

# A Statistical Model for Word Discovery in Transcribed Speech

Anand Venkataraman\*

*A statistical model for segmentation and word discovery in continuous speech is presented. An incremental unsupervised learning algorithm to infer word boundaries based on this model is described. Results are also presented of empirical tests showing that the algorithm is competitive with other models that have been used for similar tasks.*

## 1. Introduction

English speech lacks the acoustic analog of blank spaces that people are accustomed to seeing between words in written text. Discovering words in continuous spoken speech is thus an interesting problem and one that has been treated at length in the literature. The problem of identifying word boundaries is particularly significant in the parsing of written text in languages that do not explicitly include spaces between words. In addition, if we assume that children start out with little or no knowledge of the inventory of words the language possesses identification of word boundaries is a significant problem in the domain of child language acquisition.<sup>1</sup> Although speech lacks explicit demarcation of word boundaries, it is undoubtedly the case that it nevertheless possesses significant other cues for word discovery. However, it is still a matter of interest to see exactly how much can be achieved without the incorporation of these other cues; that is, we are interested in the performance of a *bare-bones* language model. For example, there is much evidence that stress patterns (Jusczyk, Cutler, and Redanz 1993; Cutler and Carter 1987) and phonotactics of speech (Mattys and Jusczyk 1999) are of considerable aid in word discovery. But a bare-bones statistical model is still useful in that it allows us to quantify precise improvements in performance upon the integration of each specific cue into the model. We present and evaluate one such statistical model in this paper.<sup>2</sup>

The main contributions of this study are as follows: First, it demonstrates the applicability and competitiveness of a conservative, traditional approach for a task for which nontraditional approaches have been proposed even recently (Brent 1999; Brent and Cartwright 1996; de Marcken 1995; Elman 1990; Christiansen, Allen, and Seidenberg 1998). Second, although the model leads to the development of an algorithm that learns the lexicon in an unsupervised fashion, results of partial supervision are presented, showing that its performance is consistent with results from learning theory. Third, the study extends previous work to higher-order  $n$ -grams, specifically up to

---

\* STAR Lab, SRI International, 333 Ravenswood Ave., Menlo Park, CA 94025. E-mail: anand@speech.sri.com

1 See, however, work in Jusczyk and Hohne (1997) and Jusczyk (1997) that presents strong evidence in favor of a hypothesis that children already have a reasonably powerful and accurate lexicon at their disposal as early as 9 months of age.

2 Implementations of all the programs discussed in this paper and the input corpus are readily available upon request from the author. The programs (totaling about 900 lines) have been written in C++ to compile under Unix/Linux. The author will assist in porting it to other architectures or to versions of Unix other than Linux or SunOS/Solaris if required.

trigrams, and discusses the results in their light. Finally, results of experiments suggested in Brent (1999) regarding different ways of estimating phoneme probabilities are also reported. Wherever possible, results are averaged over 1000 repetitions of the experiments, thus removing any potential advantages the algorithm may have had due to ordering idiosyncrasies within the input corpus.

Section 2 briefly discusses related literature in the field and recent work on the same topic. The model is described in Section 3. Section 4 describes an unsupervised learning algorithm based directly on the model developed in Section 3. This section also describes the data corpus used to test the algorithms and the methods used. Results are presented and discussed in Section 5. Finally, the findings in this work are summarized in Section 6.

## 2. Related Work

While there exists a reasonable body of literature regarding text segmentation, especially with respect to languages such as Chinese and Japanese that do not explicitly include spaces between words, most of the statistically based models and algorithms tend to fall into the supervised learning category. These require the model to be trained first on a large corpus of text before it can segment its input.<sup>3</sup> It is only recently that interest in unsupervised algorithms for text segmentation seems to have gained ground. A notable exception in this regard is the work by Ando and Lee (1999) which tries to infer word boundaries from character  $n$ -gram statistics of Japanese Kanji strings. For example, a decision to insert a word boundary between two characters is made solely based on whether character  $n$ -grams adjacent to the proposed boundary are relatively more frequent than character  $n$ -grams that straddle it. This algorithm, however, is not based on a formal statistical model and is closer in spirit to approaches based on transitional probability between phonemes or syllables in speech. One such approach derives from experiments by Saffran, Newport, and Aslin (1996) suggesting that young children might place word boundaries between two syllables where the second syllable is *surprising* given the first. This technique is described and evaluated in Brent (1999). Other approaches not based on explicit probability models include those based on information theoretic criteria such as minimum description length (Brent and Cartwright 1996; de Marcken 1995) and simple recurrent networks (Elman 1990; Christiansen, Allen, and Seidenberg 1998). The maximum likelihood approach due to Olivier (1968) is probabilistic in the sense that it is geared toward explicitly calculating the most probable segmentation of each block of input utterances (see also Batchelder 1997). However, the algorithm involves heuristic steps in periodic purging of the lexicon and in the creation in the lexicon of new words. Furthermore, this approach is again not based on a formal statistical model.

Model Based Dynamic Programming, hereafter referred to as MBDP-1 (Brent 1999), is probably the most recent work that addresses exactly the same issue as that considered in this paper. Both the approach presented in this paper and Brent's MBDP-1 are unsupervised approaches based on explicit probability models. Here, we describe only Brent's MBDP-1 and direct the interested reader to Brent (1999) for an excellent review and evaluation of many of the algorithms mentioned above.

### 2.1 Brent's model-based dynamic programming method

Brent (1999) describes a **model-based approach** to inferring word boundaries in child-directed speech. As the name implies, this technique uses dynamic programming to

---

<sup>3</sup> See, for example, Zimin and Tseng (1993).

infer the best segmentation. It is assumed that the entire input corpus, consisting of a concatenation of all utterances in sequence, is a single event in probability space and that the best segmentation of each utterance is implied by the best segmentation of the corpus itself. The model thus focuses on explicitly calculating probabilities for every possible segmentation of the entire corpus, and subsequently picking the segmentation with the maximum probability. More precisely, the model attempts to calculate

$$P(\bar{w}_m) = \sum_n \sum_L \sum_f \sum_s P(\bar{w}_m | n, L, f, s) \cdot P(n, L, f, s)$$

for each possible segmentation of the input corpus where the left-hand side is the exact probability of that particular segmentation of the corpus into words  $\bar{w}_m = w_1 w_2 \cdots w_m$ ; and the sums are over all possible numbers of words  $n$ , in the lexicon, all possible lexicons  $L$ , all possible frequencies  $f$ , of the individual words in this lexicon and all possible orders of words  $s$ , in the segmentation. In practice, the implementation uses an incremental approach that computes the best segmentation of the entire corpus up to step  $i$ , where the  $i$ th step is the corpus up to and including the  $i$ th utterance. Incremental performance is thus obtained by computing this quantity anew after each segmentation  $i - 1$ , assuming however, that segmentations of utterances up to but not including  $i$  are fixed.

There are two problems with this approach. First, the assumption that the entire corpus of observed speech should be treated as a single event in probability space appears rather radical. This fact is appreciated even in Brent (1999), which states “From a cognitive perspective, we know that humans segment each utterance they hear without waiting until the corpus of all utterances they will ever hear becomes available” (p. 89). Thus, although the incremental algorithm in Brent (1999) is consistent with a developmental model, the formal statistical model of segmentation is not.

Second, making the assumption that the corpus is a single event in probability space significantly increases the computational complexity of the incremental algorithm. The approach presented in this paper circumvents these problems through the use of a conservative statistical model that is directly implementable as an incremental algorithm. In the following section, we describe the model and how its 2-gram and 3-gram extensions are adapted for implementation.

