

1 International Journal of Semantic Computing
 Vol. 1, No. 3 (2007) 1–12
 3 © World Scientific Publishing Company



5 **CROSS-GENRE FEATURE COMPARISONS FOR SPOKEN**
 6 **SENTENCE SEGMENTATION**

7 SEBASTIEN CUENDET^{*,‡}, DILEK HAKKANI-TUR^{*,§}, JAMES FUNG^{*,¶},
 BENOIT FAVRE^{*,||} and ELIZABETH SHRIBERG^{*,†}

**International Computer Science Institute
 Berkeley, CA 94704, USA
 www.icsi.berkeley.edu*

*†SRI International, Menlo Park, CA 94025, USA
 www.speech.sri.com*

‡ees@icsi.berkeley.edu

‡cuendet@icsi.berkeley.edu

§dilek@icsi.berkeley.edu

¶jgf@icsi.berkeley.edu

||favre@icsi.berkeley.edu

19 Automatic sentence segmentation of spoken language is an important precursor to down-
 stream natural language processing. Previous studies combine lexical and prosodic fea-
 21 tures, but can impose significant computational challenges because of the large size of
 feature sets. Little is understood about which features most benefit performance, partic-
 23 ularly for speech data from different speaking styles. We compare sentence segmentation
 for speech from broadcast news versus natural multi-party meetings, using identical
 25 lexical and prosodic feature sets across genres. Results based on boosting and forward
 selection for this task show that (1) features sets can be reduced with little or no loss in
 27 performance, and (2) the contribution of different feature types differs significantly by
 genre. We conclude that more efficient approaches to sentence segmentation and similar
 tasks can be achieved, especially if genre differences are taken into account.

29 *Keywords:* Sentence segmentation; prosody; feature selection.

31 **1. Introduction**

Recent speech processing tasks have focused on a range of genres that differ in
 33 speaking style — including news broadcasts, telephone conversations, lectures, and
 meetings. Such genres differ in many aspects, including vocabulary, syntax, turn-
 35 taking, discourse phenomena, disfluencies, paralinguistic effects, and prosody [1]. A
 typical approach to language processing tasks is to apply features and approaches
 37 developed for one genre, to another genre, using genre-specific training data where
 available. Other approaches use explicit adaptation techniques [2].

However, when matched training data is not available or in contexts in which
 speed and computational expense are important, it can be worthwhile to investigate

2 *S. Cuendet et al.*

1 which features contribute most to a task and whether or not feature utility depends
on the speaking style.

3 In this study we investigate the role of identically-defined lexical and prosodic
4 features, when applied to the same task across two very different speaking styles
5 — broadcast news and face-to-face multi-party meetings. We focus on the task
6 of automatic sentence segmentation, or finding boundaries of sentence units in the
7 otherwise unannotated (devoid of punctuation, capitalization, or formatting) stream
8 of words output by a speech recognizer. Sentence segmentation is of particular
9 importance for speech understanding applications, because techniques aimed at
10 semantic processing of speech input — such as parsing, machine translation, and
11 information extraction — are often developed for text-based applications and thus
12 assume the presence of overt sentence boundaries in their input [6, 5]. Sentence
13 boundary annotation is also important for aiding human readability of the output
of automatic speech recognition systems [3].

15 Previous approaches to sentence boundary detection in speech have combined
16 lexical with prosodic features (such as pause, pitch, and energy features), using
17 various machine learning techniques. One practical concern with such approaches
18 is that although the gain from inclusion of prosodic features is considerable, these
19 prosodic features require additional initial human effort and computational expense.
20 Thus, research is needed to determine which features are most useful, and how
21 feature utility differs for different styles of speech.

22 The goal of this study is to explore these questions for the task of sentence
23 segmentation in news versus meeting speech. Specifically, we ask:

- 24 (1) Can we achieve similar sentence boundary classification performance to an all-
25 features performance result using only a small set of prosodic features?
- 26 (2) Do the different speaking styles differ in terms of which prosodic features are
27 most useful for this task?

