.

1046 ANGIE lexicon of about 100,000 personal names,
1047 originally obtained from the Web, and have
1048 trained ANGIE parse trees on this corpus to pro-
1049 duce a model mapping phonemes to letters.

1050 In an initial experiment (Chung and Seneff,
1051 2002), we developed a recognition system that is
1052 able to integrate information from a keypad input
1053 of the spelling of the word with information culled
1054 from a spoken pronunciation of the word, as
1055 schematized in Fig. 9. We defined the search space
1056 by composing an FST mapping phonetics to
1057 graphemics with an FST specifying all possible
1058 pronunciations obtainable from the keypad in-
1059 puts. We also incorporated a morph bigram for
1060 increased linguistic support, where the possible
1061 organizations of letters into morph units are
1062 determined by the parsing grammar.

1063 To evaluate this idea, we conducted experi-
1064 ments on both the OGI name corpus (Cole et al.,
1065 1992) and a set of enrollment data obtained from
1066 our Mercury system (Seneff and Polifroni, 2000).
1067 In both cases, about 16% of the names were not
1068 present in our lexicon, an indication that the un-
1069 known word problem would be unavoidable in a
1070 personal name recognition task. The OGI set
1071 contains isolated first and last names, whereas the
1072 Mercury data are utterances containing both first
1073 and last name spoken sequentially.

1074 Results are summarized in Table 4. The system
1075 performs very well on letter accuracy for the in-

Table 4

Performance results for an experiment integrating telephone keypad inputs with spoken names, to produce hypothesized spellings for the names

Test set	IV subset (84%)		OOV subset (16%)	
	LER (%)	WER (%)	LER (%)	WER (%)
Mercury	1.7	8.1	12.0	43.2
OGI	1.8	8.1	13.3	57.3

Letter error rates (LER) and word error rates (WER) are reported for the in-vocabulary (IV) and out-of-vocabulary (OOV) portions of the Mercury and OGI test sets. Both sets have about a 16% OOV rate.

vocabulary portion of both sets, as might be ex- 1076
pected. However, it should be pointed out that this 1077
system has no explicit knowledge of the vocabu- 1078
lary that it was trained on, such as a word lexicon. 1079
It is encouraging that the system was able to ob- 1080
tain a perfect spelling for nearly half of the un- 1081
known words. If the search were restricted by a 1082
word lexicon, clearly none of the OOV words 1083
would have obtained a correct spelling. For fur- 1084
ther information on this topic, please see Chung 1085
and Seneff (2002). 1086

An extension of this work resulted in a system 1087
that can recognize a *spoken* spelling of a word 1088
jointly with the corresponding *pronounced* word, 1089
using an integrated solution that improves the 1090
recognition of the spelled letters by incorporating 1091
the constraints of a sound-to-letter model applied 1092
to the pronounced word. We have integrated this 1093
technology into a dialogue system that can learn 1094
new words by prompting a user to speak and spell 1095
the word in a single turn (Chung et al., in press). 1096
We have thus far only incorporated this capability 1097
into a user enrollment phase in the Orion task 1098
delegation system (Seneff et al., 2000), but we ex- 1099
pect it to be much more generally applicable. Ta- 1100
ble 5 gives its letter and word error rates on a 1101
corpus of telephone quality “speak-and-spell” 1102
data, divided into in-vocabulary and OOV subsets. 1103

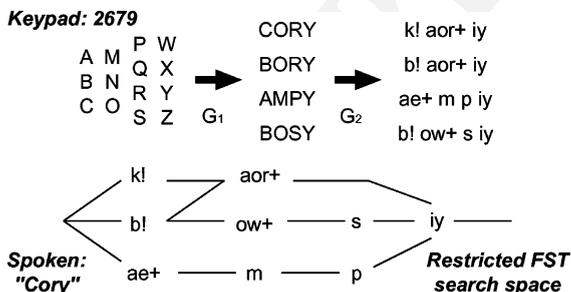


Fig. 9. A schematic for integrating keypad input with phonetic recognition, to produce a hypothesized spelling and pronunciation for the name “Cory.” Entered at the keypad is the sequence “2679,” producing a total of 144 possible four-letter sequences. Using a subword language model, FST G_1 sets out probable names, and FST G_2 maps the letters to phonemes probabilistically, based on a grammar encoding letter-to-sound rules.