### 3. Model Description

The language model described here is a fairly standard one. The interested reader is referred to Jelinek (1997, 57–78), where a detailed exposition can be found. Basically, we seek

$$\hat{\mathbf{W}} = \operatorname{argmax}_{\mathbf{W}} P(\mathbf{W}) \quad (1)$$

$$= \operatorname{argmax}_{\mathbf{W}} \prod_{i=1}^n P(w_i | w_1, \dots, w_{i-1}) \quad (2)$$

$$= \operatorname{argmin}_{\mathbf{W}} \sum_{i=1}^n -\log P(w_i | w_1, \dots, w_{i-1}) \quad (3)$$

where  $\mathbf{W} = w_1, \dots, w_n$  with  $w_i \in \mathbf{L}$  denotes a particular string of  $n$  words belonging to a lexicon  $\mathbf{L}$ .

The usual  $n$ -gram approximation is made by grouping histories  $w_1, \dots, w_{i-1}$  into equivalence classes, allowing us to collapse contexts into histories at most  $n - 1$  words

backwards (for  $n$ -grams). Estimations of the required  $n$ -gram probabilities are then done with relative frequencies using **back-off** to lower-order  $n$ -grams when a higher-order estimate is not reliable enough (Katz 1987). Back-off is done using the Witten and Bell (1991) technique, which allocates a probability of  $N_i/(N_i + S_i)$  to unseen  $i$ -grams at each stage, with the final back-off from unigrams being to an open vocabulary where word probabilities are calculated as a normalized product of phoneme or letter probabilities. Here,  $N_i$  is the number of distinct  $i$ -grams and  $S_i$  is the sum of their frequencies. The model can be summarized as follows:

$$P(w_i | w_{i-2}, w_{i-1}) = \begin{cases} \frac{S_3}{N_3 + S_3} \frac{C(w_{i-2}, w_{i-1}, w_i)}{C(w_{i-1}, w_i)} & \text{if } C(w_{i-2}, w_{i-1}, w_i) > 0 \\ \frac{N_3}{N_3 + S_3} P(w_i | w_{i-1}) & \text{otherwise} \end{cases} \quad (4)$$

$$P(w_i | w_{i-1}) = \begin{cases} \frac{S_2}{N_2 + S_2} \frac{C(w_{i-1}, w_i)}{C(w_i)} & \text{if } C(w_{i-1}, w_i) > 0 \\ \frac{N_2}{N_2 + S_2} P(w_i) & \text{otherwise} \end{cases} \quad (5)$$

$$P(w_i) = \begin{cases} \frac{C(w_i)}{N_1 + S_1} & \text{if } C(w_i) > 0 \\ \frac{N_1}{N_1 + S_1} P_\Sigma(w_i) & \text{otherwise} \end{cases} \quad (6)$$

$$P_\Sigma(w_i) = \frac{r(\#) \prod_{j=1}^{k_i} r(w_i[j])}{1 - r(\#)} \quad (7)$$

where  $C()$  denotes the count or frequency function,  $k_i$  denotes the length of word  $w_i$ , excluding the sentinel character  $\#$ ,  $w_i[j]$  denotes its  $j$ th phoneme, and  $r()$  denotes the relative frequency function. The normalization by dividing using  $1 - r(\#)$  in Equation (7) is necessary because otherwise

$$\sum_w P(w) = \sum_{i=1}^{\infty} (1 - P(\#))^i P(\#) \quad (8)$$

$$= 1 - P(\#) \quad (9)$$

Since we estimate  $P(w[j])$  by  $r(w[j])$ , dividing by  $1 - r(\#)$  will ensure that  $\sum_w P(w) = 1$ .

#### 4. Method

As in Brent (1999), the model described in Section 3 is presented as an incremental learner. The only knowledge built into the system at start-up is the phoneme table, with a uniform distribution over all phonemes, including the sentinel phoneme. The learning algorithm considers each utterance in turn and computes the most probable segmentation of the utterance using a Viterbi search (Viterbi 1967) implemented as a dynamic programming algorithm, as described in Section 4.2. The most likely placement of word boundaries thus computed is committed to before the algorithm considers the next utterance presented. Committing to a segmentation involves learning unigram, bigram, and trigram frequencies, as well as phoneme frequencies, from the inferred words. These are used to update the respective tables.

To account for effects that any specific ordering of input utterances may have on the segmentations that are output, the performance of the algorithm is averaged over 1000 runs, with each run receiving as input a random permutation of the input corpus.

##### 4.1 The input corpus

The corpus, which is identical to the one used by Brent (1999), consists of orthographic transcripts made by Bernstein-Ratner (1987) from the CHILDES collection (MacWhin-

**Table 1**

Twenty randomly chosen utterances from the input corpus with their orthographic transcripts. See the appendix for a list of the ASCII representations of the phonemes.

Phonemic Transcription	Orthographic English text
hQ slli 6v mi	How silly of me
IUk D*z D6 b7 wIT hIz h&t	Look, there's the boy with his hat
9 TINK 9 si 6nADR bUk	I think I see another book
tu	Two
DIs wAn	This one
r9t WEn De wOk	Right when they walk
huz an D6 tEl6fon &IIs	Who's on the telephone, Alice?
sIt dQn	Sit down
k&n yu fid It tu D6 dOgi	Can you feed it to the doggie?
D*	There
du yu si hIm h(	Do you see him here?
IUk	Look
yu want It In	You want it in
W* dId It go	Where did it go?
&nd WAt # Doz	And what are those?
h9 m6ri	Hi Mary
oke Its 6 cIk	Okay it's a chick
y& IUk WAt yu dId	Yeah, look what you did
oke	Okay
tek It Qt	Take it out

ney and Snow 1985). The speakers in this study were nine mothers speaking freely to their children, whose ages averaged 18 months (range 13–21). Brent and his colleagues transcribed the corpus phonemically (using the ASCII phonemic representation in the appendix to this paper) ensuring that the number of subjective judgments in the pronunciation of words was minimized by transcribing every occurrence of the same word identically. For example, “look”, “drink”, and “doggie” were always transcribed “IUk”, “drINK”, and “dOgi” regardless of where in the utterance they occurred and which mother uttered them in what way. Thus transcribed, the corpus consists of a total of 9790 such utterances and 33,399 words, and includes one space after each word and one newline after each utterance. For purposes of illustration, Table 1 lists the first 20 such utterances from a random permutation of this corpus.

It should be noted that the choice of this particular corpus for experimentation is motivated purely by its use in Brent (1999). As has been pointed out by reviewers of an earlier version of this paper, the algorithm is equally applicable to plain text in English or other languages. The main advantage of the CHILDES corpus is that it allows for ready comparison with results hitherto obtained and reported in the literature. Indeed, the relative performance of all the algorithms discussed is mostly unchanged when tested on the 1997 Switchboard telephone speech corpus with disfluency events removed.