28 Results have implications not only for the task of sentence boundary detection, but
29 more generally for prosodic modeling for natural language understanding across
genres.

31 The following section describes the data set, features, and approach. Section 3
32 reports on experiments with lexical and prosodic features using boosting and forward
33 selection of features, and provides further analysis of lexical and prosodic
34 feature usage differences in the two different corpora. A summary and conclusions
35 are provided in Section 4.

36 2. Method

37 2.1. Data

38 To study the differences between the meetings and BN speech for the task of sen-
39 tence segmentation, we use the ICSI Meetings [9] and the TDT4 English Broadcast
News [11] corpora. The ICSI Meeting Corpus is a collection of 75 meetings, including

1 simultaneous multi-channel audio recordings, word-level orthographic transcrip-
 2 tions. The meetings range in length from 17 to 103 minutes, but generally run
 3 just under an hour each, totalling 72 hours. We use a 73 meeting subset of this
 4 corpus that was also used in the previous research [9] with the same split into
 5 training, held-out and test sets. The TDT4 Corpus was collected by the Linguis-
 6 tic Data Consortium (LDC) and includes multilingual raw material, newswire and
 7 other electronic text, web audio, broadcast radio and television. We use a subset of
 8 TDT4 English broadcast radio and television data in this study.

9 In the experiments to follow, classification models are trained on a set of data,
 10 tuned on a held-out set, and tested on an unseen test set, within each genre. Statis-
 11 tics on these data sets are shown in Table 1. The statistics in the tables are computed
 12 using the forced alignments between audio and reference transcriptions.

13 As shown in Table 1, the two different speaking styles differ significantly in
 14 mean sentence length, with sentences in meetings being only about half the length
 15 on average as those in broadcast news. Meetings (and conversational speech in
 16 general) tend to contain syntactically simpler sentences and significant pronominal-
 17 ization. News speech is typically read from a transcript, and more closely resembles
 18 written text. It contains for example appositions, center embeddings, and proper
 19 noun compounds, among other characteristics that contribute to longer sentences.
 20 Such differences are the result of both the more formal context of news speech and
 21 the speakers being professional speakers, as opposed to meetings where speakers are
 22 “common” people. Discourse phenomena also obviously differ across corpora, with
 23 meetings containing more turn exchanges, incomplete sentences, and higher rates of
 24 short backchannels (such as “yeah” and “uhhuh”) than speech in news broadcasts.

25 2.2. Features

26 Sentence segmentation can be seen as a binary classification problem, in which each
 27 word boundary must be labeled as a sentence boundary or as a non-sentence bound-
 28 ary.^a We define a large set of lexical and prosodic features, computed automatically
 29 based on the output of a speech recognizer, as described further, below.

Table 1. Data set statistics. Values are given in number of words, based on forced alignments.

	MRDA	TDT4
Training set size	90,000	150,000
Test set size	88,537	50,116
Held-out set size	110,851	23,363
Vocabulary size	10,887	18,697
Mean sentence length	6.54	14.69

^aMore detailed models may distinguish questions from statements, or complete from incomplete sentences.

4 *S. Cuendet et al.*

1 **Automatic speech recognition.** Automatic speech recognition results for the
2 ICSI Meetings data and TDT4 data were obtained using the state-of-the-art SRI
3 CTS system [15] and SRI BN system [12], respectively. The meetings recognizer was
4 trained using no acoustic data or transcripts from the analyzed meetings corpus.
5 The word error rate for the recognizer output of the complete meetings corpus
6 is 38.2%.

7 Recognition scores for the TDT4 corpus is not easily definable as only closed
8 captions are available that frequently do not match well the actual words of the
9 broadcast news shows. The estimated word error rate lies between 17% and 19%.