3.5. Novel FST configurations

All of the research efforts described thus far 1105
involving an FST formulation of ANGIE’s parse 1106
structure have taken the point of view that the 1107
parse trees are decomposed into a simple column- 1108

Table 5

Letter error rates (LER) and word error rates (WER) for speak-and-spell utterances for peoples' first and/or last names

Set	No. of Utterances	LER (%)	WER (%)
In-Voc	416	8.4	27.4
OOV	219	12.4	46.1

Results are reported separately for words that were in the 100,000 word training lexicon (In-Voc) and words that were not part of the training lexicon (OOV). See (Chung et al., in press) for details.

1109 column transition matrix. This approach depends
 1110 upon direct observation in the training set of every
 1111 unique column pattern in every appropriate col-
 1112 umn context, and therefore can suffer from sparse
 1113 data problems. Recent research by Mou et al.
 1114 (2001) has been able to successfully encode the
 1115 *entire parsing mechanism* into an FST formulation,
 1116 thus retaining the generality that the rules achieve
 1117 directly in the FST representation. His strategy is
 1118 to use a recursive transition network formulation
 1119 to encode the context free rule component, which
 1120 can produce as outputs a detailed parse of the
 1121 input phone sequence. Subsequently, each layer
 1122 can be separately modelled so as to ignore all of
 1123 the elements that are irrelevant to that layer,
 1124 inserting just the portion of the probability model
 1125 that is provided for that layer. By jointly com-
 1126 posing the FSTs representing all of the layers, the
 1127 entire probability model can be inserted into the
 1128 resulting FST composition. While the resulting
 1129 FST is substantially larger than the FSTs obtained
 1130 from the column bigram approach, this formula-
 1131 tion is attractive because it provides a detailed
 1132 parse of each word, and because it permits us to
 1133 explore a variety of different probability formula-
 1134 tions to help identify which aspects of ANGIE's
 1135 probability model are most crucial.

1136 **4. Summary and future work**

1137 This paper describes a framework for acquiring
 1138 subword linguistic knowledge by parsing letters
 1139 and/or spoken pronunciations into words, via a
 1140 context-free grammar, and acquiring a supporting
 1141 probability model from a corpus of observations
 1142 within the domain of interest. We have identified

1143 several different ways in which such a framework
 1144 has utility in tasks related mainly to speech rec-
 1145 ognition. These include letter-to-sound and sound-
 1146 to-letter modeling, acquiring a probability model
 1147 for phonological realizations of words in fluent
 1148 speech, developing a duration model that takes
 1149 into account the hierarchy and also produces as a
 1150 side-effect a word-by-word speaking rate estimate,
 1151 acquiring a model for unknown words by gener-
 1152 alizing from the observed known words, and
 1153 obtaining a high quality phonetic graph in the first
 1154 stage of a two-stage large-vocabulary recognition
 1155 task. The ideas described here encompass research
 1156 that I have been conducting over the last several
 1157 years, collaboratively with both students and
 1158 researchers in the SLS group.

1159 We are encouraged by the results of the pho-
 1160 nological modeling, which demonstrated signifi-
 1161 cant reductions in understanding error rates in our
 1162 Mercury flight reservation domain. Probability
 1163 modeling for phonology might have even higher
 1164 payoff in recognition involving human-human
 1165 dialogues, such as the Switchboard corpus (God-
 1166 frey et al., 1992), where the speaking style is likely
 1167 to be considerably more casual than that used in
 1168 spoken interactions with a computer dialogue
 1169 system.

1170 We anticipate that the ANGIE hierarchical
 1171 representations can play a role in subword mod-
 1172 eling for speech synthesis. For instance, a correct
 1173 durational model is more critical in speech syn-
 1174 thesis, and it is known that phoneme durations
 1175 depend significantly on the position of the pho-
 1176 neme in the syllable and of the syllable in the word.
 1177 Furthermore, ANGIE's hierarchical framework
 1178 might provide a convenient mechanism to aid in
 1179 unit selection for concatenative speech synthesis.

1180 The main original motivation for characterizing
 1181 word substructure was to be able to model un-
 1182 known words as derivative from the substructure
 1183 of known words. The ability to support the auto-
 1184 matic acquisition of new words to both the rec-
 1185 ognition and understanding components of a
 1186 spoken conversational system will likely lead to a
 1187 breakthrough in dialogue system design. A sys-
 1188 tem's ability to immediately augment its working
 1189 vocabulary with a list of names obtained from a
 1190 Web page being presented to the user will greatly

1191 enhance the set of services it can offer to its user
 1192 population. Furthermore, if the user can simply
 1193 speak and spell a word they would like to see ad-
 1194 ded, they are empowered to configure the system
 1195 in ways that will be of much greater use to them in
 1196 the future. Ongoing research is aimed at develop-
 1197 ing conversational systems with flexible vocabu-
 1198 laries, where proper nouns presented to the user in
 1199 Web pages are automatically added to the system's
 1200 working vocabulary, and the user is empowered to
 1201 personalize the system to their own favored
 1202 information sources through natural spoken
 1203 interaction.