#### 4.2 Algorithm

The dynamic programming algorithm finds the most probable word sequence for each input utterance by assigning to each segmentation a score equal to its probability, and committing to the segmentation with the highest score. In practice, the implementation computes the negative logarithm of this score and thus commits to the segmentation with the least negative logarithm of its probability. The algorithm for the unigram

```

BEGIN
  Input (by ref) utterance u[0..n] where u[i] are the characters in it.

  bestSegpoint := n;
  bestScore := evalWord(u[0..n]);
  for i from 0 to n-1; do
    subUtterance := copy(u[0..i]);
    word := copy(u[i+1..n]);
    score := evalUtterance(subUtterance) + evalWord(word);
    if (score < bestScore); then
      bestScore = score;
      bestSegpoint := i;
    fi
  done
  insertWordBoundary(u, bestSegpoint)
  return bestScore;
END

```

**Figure 1. Algorithm: evalUtterance**

Recursive optimization algorithm to find the best segmentation of an input utterance using the unigram language model described in this paper.

```

BEGIN
  Input (by reference) word w[0..k] where w[i] are the phonemes in it.

  score = 0;
  if L.frequency(word) == 0; then {
    escape = L.size()/(L.size()+L.sumFrequencies())
    P_0 = phonemes.relativeFrequency('#');
    score = -log(escape) -log(P_0/(1-P_0));
    for each w[i]; do
      score -= log(phonemes.relativeFrequency(w[i]));
    done
  } else {
    P_w = L.frequency(w)/(L.size()+L.sumFrequencies());
    score = -log(P_w);
  }
  return score;
END

```

**Figure 2. Function: evalWord**

The function to compute  $-\log P(w)$  of an input word  $w$ .  $L$  stands for the lexicon object. If the word is novel, then it backs off to using a distribution over the phonemes in the word.

language model is presented in recursive form in Figure 1 for readability. The actual implementation, however, used an iterative version. The algorithm to evaluate the back-off probability of a word is given in Figure 2. Algorithms for bigram and trigram language models are straightforward extensions of that given for the unigram model. Essentially, the algorithm description can be summed up semiformaly as follows: For each input utterance  $u$ , we evaluate every possible way of segmenting it as  $u = u' + w$  where  $u'$  is a subutterance from the beginning of the original utterance up to some point within it and  $w$ —the lexical difference between  $u$  and  $u'$ —is treated as a word. The subutterance  $u'$  is itself evaluated recursively using the same algorithm. The base case for recursion when the algorithm rewinds is obtained when a subutterance cannot be split further into a smaller component subutterance and word, that is, when its length is zero. Suppose for example, that a given utterance is  $abcde$ , where the letters represent phonemes. If  $\text{seg}(x)$  represents the best segmentation of the utterance  $x$  and

$\mathbf{word}(x)$  denotes that  $x$  is treated as a word, then

$$\mathbf{seg}(abcde) = \mathbf{best\ of} \begin{cases} \mathbf{word}(abcde) \\ \mathbf{seg}(a) + \mathbf{word}(bcde) \\ \mathbf{seg}(ab) + \mathbf{word}(cde) \\ \mathbf{seg}(abc) + \mathbf{word}(de) \\ \mathbf{seg}(abcd) + \mathbf{word}(e) \end{cases}$$

The **evalUtterance** algorithm in Figure 1 does precisely this. It initially assumes the entire input utterance to be a word on its own by assuming a single segmentation point at its right end. It then compares the log probability of this segmentation successively to the log probabilities of segmenting it into all possible subutterance-word pairs.

The implementation maintains four separate tables internally, one each for unigrams, bigrams, and trigrams, and one for phonemes. When the procedure is initially started, all the internal  $n$ -gram tables are empty. Only the phoneme table is populated with equipossible phonemes. As the program considers each utterance in turn and commits to its best segmentation according to the **evalUtterance** algorithm, the various internal  $n$ -gram tables are updated correspondingly. For example, after some utterance  $abcde$  is segmented into  $a\ bc\ de$ , the unigram table is updated to increment the frequencies of the three entries  $a$ ,  $bc$ , and  $de$ , each by 1, the bigram table to increment the frequencies of the adjacent bigrams  $a\ bc$  and  $bc\ de$ , and the trigram table to increment the frequency of the trigram  $a\ bc\ de$ .<sup>4</sup> Furthermore, the phoneme table is updated to increment the frequencies of each of the phonemes in the utterance, including one sentinel for each word inferred.<sup>5</sup> Of course, incrementing the frequency of a currently unknown  $n$ -gram is equivalent to creating a new entry for it with frequency 1. Note that the very first utterance is necessarily segmented as a single word. Since all the  $n$ -gram tables are empty when the algorithm attempts to segment it, all probabilities are necessarily computed from the level of phonemes up. Thus, the more words in the segmentation of the first utterance, the more sentinel characters will be included in the probability calculation, and so the lesser the corresponding segmentation probability will be. As the program works its way through the corpus,  $n$ -grams inferred correctly by virtue of their relatively greater preponderance compared to noise tend to dominate their respective  $n$ -gram distributions and thus dictate how future utterances are segmented.

One can easily see that the running time of the program is  $O(mn^2)$  in the total number of utterances ( $m$ ) and the length of each utterance ( $n$ ), assuming an efficient implementation of a hash table allowing nearly constant lookup time is available. Since individual utterances typically tend to be small, especially in child-directed speech as evidenced in Table 1, the algorithm practically approximates to a linear time procedure. A single run over the entire corpus typically completes in under 10 seconds on a 300 MHz i686-based PC running Linux 2.2.5-15.

Although the algorithm is presented as an unsupervised learner, a further experiment to test the responsiveness of each algorithm to training data is also described here: The procedure involves reserving for training increasing amounts of the input corpus, from 0% in steps of approximately 1% (100 utterances). During the training period, the algorithm is presented with the correct segmentation of the input utterance, which it uses to update trigram, bigram, unigram, and phoneme frequencies as

<sup>4</sup> Amending the algorithm to include special markers for the start and end of each utterance was not found to make a significant difference in its performance.

<sup>5</sup> In this context, see also Section 5.2 regarding experiments conducted to investigate different ways of estimating phoneme probabilities.

required. After the initial training segment of the input corpus has been considered, subsequent utterances are then processed in the normal way.

### 4.3 Scoring

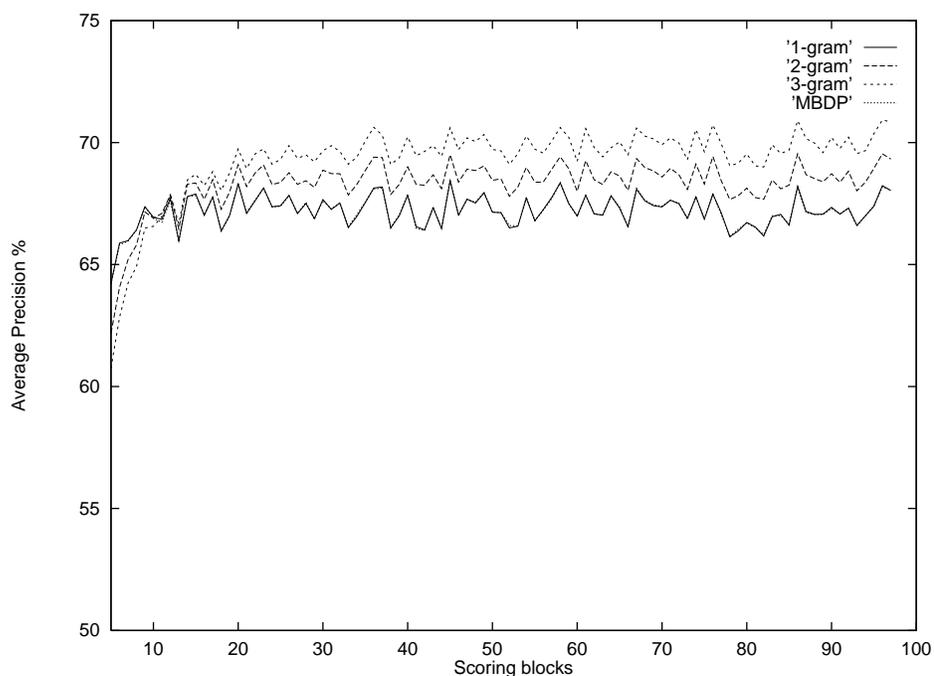
In line with the results reported in Brent (1999), three scores are reported — **precision**, **recall**, and **lexicon precision**. Precision is defined as the proportion of predicted words that are actually correct. Recall is defined as the proportion of correct words that were predicted. Lexicon precision is defined as the proportion of words in the predicted lexicon that are correct. In addition to these, the number of correct and incorrect words in the predicted lexicon were computed, but they are not graphed here because lexicon precision is a good indicator of both.