10 **Lexical features.** Previous work on sentence segmentation in broadcast news
11 speech and in telephone conversations has used lexical and prosodic information
12 [10, 4]. Additional work has studied the contribution of syntactic information [7].
13 Lexical features are usually represented as N -grams of words. In this work, lexical
14 information is represented by 6 N -gram features for each word boundary: 3 uni-
15 grams, 2 bigrams and 1 trigram. Naming the word preceding the word boundary of
16 interest as the *current* word, and the preceding and following words as the *previous*
17 and *next* word respectively, the 6 lexical features are as follows:

- 18 • unigrams: {previous}, {current}, {next},
- 19 • bigrams: {current, next}, {previous, current},
- 20 • trigram: {previous, current, next}.

21 **Prosodic features.** Prosodic information is represented using mainly continuous
22 values. We use 59 prosodic features, defined for and extracted from the regions
23 around each inter-word boundary. The features include the pause duration at the
24 boundary, normalized phone durations of the word preceding the boundary, and a
25 variety of speaker-normalized pitch features and energy features preceding, follow-
26 ing, and across the boundary. Features are an extension of similar features described
27 in [10]. The extraction region around the boundary focuses on either one-word win-
28 dows or brief time windows around the boundary. Measures include the maximum,
29 minimum or average value in this time range. Pitch features are normalized by
30 speaker, using the method to estimate a speaker's baseline pitch values described
31 in [10]. Duration features, which measure the duration of the last vowel and the last
32 rhyme in the word before the word boundary of interest, are normalized by statistics
33 on the relevant phones in the training data. We also include "turn" features based
34 on speaker changes. The turn features are computed differently on the two corpora.
35 In TDT4, the speaker turns are determined by an alignment between the output of
36 an external diarization system [13] and the words. Meetings are already broken up
37 by channel with one channel per speaker and thus do not need diarization. A turn
38 is added within a channel when the pause between two words is greater than 0.5
39 seconds. The reason for this is to match with the diarization system, where even if
40 there is no speaker change, a non-speech region greater than 0.5 second segments
41 the speaker turn.

1 **Boosting and forward selection.** For classification of word boundaries, we use
2 the AdaBoost algorithm [8], which has been shown to be one of the best classifiers
3 for this task [16]. Boosting aims to combine weak base classifiers to come up with
4 a strong classifier. The learning algorithm is iterative. In each iteration, a different
5 distribution or weighting over the training examples is used to give more emphasis
6 to examples that are often misclassified by the preceding weak classifiers. For this
7 approach we use the BoosTexter tool described in [8]. BoosTexter handles both
8 discrete and continuous features, which allows for a convenient incorporation of the
9 prosodic features described above (no binning is needed). The weak learners are
10 one-level decision trees (stumps).

11 To analyze the difference in prosodic feature importance to sentence segmenta-
12 tion in the two genres, we rank features according to the forward selection algorithm
13 (FSA). The FSA is an iterative algorithm that begins with an empty set of features.
14 At each iteration, every feature that has not yet been selected is evaluated together
15 with the previously-selected features. The feature that yields the best performance
16 is then added to the set of selected features and a new iteration, which consid-
17 ers the remaining features, begins. Although computationally expensive the FSA
18 has the advantage of being intuitive and of capturing the correlation between two
19 similar features. Indeed, once a feature has been selected, features with which it is
20 highly correlated are less likely to be picked, since they would bring few additional
21 knowledge to the classifier.

3. Experiments and Results

23 3.1. Overall results

24 Metrics Sentence segmentation quality is usually computed using one of two mea-
25 sures — F-measure or NIST error. F-measure is the harmonic mean of the recall
26 and precision measures of the sentence boundaries hypothesized by the classifier to
27 those assigned by human labelers. The NIST error rate is the ratio of the number
28 of incorrect hypotheses made by the classifier to the number of reference sentence
29 boundaries. If no boundaries are marked by sentence segmentation, this metric is
30 100%, but it can exceed 100%; the maximum error rate metric is the ratio of number
31 of words to the number of correct boundaries. In this work, we report performance
32 using only F-Measure.