1204 Acknowledgements

1205 The research described in this paper would not
 1206 have been possible without the contributions made
 1207 by many former and current members of the
 1208 Spoken Language Systems Group. They include,
 1209 in alphabetical order, Grace Chung, Vladislav
 1210 Gabovich, Raymond Lau, Karen Livescu, Helen
 1211 Meng, Xiaolong Mou, Aarati Parmar, and Chao
 1212 Wang. I am especially indebted to Grace Chung,
 1213 whose doctoral thesis made significant advances
 1214 towards the goals that were originally conceived
 1215 for the ANGIE framework. Victor Zue read drafts
 1216 of the manuscript and made many helpful sug-
 1217 gestions.

1218 References

- 1219 Allen, J., Hunnicutt, M.S., Klatt, D., 1987. From Text to
 1220 Speech: The MITalk System. Cambridge Studies in Speech
 1221 and Science Communication. Cambridge University Press,
 1222 Cambridge.
 1223 Bacchiani, M., Ostendorf, M., 1998. Using automatically
 1224 derived acoustic subword units in large vocabulary speech
 1225 recognition. In: Proc. ICSLP'98, Vol. IV, Sydney, Australia,
 1226 November 1998, pp. 1319–1322.
 1227 Bazzi, I., Glass, J., xxxx. Learning units for domain-indepen-
 1228 dent out-of-vocabulary word modeling. In: Proc. Euro-
 1229 speech, Aalborg, Denmark, pp. 61–64.
 1230 Chomsky, N., Halle, M., 1968. The Sound Pattern of English.
 1231 Harper & Row, New York, NY (Republished in Paperback,
 1232 MIT Press, Cambridge, MA, 1991).

- Chung, G., 1997. Hierarchical Duration Modeling for a Speech 1233
 Recognition System. S.M. Thesis, MIT Department of 1234
 Electrical Engineering and Computer Science. 1235
 Chung, G., 2000a. Automatically incorporating unknown 1236
 words in Jupiter. In: Proc. ICSLP, Beijing, China, October 1237
 2000, pp. 520–523. 1238
 Chung, G., 2000b. A three-stage solution for flexible vocabulary 1239
 speech understanding. In: Proc. ICSLP 2000, Beijing, 1240
 China, October 2000, pp. 266–269. 1241
 Chung, G., 2001. Towards multi-domain speech understanding 1242
 with flexible and dynamic vocabulary. Ph.D. Dissertation, 1243
 MIT Department of Electrical Engineering and Computer 1244
 Science, June 2001. 1245
 Chung, G., Seneff, S., 1997. Hierarchical duration modeling for 1246
 speech recognition using the ANGIE framework. In: Proc. 1247
 EUROSPEECH'97, Rhodes, Greece, September 1997, pp. 1248
 1475–1478. 1249
 Chung, G., Seneff, S., 1998. Improvements in speech under- 1250
 standing accuracy through the integration of hierarchical 1251
 linguistic, prosodic, and phonological constraints in the 1252
 Jupiter domain. In: ICSLP'98, Sydney, Australia, December 1253
 1998, pp. 935–939. 1254
 Chung, G., Seneff, S., 2002. Integrating speech with keypad 1255
 input for automatic entry of spelling and pronunciation of 1256
 new words. In: Proc. ICSLP '02, Vol. III, Denver, CO, 1257
 September 2002, pp. 2061–2064. 1258
 Chung, G., Seneff, S., Wang, C., in press. Automatic acquisition 1259
 of names using speak and spell mode in spoken dialogue 1260
 systems. In: Proc. HLT-NAACL, Edmonton, Canada, May 1261
 2003. 1262
 Church, K.W., 1983. Phrase-structure parsing: a method for 1263
 taking advantage of allophonic constraints. Ph.D. Thesis, 1264
 Department of Electrical Engineering and Computer Sci- 1265
 ence, MIT, Cambridge, MA. 1266
 Cohen, M.H., 1989. Phonological structures for speech recog- 1267
 nition. Ph.D. Dissertation, University of California, Berke- 1268
 ley, CA. 1269
 Cole, R. et al., 1992. A telephone speech database of spelled and 1270
 spoken names. In: Proc. ICSLP'92, Banff, Canada, October 1271
 1992. 1272
 Cook, G., Robinson, T., Christie, J., 1998. Real-time recogni- 1273
 tion of broadcast news. In: Proc. ICSLP'98, Vol. IV, 1274
 Sydney, Australia, November 1998, pp. 1319–1322. 1275
 Dupont, S., Boulard, H., Deroo, O., Fontaine, V., Boite, J.- 1276
 M., 1990. Hybrid HMM/ANN systems for training inde- 1277
 pendent tasks: experiments on phonebook and related 1278
 improvements. In: Proc. ICASSP, Albuquerque, NM. 1279
 Gabovich, V.Y., 2002. A multi-stage sound-to-letter recognizer. 1280
 M.Eng Thesis, MIT, May 2002. 1281
 Gauvain, J.L., Lamel, L.F., Adda, G., Adda-Decker, M., 1993. 1282
 Speaker independent continuous speech dictation. In: Proc. 1283
 EUROSPEECH'93, Berlin, Germany, September 1993, pp. 1284
 125–128. 1285
 Glass, J., Chang, J., McCandless, M., 1996. A probabilistic 1286
 framework for feature-based speech recognition. In: Proc. 1287
 ICSLP'96, Philadelphia, PA, October 1996, pp. 2277–2280. 1288