Precision and recall scores were computed incrementally and cumulatively within scoring blocks, each of which consisted of 100 utterances. These scores were computed and averaged only for the utterances within each block scored, and thus represent the performance of the algorithm only on the block scored, occurring in the exact context among the other scoring blocks. Lexicon scores carried over blocks cumulatively. In cases where the algorithm used varying amounts of training data, precision, recall, and lexical precision scores are computed over the entire corpus. All scores are reported as percentages.

## 5. Results

Figures 3–5 plot the precision, recall, and lexicon precision of the proposed algorithm for each of the unigram, bigram, and trigram models against the MBDP-1 algorithm. Although the graphs compare the performance of the algorithm with only one published result in the field, comparison with other related approaches is implicitly available. Brent (1999) reports results of running the algorithms due to Elman (1990) and Olivier (1968), as well as algorithms based on mutual information and transitional probability between pairs of phonemes, over exactly the same corpus. These are all shown to perform significantly worse than Brent's MBDP-1. The random baseline algorithm in Brent (1999), which consistently performs with under 20% precision and recall, is not graphed for the same reason. This baseline algorithm offers an important advantage: It knows the exact number of word boundaries, even though it does not know their locations. Brent argued that if MBDP-1 performs as well as this random baseline, then at the very least, it suggests that the algorithm is able to infer information equivalent to knowing the right number of word boundaries. A second important reason for not graphing the algorithms with worse performance was that the scale on the vertical axis could be expanded significantly by their omission, thus allowing distinctions between the plotted graphs to be seen more clearly.

The plots originally given in Brent (1999) are over blocks of 500 utterances. However, because they are a result of running the algorithm on a single corpus, there is no way of telling whether the performance of the algorithm was influenced by any particular ordering of the utterances in the corpus. A further undesirable effect of reporting results of a run on exactly one ordering of the input is that there tends to be too much variation between the values reported for consecutive scoring blocks. To mitigate both of these problems, we report averaged results from running the algorithms on 1000 random permutations of the input data. This has the beneficial side effect of allowing us to plot with higher granularity, since there is much less variation in the precision and recall scores. They are now clustered much closer to their mean values in each block, allowing a block size of 100 to be used to score the output. These plots are thus much more readable than those obtained without such averaging of the results.

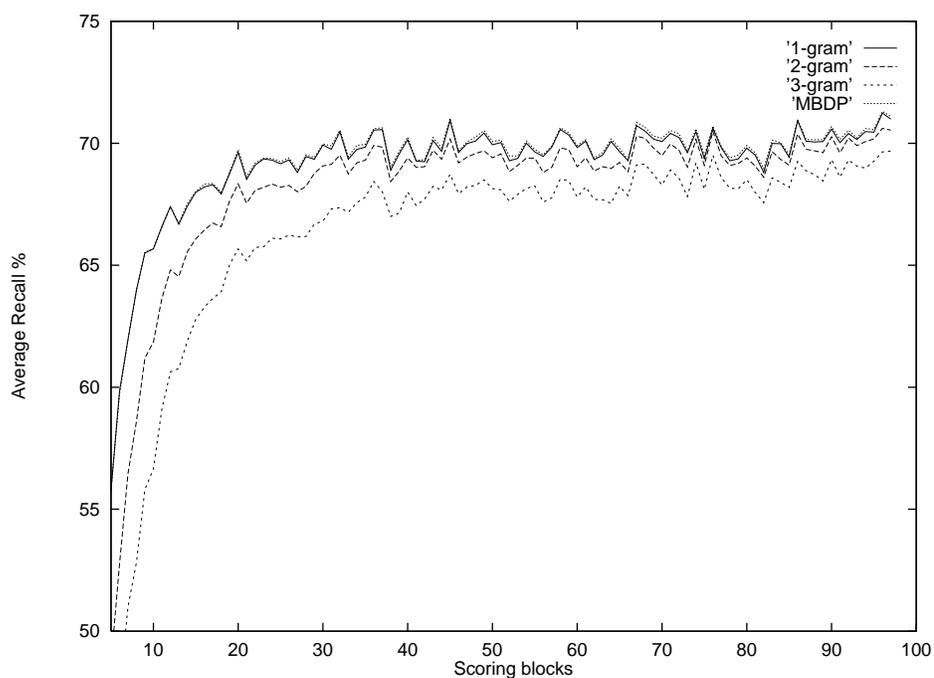


**Figure 3**

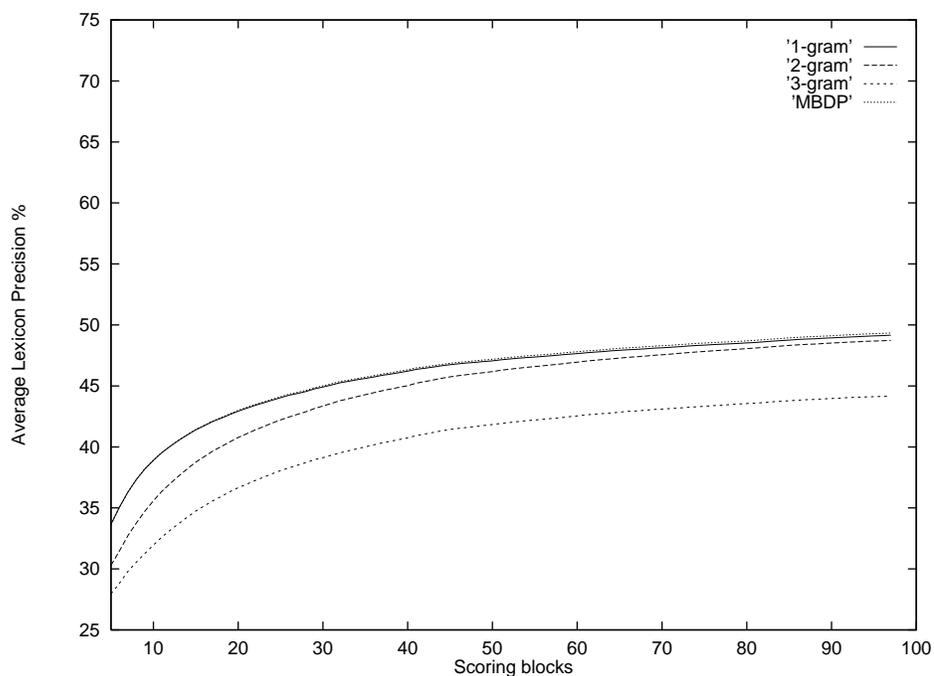
Averaged precision. This is a plot of the segmentation precision over 100 utterance blocks averaged over 1000 runs, each using a random permutation of the input corpus. Precision is defined as the percentage of identified words that are correct, as measured against the target data. The horizontal axis represents the number of blocks of data scored, where each block represents 100 utterances. The plots show the performance of the 1-gram, 2-gram, 3-gram, and MBDP-1 algorithms. The plot for MBDP-1 is not visible because it coincides almost exactly with the plot for the 1-gram model. Discussion of this level of similarity is provided in Section 5.5. The performance of related algorithms due to Elman (1990), Olivier (1968) and others is implicitly available in this and the following graphs since Brent (1999) demonstrates that these all perform significantly worse than MBDP-1.

One may object that the original transcripts carefully preserve the order of utterances directed at children by their mothers, and hence randomly permuting the corpus would destroy the fidelity of the simulation. However, as we argued, the permutation and averaging does have significant beneficial side effects, and if anything, it only eliminates from the point of view of the algorithms the important advantage that real children may be given by their mothers through a specific ordering of the utterances. In any case, we have found no significant difference in performance between the permuted and unpermuted cases as far as the various algorithms are concerned.