33 3.2. Lexical N-grams

34 To characterize lexical differences across the two genres, we follow the comparative
35 study reported in [14] in the context of text categorization, and utilize the widely
36 used information gain (IG) metric. Given a term, *information gain* measures the
37 amount of information obtained for the class prediction from the presence/absence
38 of the term. In the case of a binary classification, the definition of the information

6 *S. Cuendet et al.*

gain of a term t is a simplification of the definition presented in [14]:

$$\begin{aligned}
 G(t) = & -p(N) \log p(N) - p(S) \log p(S) \\
 & + p(t) [p(N|t) \log p(N|t) + p(S|t) \log p(S|t)] \\
 & + p(\bar{t}) [p(N|\bar{t}) \log p(N|\bar{t}) + p(S|\bar{t}) \log p(S|\bar{t})]
 \end{aligned}$$

1 where S and N are the classes that designate a sentence boundary and a non-
 2 sentence boundary, respectively. Note that the IG score takes into account both
 3 classes, and we therefore do not need to take the average of the two classes. The χ^2
 4 statistic described in [14] is also useful to isolate the information of a term together
 5 with a particular class. However, in a two-class problem such as that examined here,
 6 the computation is symmetric, and therefore results are similar to those obtained
 7 using the IG score.

8 We consider each feature separately and compute the IG for each term that
 9 occurs in the feature vector (a term being a word in the case of the unigram feature,
 10 a bigram for the features represented by bigrams, etc.). For each genre and each of 6
 11 lexical features, we extract the 10 terms that have the highest score. By doing that,
 12 we isolate the words that have a strong correlation with the occurrence of sentence
 13 boundaries. The underlying assumption is that if two genres are similar, the terms
 14 that are the best indicators of the beginning or the end of sentences should be
 15 similar in both genres.

16 Tables 2, 3 show the top 8 terms according to their IG score for 2 of the 3
 17 unigram features and the bigram feature. IG values for the lexical features were
 18 computed on the held-out sets.

19 The tables show clear differences in word associations with sentence boundaries
 20 across genres. In meeting speech, as noted earlier, there are high rates of single-word
 21 backchannels such as **yeah**, **uhhuh**, and **right**. Since backchannels are treated as
 22 individual sentences in the annotation of this corpus, the presence of a backchannel
 23 word is a strong indicator for both a preceding and a following sentence boundary.
 (Note that not all cases of, for example, **right** are backchannels, since the word can

Table 2. Most frequent words for pre- and post-boundary unigram features.

Pre-boundary unigram		Post-boundary unigram	
MRDA	TDT4	MRDA	TDT4
yeah	the	yeah	i
uhhuh	to	so	of
okay	and	uhhuh	uh
right	of	and	to
the	a	but	but
huh	in	okay	he
i	washington	right	we
um	for	oh	i'm

Table 3. Most frequent words for pre-boundary — post-boundary bigram feature.

Pre-boundary — Post-boundary bigram	
MRDA	TDT4
yeah yeah	of the
uhhuh uhhuh	in the
yeah so	com the
yeah uhhuh	to the
yeah i	for the
okay so	court the
uhhuh yeah	glascoff coming
right so	at the

1 be used in other contexts. But in this data most backchannels use words that are
 2 more frequent in backchannels than in other contexts).

3 A quite different pattern is observed for the TDT4 corpus. In this case, words
 4 that have the highest IG score show no obvious correlation with the sentence bound-
 5 ary class. The explanation is that in BN, given the size of the data sets used, very
 6 few specific words appear repeatedly at sentence boundaries. Backchannels, fillers,
 7 and discourse markers are relatively rare, and a much larger set of words (includ-
 8 ing proper nouns) appear at sentence edges. As a consequence, words that obtain
 9 the best IG score for the TDT4 corpus are those that are highly correlated with
 10 the *non*-sentence boundary class distribution, i.e. words that are unlikely to end a
 11 sentence, such as **the**, **to**, or **a**. Note that this analysis does not hold for the *next*
 12 *word* feature, since words that begin a sentence have a pattern, even in the case of
 13 BN, as shown by the presence of **i**, **and**, and **but**.