- 1289 Glass, J.R., Hazen, T.J., Telephone-based conversational
1290 speech recognition in the Jupiter domain. In: ICSLP'98,
1291 Sydney, Australia, December 1998, pp. 1327–1330.
- 1292 Godfrey, J., Holliman, E., McDaniel, J., 1992. SWITCHBOARD:
1293 telephone speech corpus for research and development. In:
1294 Proc. ICASSP-92, San Francisco, pp. 517–520.
- 1295 Greenberg, S., 1999. Speaking in shorthand—a syllable-centric
1296 perspective for understanding pronunciation variation.
1297 *Speech Commun.* 29, 159–176.
- 1298 Hayes, B., 1995. *Metrical Stress Theory: Principles and Case*
1299 *Studies*. University of Chicago Press, Chicago.
- 1300 Hetherington, I.L., 1994. The problem of new, out-of-vocabu-
1301 lary words in spoken language systems. Ph.D. Thesis,
1302 Department of Electrical Engineering and Computer Sci-
1303 ence, MIT, Cambridge, MA, October 1994.
- 1304 Hetherington, I.L., 2001. An efficient implementation of
1305 phonological rules using finite state transducers. In: Proc.
1306 Eurospeech 2001, Aalborg, Denmark, September 2001, pp.
1307 1599–1602.
- 1308 Hetherington, I.L., Zue, V., 1993. New words: implications for
1309 continuous speech recognition. In: Proc. Third European
1310 Conf. on Speech Communication and Technology, Berlin,
1311 Germany, September 1993.
- 1312 Kahn, D., 1976. Syllable-based generalizations in English
1313 phonology. Ph.D. Thesis, Department of Linguistics and
1314 Philosophy, MIT, Cambridge, MA.
- 1315 Kucera, H., Francis, W., 1967. *Computational Analysis of*
1316 *Present-Day American English*. Brown University Press.
- 1317 Lau, R., 1998. Subword lexical modeling for speech recogni-
1318 tion. Ph.D. Thesis, Department of Electrical Engineering
1319 and Computer Science, Massachusetts Institute of Technol-
1320 ogy, Cambridge, MA, May 1998.
- 1321 Lau, R., Seneff, S., 1997. Providing sublexical constraints for
1322 word spotting within the ANGIE framework. In: Proc.
1323 EUROSPEECH'97, Rhodes, Greece, September 22–25,
1324 1997, pp. 263–266.
- 1325 Lau, R., Seneff, S., 1998. A unified system for sublexical and
1326 linguistic modeling using ANGIE and TINA. In: ICSLP'98,
1327 Sydney, Australia, December 1998, pp. 2443–2446.
- 1328 Meng, H., 1995. Phonological parsing for bi-directional letter-
1329 to-sound/sound-to-letter generation. Ph.D. Thesis, Depart-
1330 ment of Electrical Engineering and Computer Science, MIT,
1331 Cambridge, MA, June 1995.
- 1332 Meng, H., Hunnicutt, S., Seneff, S., Zue, V., 1996. Reversible
1333 letter-to-sound/sound-to-letter generation based on parsing
1334 word morphology. *Speech Commun.* 18, 47–63.
- 1335 Mou, X., Seneff, S., Zue, V., 2001. Context-dependent proba-
1336 bilistic hierarchical sub-lexical modeling using finite state
1337 transducers. In: Proc. EUROSPEECH 2001, Aalborg,
1338 Denmark, September 2001, pp. 451–454.
- 1339 Nguyen, L., Anastasakos, T., Kubala, F., LaPre, C., Makhoul,
1340 J., Schwartz, R., Yuan, N., Zavaliagkos, G., Zhao, Y., 1995.
1341 The 1994 BBN/BYBLOS speech recognition system. In:
1342 Proc. ARPA Spoken Language Systems Technology Work-
1343 shop'95, Austin, TX, January 1995, pp. 693–696.