In this context, it would be interesting to see how the algorithms would fare if the utterances were in fact favorably ordered, that is, in order of increasing length. Clearly, this is an important advantage for all the algorithms concerned. Section 5.3 presents the results of an experiment based on a generalization of this situation, where instead of ordering the utterances favorably, we treat an initial portion of the corpus as a training component, effectively giving the algorithms free word boundaries after each word.



**Figure 4**  
Averaged recall over 1000 runs, each using a random permutation of the input corpus.



**Figure 5**  
Averaged lexicon precision over 1000 runs, each using a random permutation of the input corpus.

### 5.1 Discussion

Clearly, the performance of the present model is competitive with MBDP-1 and, as a consequence, with other algorithms evaluated in Brent (1999). However, note that the model proposed in this paper has been developed entirely along conventional lines and has not made the somewhat radical assumption that the entire observed corpus is a single event in probability space. Assuming that the corpus consists of a single event, as Brent does, requires the explicit calculation of the probability of the lexicon in order to calculate the probability of any single segmentation. This calculation is a nontrivial task since one has to sum over all possible orders of words in  $L$ . This fact is recognized in Brent (1999, Appendix), where the expression for  $P(L)$  is derived as an approximation. One can imagine then that it would be correspondingly more difficult to extend the language model in Brent (1999) beyond the case of unigrams. In practical terms, recalculating lexicon probabilities before each segmentation increases the running time of an implementation of the algorithm. Although all the algorithms discussed tend to complete within one minute on the reported corpus, MBDP-1's running time is quadratic in the number of utterances, while the language models presented here enable computation in almost linear time. The typical running time of MBDP-1 on the 9790-utterance corpus averages around 40 seconds per run on a 300 MHz i686 PC while the 1-gram, 2-gram, and 3-gram models average around 7, 10, and 14 seconds, respectively.

Furthermore, the language models presented in this paper estimate probabilities as relative frequencies, using commonly used back-off procedures, and so they do not assume any priors over integers. However, MBDP-1 requires the assumption of two distributions over integers, one to pick a number for the size of the lexicon and another to pick a frequency for each word in the lexicon. Each is assumed such that the probability of a given integer  $P(i)$  is given by  $\frac{6}{\pi^2 i^2}$ . We have since found some evidence suggesting that the choice of a particular prior does not offer any significant advantage over the choice of any other prior. For example, we have tried running MBDP-1 using  $P(i) = 2^{-i}$  and still obtained comparable results. It should be noted, however, that no such subjective prior needs to be chosen in the model presented in this paper.

The other important difference between MBDP-1 and the present model is that MBDP-1 assumes a uniform distribution over all possible word orders. That is, in a corpus that contains  $n_k$  distinct words such that the frequency in the corpus of the  $i$ th distinct word is given by  $f_k(i)$ , the probability of any one ordering of the words in the corpus is

$$\frac{\prod_{i=1}^{n_k} f_k(i)!}{k!}$$

because the number of unique orderings is precisely the reciprocal of the above quantity. Brent (1999) mentions that there may well be efficient ways of using  $n$ -gram distributions in the MBDP-1 model. The framework presented in this paper is a formal statement of a model that lends itself to such easy  $n$ -gram extendibility using the back-off scheme proposed here. In fact, the results we present are direct extensions of the unigram model into bigrams and trigrams.

In this context, an intriguing feature of the results is worth discussing here. Note that while, with respect to precision, the 3-gram model is better than the 2-gram model, which in turn is better than the 1-gram model, with respect to recall, their performance is exactly the opposite. A possible explanation of this behavior is as follows: Since the 3-gram model places greatest emphasis on word triples, which are relatively less frequent than words and word pairs, it has, of all the models, the least evidence available

to infer word boundaries from the observed data. Even though back-off is performed for bigrams when a trigram is not found, there is a cost associated with such backing off—this is the extra fractional factor  $N_3/(N_3 + S_3)$  in the calculation of the segmentation's probability. Consequently, the 3-gram model is the most conservative in its predictions. When it does infer a word boundary, it is likely to be correct. This contributes to its relatively higher precision since precision is a measure of the proportion of inferred boundaries that were correct. More often than not, however, when the 3-gram model does not have enough evidence to infer words, it simply outputs the default segmentation, which is a single word (the entire utterance) instead of more than one incorrectly inferred one. This contributes to its poorer recall since recall is an indicator of the number of words the model fails to infer. Poorer lexicon precision is likewise explained. Because the 3-gram model is more conservative, it infers new words only when there is strong evidence for them. As a result many utterances are inserted as whole words into its lexicon, thereby contributing to decreased lexicon precision. The framework presented here thus provides a systematic way of trading off precision against recall or vice-versa. Models utilizing higher-order  $n$ -grams give better recall at the expense of precision.

## 5.2 Estimation of phoneme probabilities

Brent (1999, 101) suggests that it might be worthwhile to study whether learning phoneme probabilities from distinct lexical entries yields better results than learning these probabilities from the input corpus. That is, rather than inflating the probability of the phoneme “th” in *the* by the preponderance of *the* and *the*-like words in actual speech, it is better to control it by the number of such distinct words. Presented below are an initial analysis and experimental results in this regard.

Assume the existence of some function  $\Psi_X: \mathbf{N} \rightarrow \mathbf{N}$  that maps the size,  $n$ , of a corpus  $\mathbf{C}$ , onto the size of some subset  $\mathbf{X}$  of  $\mathbf{C}$  we may define. If this subset  $\mathbf{X} = \mathbf{C}$ , then  $\Psi_{\mathbf{C}}$  is the identity function, and if  $\mathbf{X} = \mathbf{L}$  is the set of distinct words in  $\mathbf{C}$ , we have  $\Psi_{\mathbf{L}}(n) = |\mathbf{L}|$ .

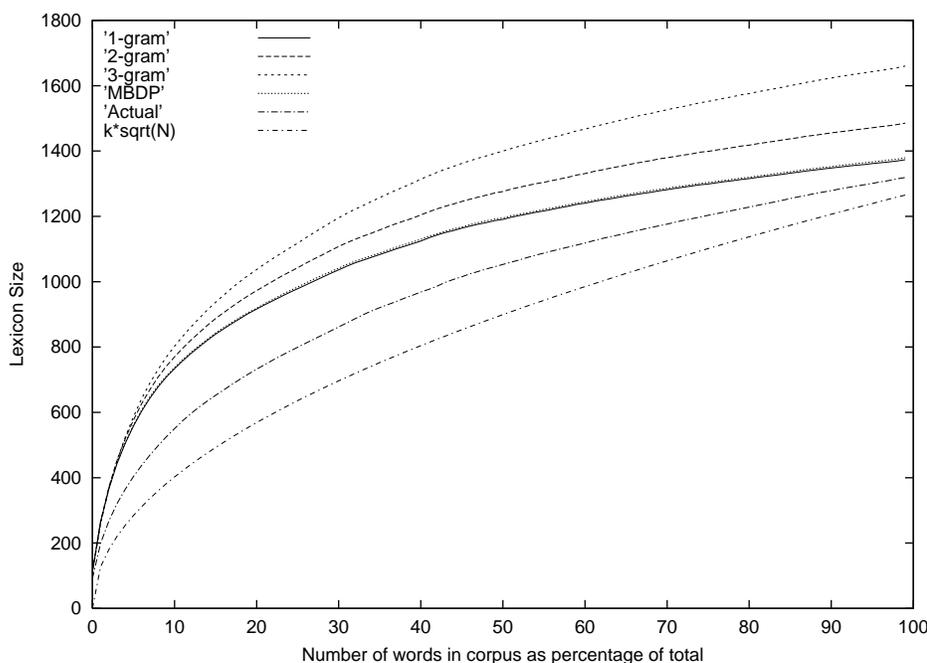
Let  $l_X$  be the average number of phonemes per word in  $\mathbf{X}$  and let  $E_{aX}$  be the average number of occurrences of phoneme  $a$  per word in  $\mathbf{X}$ . Then we may estimate the probability of an arbitrary phoneme  $a$  from  $\mathbf{X}$  as follows.