14 Comparing the two lists (Tables 2 and 3) for the *current word* and the *next word*
 15 feature in the case of MRDA reveals the double usage of certain words like **yeah**
 16 or **okay**. In conversational speech, such words can be used either as backchannels,
 17 which make the rank high in the *previous word* table. On the other hand, they are
 18 also used to start new sentences, which explains why they are so well ranked in the
 19 *next word* table.

20 The symmetry between the two classes in the IG computation allows some
 21 bigrams highly correlated with non-sentence boundary to have a high score for
 22 TDT4. For example, for the *current word* — *next word* feature, the bigram **of the**
 23 has the highest score, since it appears 536 times, but only twice with a sentence
 boundary, and thus 534 times with a non-sentence boundary.

25 3.3. Prosodic features

26 To rank the prosodic features according to their importance, we ran the FSA for 20
 27 iterations. We used the BoosTexter tool [8] to train a classifier on the training data
 and evaluated the performance of the sentence segmentation on the held-out set. The

Table 4. Features selected for MRDA and TDT4; columns 2 and 3 show the F-Measure on the held-out set and the relative improvement from one feature to the next one, respectively. Column 4 shows the F-Measure on the test set. The last column is the F-Measure when using the feature alone.

Feature name	Held-out	Rel. impr.	Test	Alone
MRDA				
PAU-DUR	60.0	—	62.2	60.0
FOK-WRD-DIFF-LOLO-N	61.0	1.7%	62.8	22.8
LAST-RHYME-NORM-DUR-PH	61.8	1.3%	63.6	12.6
PAU-DUR-PREV	62.5	1.2%	64.3	11.1
CROSS-SPKR PAUSE	63.0	0.7%	64.6	47.8
ENERGY-WIN-DIFF-HIHI-N	63.2	0.4%	65.3	20.4
LAST-RHYME-DUR-PH	63.8	0.5%	65.3	11.8
LAST-VOW-DUR-Z	63.8	0.4%	65.6	6.7
TDT4				
PAU-DUR	56.0	—	55.4	56.0
FOK-DIFF-LAST-KBASELN	58.2	3.8%	57.0	34.4
FOK-WIN-DIFF-LOHI-N	59.4	2.1%	57.9	18.9
TURN-F	60.0	1.0%	58.5	35.5
PAU-DUR-PREV	60.2	0.4%	58.7	0.0
FOK-LR-MEAN-KBASELN	60.4	0.3%	58.7	0.0
FOK-DIFF-MNMN-N	60.5	0.2%	59.0	15.9
SLOPE-LAST-N	60.5	0.1%	58.9	0.8

1 feature with the best F-Measure was selected at each iteration. A classifier was then
 2 built on the training set and evaluated on the test set for each feature set. Table 4
 3 reports the features that were selected until the F-Measure stopped increasing, and
 4 the corresponding performance on the development and the test sets.

5 The two sets each make significant use of the pause duration feature. The “pau-
 6 dur-prev” or duration of the pause one boundary earlier than the boundary of
 7 consideration is useful in MRDA in part because of the prevalence of single-word
 8 sentences such as backchannels, as described in Sec. 2.1. Both corpora make ample
 9 use of pitch features. TDT4 makes more use of “baseline” normalized pitch features
 10 that compare the location of a particular preboundary word in a speaker’s pitch
 11 range to the value of a speaker’s estimated baseline pitch. The closer the local pitch
 12 value is to the speaker’s baseline, the more likely it is that the speaker is near a
 13 sentence end. This makes sense in that news (and read) speech is more careful and
 14 regular in intonation, whereas meeting speech is more informal and involves par-
 15 alinguistic variation that can shift ending pitch values. MRDA makes use of pitch in
 16 the second feature selected, but this feature compares the pitch in words before and
 17 after the boundary, rather than the current pitch value to the speaker’s estimated
 18 pitch floor. In this case, a large value of the feature indicates a sentence boundary,
 19 consistent with a large pitch reset. One interesting finding, perhaps counterintuitive
 20 at first, is that meetings make more use of duration features (of vowels or syllable
 21 rhymes) than do news broadcasts. Typically, there is a correlation between ends
 of major phrases and pre-boundary lengthening. Separate analysis revealed that