- Onishi, S., Yamamoto, H., Sagisaka, Y., 2001. Structured
1344 language model for class identification of out-of-vocabulary
1345 words arising from multiple word classes. In: EURO-
1346 SPEECH 2001, Aalborg, Denmark, September 2001, pp.
1347 693–696. 1348
- Parmar, A.D., 1997. A semi-automatic system for the syllab-
1349 ification and stress assignment of large lexicons. M.Eng
1350 Thesis, Department of Electrical Engineering and Computer
1351 Science, MIT, Cambridge, MA, June 1997. 1352
- Randolph, M.A., 1989. Syllable-based constraints on properties
1353 of English sounds. Ph.D. Thesis, Department of Electrical
1354 Engineering and Computer Science, MIT, Cambridge, MA,
1355 September 1989. 1356
- Scalise, S., 1986. *Generative Morphology*. Foris Publications,
1357 Dordrecht, Netherlands. 1358
- Seneff, S., 1992. TINA: a natural language system for spoken
1359 language applications. *Comput. Linguist.* 18 (1), 61–86. 1360
- Seneff, S., 1998. The use of linguistic hierarchies in speech
1361 understanding. Keynote Address. In: ICSLP'98, Sydney,
1362 Australia, December 1998, pp. 3321–3330. 1363
- Seneff, S., Polifroni, J., 2000. Dialogue management in the
1364 mercury flight reservation system. In: Proc. ANLP-NAACL
1365 2000, Satellite Workshop, Seattle, Washington, May 2000,
1366 pp. 1–6. 1367
- Seneff, S., Wang, C., 2002. Modelling phonological rules
1368 through corpus-trained linguistic hierarchies. In: ISCA
1369 Tutorial and Research Workshop on Pronunciation Mod-
1370 eling and Lexicon Adaptation for Spoken Language, Estes
1371 Park, CO, September 2002, pp. 71–76. 1372
- Seneff, S., Lau, R., Meng, H., 1996. ANGIE: A new framework
1373 for speech analysis based on morpho-phonological model-
1374 ing. In: Proc. ICSLP'96, Vol. 1, Philadelphia, PA, October
1375 1996, pp. 110–113. Available from <http://www.sls.lcs.mit.edu/raylau/icslp96_angie.pdf>. 1376 1377
- Seneff, S., Chuu, C., Cyphers, D.S., 2000. ORION: from on-line
1378 interaction to off-line delegation. In: Proc. ICSLP '00, Vol.
1379 II, Beijing, China, October 2000, pp. 142–145. 1380
- Selkirk, E., 1982. Syllables. In: van der Hulst, Harry, Smith,
1381 Norval (Eds.), *The Structure of Phonological Representa-*
1382 *tions*. Foris, Dordrecht, pp. 337–383. 1383
- Weintraub, M., Murveit, H., Cohen, M., Price, P., Bernstein, J.,
1384 Baldwin, G., Bell, D., 1989. Linguistic constraints in hidden
1385 Markov model based speech recognition. In: Proc. ICASSP
1386 '89, Glasgow, Scotland, May 1989, pp. 699–702. 1387
- Woodland, P.C., Young, S.J., 1993. The HTK tied-state
1388 continuous speech recognizer. In: Proc. European Conf.
1389 on Speech Commun. and Technology, Vol. 3, pp. 2207–
1390 2210. 1391
- Zue, V., 1983. The use of Phonetic rules in automatic speech
1392 recognition. *Speech Commun.* 2, 181–186. 1393
- Zue, V., Glass, J., Goodine, D., Hirschman, L., Leung, H.,
1394 Phillips, M., Polifroni, J., Seneff, S., 1991. The MIT ATIS
1395 system: preliminary development, spontaneous speech data
1396 collection and performance evaluation. In: Proc. European
1397 Conf. on Speech Communication and Technology, Vol. 2,
1398 Genoa, Italy, September 1991, pp. 537–541. 1399