$$\begin{aligned} P(a | \mathbf{X}) &= \frac{C(a | \mathbf{X})}{\sum_{a_i} C(a_i | \mathbf{X})} \\ &= \frac{E_{aX} \Psi_X(N)}{l_X \Psi_X(N)} \end{aligned}$$

where, as before,  $C(a | \mathbf{X})$  is the count function that gives the frequency of phoneme  $a$  in  $\mathbf{X}$ . If  $\Psi_X$  is deterministic, we can then write

$$P(a | \mathbf{X}) = \frac{E_{aX}}{l_X} \quad (10)$$

Our experiments suggest that  $E_{aL} \sim E_{aC}$  and that  $l_L \sim l_C$ . We are thus led to suspect that estimates should be roughly the same regardless of whether probabilities are estimated from  $\mathbf{L}$  or  $\mathbf{C}$ . This is indeed borne out by the results we present below. Of course, this is true only if there exists, as we assumed, some deterministic function  $\Psi_L$  and this may not necessarily be the case. There is, however, some evidence that the number of distinct words in a corpus can be related to the total number of words in



**Figure 6**

Plot shows the rate of growth of the lexicon with increasing corpus size as percentage of total size. *Actual* is the actual number of distinct words in the input corpus. *1-gram*, *2-gram*, *3-gram* and *MBDP* plot the size of the lexicon as inferred by each of the algorithms. It is interesting that the rates of lexicon growth are roughly similar to each other regardless of the algorithm used to infer words and that they may all potentially be modeled by a function such as  $k\sqrt{N}$  where  $N$  is the corpus size.

the corpus in this way. In Figure 6 the rate of lexicon growth is plotted against the proportion of the corpus size considered. The values for lexicon size were collected using the Unix filter

```
cat $*|tr ' ' \\012|awk '{print (L[$0]++)? v : ++v;}'
```

and smoothed by averaging over 100 runs, each on a separate permutation of the input corpus. The plot strongly suggests that the lexicon size can be approximated by a deterministic function of the corpus size. It is interesting that the shape of the plot is roughly the same regardless of the algorithm used to infer words, suggesting that all the algorithms segment *word-like* units that share at least some statistical properties with actual words.

Table 2 summarizes our empirical findings in this regard. For each model—namely, 1-gram, 2-gram, 3-gram and MBDP-1—we test all three of the following possibilities:

1. Always use a uniform distribution over phonemes.
2. Learn the phoneme distribution from the lexicon.
3. Learn the phoneme distribution from the corpus, that is, from all words, whether distinct or not.

**Table 2**

Summary of results from each of the algorithms for each of the following cases: Lexicon–Phoneme probabilities estimated from the lexicon, Corpus–Phoneme probabilities estimated from input corpus and Uniform–Phoneme probabilities assumed uniform and constant.

Precision				
	1-gram	2-gram	3-gram	MBDP
<b>Lexicon</b>	67.7	68.08	68.02	67
<b>Corpus</b>	66.25	66.68	68.2	66.46
<b>Uniform</b>	58.08	64.38	65.64	57.15

Recall				
	1-gram	2-gram	3-gram	MBDP
<b>Lexicon</b>	70.18	68.56	65.07	69.39
<b>Corpus</b>	69.33	68.02	66.06	69.5
<b>Uniform</b>	65.6	69.17	67.23	65.07

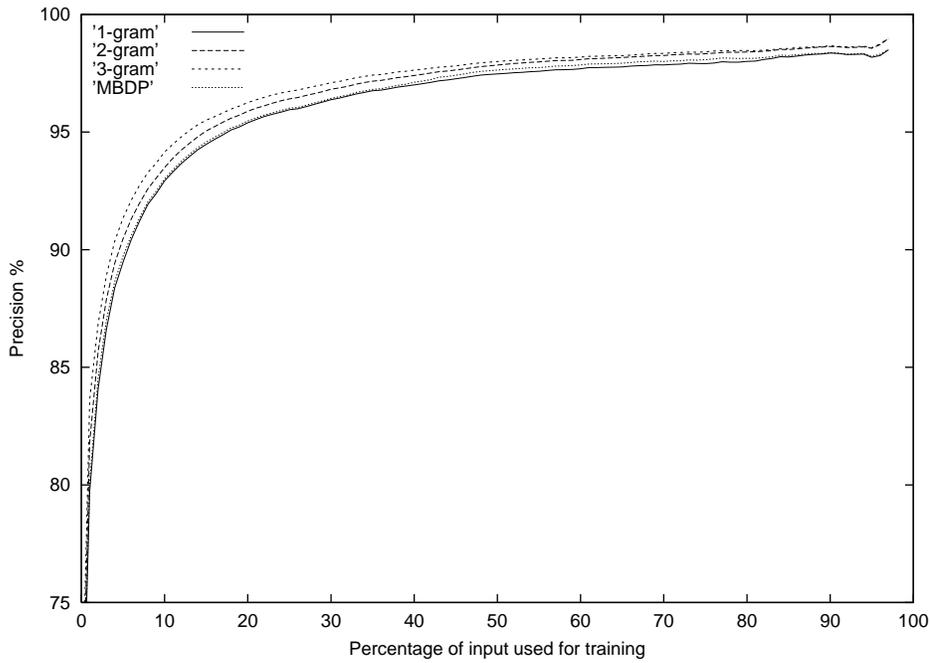
  

Lexicon Precision				
	1-gram	2-gram	3-gram	MBDP
<b>Lexicon</b>	52.85	54.45	47.32	53.56
<b>Corpus</b>	52.1	54.96	49.64	52.36
<b>Uniform</b>	41.46	52.82	50.8	40.89

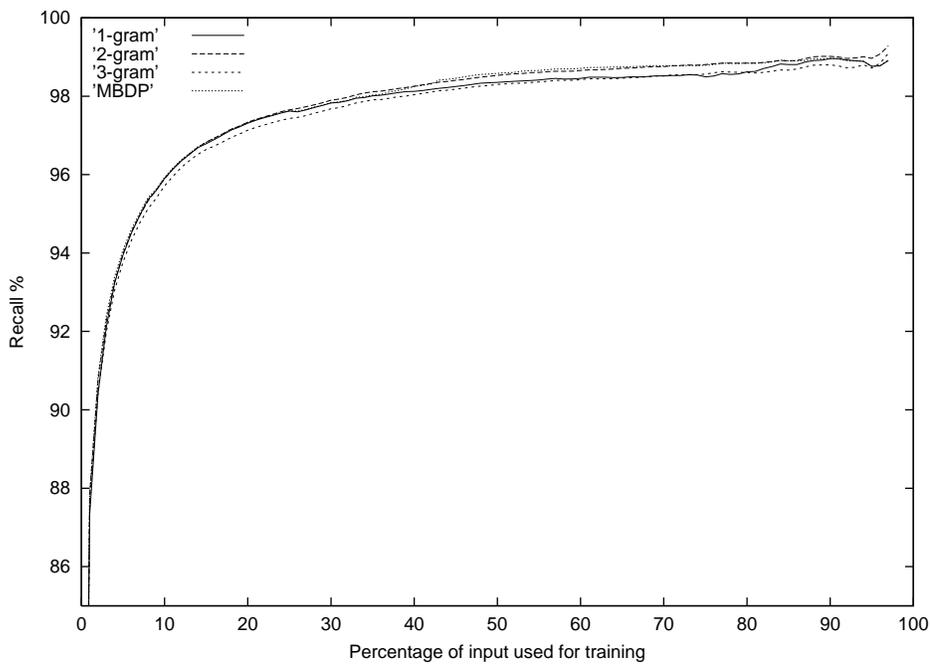
The row labeled *Lexicon* lists scores on the entire corpus from a program that learned phoneme probabilities from the lexicon. The row labeled *Corpus* lists scores from a program that learned these probabilities from the input corpus, and the row labeled *Uniform* lists scores from a program that just assumed uniform phoneme probabilities throughout. While the performance is clearly seen to suffer when a uniform distribution over phonemes is assumed, whether the distribution is estimated from the lexicon or the corpus does not seem to make any significant difference. These results lead us to believe that, from an empirical point of view, it really does not matter whether phoneme probabilities are estimated from the corpus or the lexicon. Intuitively, however, it seems that the right approach ought to be one that estimates phoneme frequencies from the corpus data since frequent words ought to have a greater influence on the phoneme distribution than infrequent ones.