1 durational lengthening is indeed present for sentence boundaries in both corpora,
 2 but that in the case of news speech, lengthening occurs frequently elsewhere as well.
 3 That is, the register used in news broadcasts tends to insert frequent prominences
 4 and sub-sentential breaks, perhaps to keep the attention of the listener. Thus dura-
 5 tion features may cause considerable false alarms in the case of broadcast news and
 6 are therefore less useful than they are for conversational speech. Finally, energy
 7 features do not appear to be as useful as pause, pitch, and duration features, across
 8 genres.

9 In MRDA, the previous pause feature, measuring the pause duration before the
 10 current word on the same channel, brings a relative improvement three times as
 11 large as in TDT4. The previous pause feature captures information about short
 12 utterances. When it is high and the current pause is high too, it suggests that the
 13 current sentence is only one word long. This is especially appropriate for conversa-
 14 tional speech in which many utterances are backchannels. Often one-word long, as
 15 already mentioned earlier in the study of lexical features. For TDT4, the improve-
 16 ment over the pause for the second feature selected is larger than for MRDA (3.8%
 17 vs. 1.7%). The smaller pauses at non-boundaries in more formal speech, as well as
 18 longer pauses between the end of a sentence and the beginning of the next one,
 19 both explain this.

20 The cross-speaker pause feature is used only in MRDA by construction, since it
 21 measures not only the pause on a single channel (as the normal pause feature), but
 22 takes into account all of the channels. In the case of MRDA, where every speaker
 23 has a microphone, the cross-speaker pause is not equivalent to the pause feature,
 24 whereas it is equivalent to that feature in TDT4. In addition to the pause and
 25 the two first pitch features, the turn feature provides significant benefit for TDT4.
 26 The turn feature is a binary feature which indicates a change of speaker. In BN, the
 27 speaker turn is automatically estimated by a diarization system, whereas on MRDA
 28 a turn is introduced every time there is a pause longer than 0.5 seconds. Thus in
 29 MRDA, the turn is highly correlated with the pause feature, whereas in TDT4 it is
 30 an independent input.

31 Further differences between MRDA and TDT4 are shown in Table 5. In MRDA,
 32 using only the lexical features results in a significantly better performance than
 33 using only the prosodic features (+4.3% absolute). On the contrary, in TDT4, the
 34 prosodic model performs better than the lexical model (+1.4% absolute). The higher
 35 performance of the prosodic model reflects the more formal speech of TDT4, both

Table 5. Comparison of the F-Measure with lexical features only, prosodic features only, and prosodic and lexical features together for MRDA and TDT4.

Corpus	Lexical	Prosodic	Both
MRDA	69.8	65.6	73.7
TDT4	58.0	59.4	61.8

Table 6. F-Measure with all 59 prosodic features and after 20 iterations of the FSA algorithm.

Corpus	Chance	All features (59)	FSA (20 iterations)
MRDA	15.8	65.6	66.0
TDT4	6.9	59.4	59.2

1 because speakers make a better use of prosody and because the lexical model is
 2 less strongly correlated with sentence boundaries than in conversational speech, as
 3 explained earlier.

4 Table 6 shows the performance of sentence segmentation for both corpora, when
 5 the classifier makes use of all the features and when it uses only the first 20 fea-
 6 tures selected by the FSA algorithm. While the performance with all the features
 7 is expected to be better than that with only a subset, one can observe that the
 8 performances are very close. On MRDA, the performance with the reduced set of
 9 20 features is actually better than when using all features. Going back to Table 4,
 10 one can see that the performance with all the prosodic features is already reached
 11 by the reduced set of prosodic features after eight iterations of the FSA. In the case
 12 of TDT4, the performance after eight iterations is 0.5% absolute less than that with
 13 all of the features, but after four iterations only, the F-Measure score is less than
 14 1% less than that with all the features. The score of the full set is reached at iter-
 15 ation 22 of the FSA on TDT4. Thus on TDT4 and MRDA, the same performance
 16 is reached by using 37% and 14% of the features, respectively.

17 Table 6 also shows the “chance performance” on one corpus. The chance perfor-
 18 mance assumes no knowledge about the data and simply classifies every example
 19 of the test set with respect to prior probability of each class in the training set.
 20 The performance reported is an average of the F-Measure over 10 runs. Comparing
 21 the chance performance with the score of the features when used alone (Table 4)
 22 shows that the performance with the first feature selected for both corpora (pause)
 23 is already four times as high as the chance performance. Some features picked later
 24 by the FSA have a performance worse than chance, but together with the previous
 25 features chosen they are able to improve the performance.

26 Reducing the set of features is important in terms of memory and CPU usage,
 27 as well as for computation time. For example, on the same machine, the training
 28 time is reduced by a factor of seven when using only eight features versus using all
 29 59 features (2 h vs. 14 h).

4. Summary and Conclusions

30 We have compared lexical and prosodic sentence segmentation features for broadcast
 31 news and meeting speech, using identical feature sets and definitions for both genres.
 32 Analysis of sentence distributions in the two corpora show significant differences
 33 in average sentence length, lexical, and prosodic features. For example, sentences
 34 in meetings are on average only half as long as those in broadcast conversations.
 35

1 Whereas important lexical N-grams for meetings are positive cues associated with
2 backchannels and various discourse phenomena, lexical N-grams for news speech
3 are negative cues, i.e. N-grams in which a sentence boundary is highly unlikely.

4 Experiments on prosodic features using forward selection show that similar or
5 even better performance can be achieved by using fewer features. Useful feature
6 types, however, depend on the corpus. While both genres make use of pause and
7 pitch information, pitch features contribute relatively more information in news
8 speech. News speech makes use of local range information, in the form of features
9 relative to the speaker's baseline, whereas pitch features in meetings capture pitch
10 resets across inter-word boundaries. Interestingly, duration features, while corre-
11 lated with sentence boundaries in both genres, are relatively more useful in meet-
12 ings. Inspection reveals that in news speech, a problem for duration features is
13 that they indicate many other locations, including prominent syllables and sub-
14 sentential boundaries. Energy features appear to be less important than pause,
15 pitch, and duration features in both genres.

16 Sentence segmentation is one of a number of tasks in which lexical and prosodic
17 features can be combined for better performance. Based on results found here,
18 we conclude that feature selection can produce similar or even better performance
19 results, but that the particular features depend on the speech genre. Although in this
20 case training data was available for both genres, information about which features
21 benefit which genre should be even more important when adapting models to data
22 for which little or no matched training data is available.

23 Acknowledgments

24 This material is based upon work supported by the Defense Advanced Research
25 Projects Agency (DARPA) under contract No. HR0011-06-C-0023 and contract
26 No. NBCHD030010. Any opinions, findings, and conclusions or recommendations
27 expressed in this material are those of the authors and do not necessarily reflect
28 the view of DARPA. The authors would like to thank Matthias Zimmermann, Yang
29 Liu, and Mathew Magimai Doss for their help and suggestions.