### 5.3 Responsiveness to training

It is interesting to compare the responsiveness of the various algorithms to the effect of training data. Figures 7–8 plot the results (precision and recall) over the whole input corpus, that is,  $\text{blocksize} = \infty$ , as a function of the initial proportion of the corpus reserved for training. This is done by dividing the corpus into two segments, with an initial training segment being used by the algorithm to learn word, bigram, trigram, and phoneme probabilities, and the second segment actually being used as the test data. A consequence of this is that the amount of data available for testing becomes progressively smaller as the percentage reserved for training grows. So the significance of the test diminishes correspondingly. We can assume that the plots cease to be meaningful and interpretable when more than about 75% (about 7500 utterances) of the corpus is used for training. At 0%, there is no training information for any



**Figure 7** Responsiveness of the algorithm to training information. The horizontal axis represents the initial percentage of the data corpus that was used for training the algorithm. This graph shows the improvement in precision with training size.



**Figure 8** Improvement in recall with training size.

**Table 3**

Errors in the output of a fully trained 3-gram language model. Erroneous segmentations are shown in boldface.

#	3-gram output	Target
3482	... in the <b>doghouse</b>	... in the dog house
5572	<b>aclock</b>	a clock
5836	that's <b>alright</b>	that's all right
7602	that's right it's a <b>hairbrush</b>	that's right it's a hair brush

algorithm and the scores are identical to those reported earlier. We increase the amount of training data in steps of approximately 1% (100 utterances). For each training set size, the results reported are averaged over 25 runs of the experiment, each over a separate random permutation of the corpus. As before, this was done both to correct for ordering idiosyncrasies, and to smooth the graphs to make them easier to interpret.

We interpret Figures 7 and 8 as suggesting that the performance of all algorithms discussed here can be boosted significantly with even a small amount of training. It is noteworthy and reassuring to see that, as one would expect from results in computational learning theory (Haussler 1988), the number of training examples required to obtain a desired value of precision  $p$ , appears to grow with  $1/(1-p)$ . The intriguing reversal in the performance of the various  $n$ -gram models with respect to precision and recall is again seen here and the explanation for this too is the same as discussed earlier. We further note, however, that the difference in performance between the different models tends to narrow with increasing training size; that is, as the amount of evidence available to infer word boundaries increases, the 3-gram model rapidly catches up with the others in recall and lexicon precision. It is likely, therefore, that with adequate training data, the 3-gram model might be the most suitable one to use. The following experiment lends support to this conjecture.

#### 5.4 Fully trained algorithms

The preceding discussion raises the question of what would happen if the percentage of input used for training was extended to the limit, that is, to 100% of the corpus. This precise situation was tested in the following way: The entire corpus was concatenated onto itself; the models were then trained on the first half and tested on the second half of the corpus thus augmented. Although the unorthodox nature of this procedure does not allow us to attach all that much significance to the outcome, we nevertheless find the results interesting enough to warrant some mention, and we thus discuss here the performance of each of the four algorithms on the test segment of the input corpus (the second half). As one would expect from the results of the preceding experiments, the trigram language model outperforms all others. It has a precision and recall of 100% on the test input, except for exactly four utterances. These four utterances are shown in Table 3, retranscribed into plain English.

Intrigued as to why these errors occurred, we examined the corpus, only to find erroneous transcriptions in the input. *dog house* is transcribed as a single word “dOghQs” in utterance 614, and as two words elsewhere. Likewise, *o'clock* is transcribed “6klAK” in utterance 5917, *alright* is transcribed “OlR9t” in utterance 3937, and *hair brush* is transcribed “h\*brAS” in utterances 4838 and 7037. Elsewhere in the corpus, these are transcribed as two words.

The erroneous segmentations in the output of the 2-gram language model are shown in Table 4. As expected, the effect of reduced history is apparent from an in-

**Table 4**

Errors in the output of a fully trained 2-gram language model. Erroneous segmentations are shown in boldface.

#	2-gram output	Target
614	you want the <b>dog house</b>	you want the doghouse
3937	thats <b>all right</b>	that's alright
5572	<b>a clock</b>	a clock
7327	look a <b>hairbrush</b>	look a hair brush
7602	that's right its a <b>hairbrush</b>	that's right its a hair brush
7681	<b>hairbrush</b>	hair brush
7849	it's called a <b>hairbrush</b>	it's called a hair brush
7853	<b>hairbrush</b>	hair brush

crease in the total number of errors. However, it is interesting to note that while the 3-gram model incorrectly segmented an incorrect transcription (utterance 5836) *that's all right* to produce *that's alright*, the 2-gram model incorrectly segmented a correct transcription (utterance 3937) *that's alright* to produce *that's all right*. The reason for this is that the bigram *that's all* is encountered relatively frequently in the corpus and this biases the algorithm toward segmenting the *all* out of *alright* when it follows *that's*. However, the 3-gram model is not likewise biased because, having encountered the exact 3-gram *that's all right* earlier, there is no back-off to try bigrams at this stage.

Similarly, it is interesting that while the 3-gram model incorrectly segments the incorrectly transcribed *dog house* into *doghouse* in utterance 3482, the 2-gram model incorrectly segments the correctly transcribed *doghouse* into *dog house* in utterance 614. In the trigram model,  $-\log P(\text{house}|\text{the}, \text{dog}) = 4.8$  and  $-\log P(\text{dog}|\text{in}, \text{the}) = 5.4$ , giving a score of 10.2 to the segmentation *dog house*. However, due to the error in transcription, the trigram *in the doghouse* is never encountered in the training data, although the bigram *the doghouse* is. Backing off to bigrams,  $-\log P(\text{doghouse}|\text{the})$  is calculated as 8.1. Hence the probability that *doghouse* is segmented as *dog house* is less than the probability that it is a single word. In the 2-gram model, however,  $-\log P(\text{dog}|\text{the})P(\text{house}|\text{dog}) = 3.7 + 3.2 = 6.9$  while  $-\log P(\text{doghouse}|\text{the}) = 7.5$ , whence *dog house* is the preferred segmentation even though the training data contained instances of all three bigrams. For errors in the output of a 1-gram model, see Table 5.

The errors in the output of Brent's fully trained MBDP-1 algorithm are not shown here because they are identical to those produced by the 1-gram model except for one utterance. This single difference is the segmentation of utterance 8999, "lItL QtlEts" (*little outlets*), which the 1-gram model segmented correctly as "lItL QtlEts", but MBDP-1 segmented as "lItL Qt lEts". In both MBDP-1 and the 1-gram model, all four words, *little*, *out*, *lets* and *outlets*, are familiar at the time of segmenting this utterance. MBDP-1 assigns a score of  $5.3 + 5.95 = 11.25$  to the segmentation *out + lets* versus a score of 11.76 to *outlets*. As a consequence, *out + lets* is the preferred segmentation. In the 1-gram language model, the segmentation *out + lets* scores  $5.31 + 5.97 = 11.28$ , whereas *outlets* scores 11.09. Consequently, it selects *outlets* as the preferred segmentation. The only thing we could surmise from this was either that this difference must have come about due to chance (meaning that this may not have occurred if certain parts of the corpus had been different in any way) or else the interplay between the different elements in the two models is too subtle to be addressed within the scope of this paper.