References

- 31 [1] D. Biber, *Variation across Speech and Writing*, Cambridge University Press, Cam-
32 bridge, 1988.
- 33 [2] S. Cuendet, D. Hakkani-Tür, and G. Tur, Model adaptation for sentence unit seg-
34 mentation from speech, in *Proceedings of SLT*, Aruba, 2006.
- 35 [3] D. Jones, W. Shen, E. Shriberg, A. Stolcke, T. Kamm and D. Reynolds, Two experi-
36 ments comparing reading with listening for human processing of conversational tele-
37 phone speech, in *Proceedings of EUROSPEECH*, 2005, pp. 1145–1148.
- 38 [4] Y. Liu, E. Shriberg, A. Stolcke, B. Peskin, J. Ang, D. Hillard, M. Ostendorf, M. Toma-
39 lin, P. Woodland and M. Harper, Structural metadata research in the EARS program,
in *Proceedings of ICASSP*, 2005.

12 *S. Cuendet et al.*

- 1 [5] J. Makhoul, A. Baron, I. Bulyko, L. Nguyen, L. Ramshaw, D. Stallard, R. Schwartz
3 and B. Xiang, The effects of speech recognition and punctuation on information
extraction performance, in *Proc. of Interspeech*, Lisbon, 2005, pp. 57–60.
- 5 [6] J. Mrozinski, E. W. D. Whittaker, P. Chatain and S. Furui, Automatic sentence seg-
mentation of speech for automatic summarization, in *Proc. ICASSP*, Vol. 1, Philadel-
7 phia, PA, 2005, pp. I–I.
- 9 [7] B. Roark, Y. Liu, M. Harper, R. Stewart, M. Lease, M. Snover, I. Shafran, B. Dorr,
J. Hale, A. Krasnyanskaya and L. Yung, Reranking for sentence boundary detection
in conversational speech, in *Proceedings of ICASSP*, Toulouse, France, 2006.
- 11 [8] R. E. Schapire and Y. Singer, Boostexter: A boosting-based system for text catego-
rization, *Machine Learning* **39**(2/3) (2000) 135–168.
- 13 [9] E. Shriberg, R. Dhillon, S. Bhagat, J. Ang and H. Carvey, The ICSI meeting recorder
dialog act (MRDA) corpus, in *Proceedings of SigDial Workshop*, Boston, MA, 2004.
- 15 [10] E. Shriberg, A. Stolcke, D. Hakkani-Tür and G. Tur, Prosody-based automatic seg-
mentation of speech into sentences and topics, *Speech Communication* **3** (2000) 2037–
2040.
- 17 [11] S. Strassel and M. Glenn, Creating the annotated TDT-4 Y2003 evaluation corpus,
in *TDT 2003 Evaluation Workshop*, NIST, 2003.
- 19 [12] A. Venkataraman, A. Stolcke, W. Wang, D. Vergyri, V. Gadde and J. Zheng, SRIs
2004 broadcast news speech to text system, in *EARS Rich Transcription 2004 Work-*
21 *shop*, Palisades, NY, 2004.
- 23 [13] C. Wooters, J. Fung, B. Peskin and X. Anguera, Towards robust speaker segmen-
tation: ICSI-SRI Fall 2004 diarization system, in *RT-04F Workshop*, 2004.
- 25 [14] Y. Yang and J. Pedersen, A comparative study on feature selection in text catego-
rization, in *Proceedings of ICML*, Nashville, US, 1997, pp. 412–420.
- 27 [15] Q. Zhu, A. Stolcke, B. Chen and N. Morgan, Using MLP features in SRIs conversa-
tional speech recognition system, in *Proceedings of INTERSPEECH*, Lisbon, Portu-
gal, 2005, pp. 2141–2144.
- 29 [16] M. Zimmermann, D. Hakkani-Tür, J. Fung, N. Mirghafori, E. Shriberg and Y. Liu,
The ICSI+ multi-lingual sentence segmentation system, in *Proceedings of ICSLP*,
31 Pittsburgh, PA, 2006.