**Table 5**  
Errors in the output of a fully trained 1-gram language model.

#	1-gram output	Target
244	brush <b>Alice</b> 's hair	brush Alice's hair
503	you're <b>in to</b> distraction ...	you're into distraction ...
1066	you <b>my trip</b> it	you might rip it
1231	this is little <b>doghouse</b>	this is little dog house
1792	stick it <b>on to</b> there	stick it onto there
3056	... so he doesn't run <b>in to</b>	... so he doesn't run into
3094	... to be in the <b>highchair</b>	... to be in the high chair
3098	... for this <b>highchair</b>	... for this high chair
3125	... <b>already</b> ...	... all ready ...
3212	... could talk <b>in to</b> it	... could talk into it
3230	can <b>heel I</b> down on them	can he lie down on them
3476	that's a <b>doghouse</b>	that's a dog house
3482	... in the <b>doghouse</b>	... in the dog house
3923	... when <b>it's nose</b>	... when it snows
3937	that's <b>all right</b>	that's alright
4484	its about <b>mealtime s</b>	its about meal times
5328	tell him to <b>way cup</b>	tell him to wake up
5572	<b>o'clock</b>	a clock
5671	where's my little <b>hairbrush</b>	where's my little hair brush
6315	that's a <b>nye</b>	that's an i
6968	okay mommy <b>take seat</b>	okay mommy takes it
7327	look a <b>hairbrush</b>	look a hair brush
7602	that's right its a <b>hairbrush</b>	that's right its a hair brush
7607	go <b>along</b> way to find it today	go a long way to find it today
7676	mom <b>put sit</b>	mom puts it
7681	<b>hairbrush</b>	hair brush
7849	its called a <b>hairbrush</b>	its called a hair brush
7853	<b>hairbrush</b>	hair brush
8990	... in the <b>highchair</b>	... in the high chair
8994	for baby's a nice <b>highchair</b>	for baby's a nice high chair
8995	that's like a <b>highchair</b> that's right	that's like a high chair that's right
9168	he has <b>along</b> tongue	he has a long tongue
9567	you wanna go in the <b>highchair</b>	you wanna go in the high chair
9594	<b>along</b> red tongue	a long red tongue
9674	<b>doghouse</b>	dog house
9688	<b>highchair</b> again	high chair again
9689	... the <b>highchari</b>	... the high chair
9708	I have <b>along</b> tongue	I have a long tongue

### 5.5 Similarities between MBDP-1 and the 1-gram Model

The similarities between the outputs of MBDP-1 and the 1-gram model are so great that we suspect they may be capturing essentially the same nuances of the domain. Although Brent (1999) explicitly states that probabilities are not estimated for words, it turns out that implementations of MBDP-1 do end up having the same effect as estimating probabilities from relative frequencies as the 1-gram model does. The *relative probability* of a familiar word is given in Equation 22 of Brent (1999) as

$$\frac{f_k(\hat{k})}{k} \cdot \left( \frac{f_k(\hat{k}) - 1}{f_k(\hat{k})} \right)^2$$

where  $k$  is the total number of words and  $f_k(\hat{k})$  is the frequency at that point in segmentation of the  $\hat{k}$ th word. It effectively approximates to the relative frequency

$$\frac{f_k(\hat{k})}{k}$$

as  $f_k(\hat{k})$  grows. The 1-gram language model of this paper explicitly claims to use this specific estimator for the unigram probabilities. From this perspective, both MBDP-1 and the 1-gram model tend to favor the segmenting out of familiar words that do not overlap. It is interesting, however, to see exactly how much evidence each needs before such segmentation is carried out. In this context, the author recalls an anecdote recounted by a British colleague who, while visiting the USA, noted that the populace in the vicinity of his institution had grown up thinking that *Damn British* was a single word, by virtue of the fact that they had never heard the latter word in isolation. We test this particular scenario here with both algorithms. The programs are first presented with the utterance *damnbritish*. Having no evidence to infer otherwise, both programs assume that *damnbritish* is a single word and update their lexicons accordingly. The interesting question now is exactly how many instances of the word *british* in isolation each program would have to see before being able to successfully segment a subsequent presentation of *damnbritish* correctly.

Obviously, if the word *damn* is also unfamiliar, there will never be enough evidence to segment it out in favor of the familiar word *damnbritish*. Hence each program is presented next with two identical utterances, *damn*. Unless two such utterances are presented, the estimated probabilities of the familiar words *damn* and *damnbritish* will be equal; and consequently, the probability of any segmentation of *damnbritish* that contains the word *damn* will be less than the probability of *damnbritish* considered as a single word.

At this stage, we present each program with increasing numbers of utterances consisting solely of the word *british* followed by a repetition of the very first utterance—*damnbritish*. We find that MBDP-1 needs to see the word *british* on its own three times before having enough evidence to disabuse itself of the notion that *damnbritish* is a single word. In comparison, the 1-gram model is more skeptical. It needs to see the word *british* on its own seven times before committing to the right segmentation. To illustrate the inherent simplicity of the model presented here, we can show that it is easy to predict this number analytically from the 1-gram model. Let  $x$  be the number of instances of *british* required. Then using the discounting scheme described, we have

$$\begin{aligned} P(\text{damnbritish}) &= 1/(x+6) \\ P(\text{damn}) &= 2/(x+6) \quad \text{and} \\ P(\text{british}) &= x/(x+6) \end{aligned}$$

We seek an  $x$  for which  $P(\text{damn})P(\text{british}) > P(\text{damnbritish})$ . Thus, we get

$$2x/(x+6)^2 > 1/(x+6) \Rightarrow x > 6$$

The actual scores for MBDP-1 when presented with *damnbritish* for a second time are  $-\log P(\text{damnbritish}) = 2.8$  and  $-\log P(\text{D\&m}) - \log P(\text{brItIS}) = 1.8 + 0.9 = 2.7$ . For the 1-gram model,  $-\log P(\text{damnbritish}) = 2.6$  and  $-\log P(\text{D\&m}) - \log P(\text{brItIS}) = 1.9 + 0.6 = 2.5$ . Note, however, that skepticism in this regard is not always a bad attribute. It is desirable for the model to be skeptical in inferring new words because a badly inferred word will adversely influence future segmentation accuracy.

## 6. Summary

In summary, we have presented a formal model of word discovery in continuous speech. The main advantages of this model over that of Brent (1999) are: First, the present model has been developed entirely by direct application of standard techniques and procedures in speech processing. Second, it makes few assumptions about the nature of the domain and remains conservative as far as possible in its development. Finally, the model can be easily extended to incorporate more historical detail. This is clearly evidenced by the extension of the unigram model to handle bigrams and trigrams. Empirical results from experiments suggest that the algorithm performs competitively with alternative unsupervised algorithms proposed for inferring words from continuous speech. We have also carried out and reported results from experiments to determine whether particular ways of estimating phoneme (or letter) probabilities may be more suitable than others.

Although the algorithm is originally presented as an unsupervised learner, we have shown the effect that training data has on its performance. It appears that the 3-gram model is the most responsive to training information with regard to segmentation precision, obviously by virtue of the fact that it *keeps* more knowledge from the utterances presented. Indeed, we see that a *fully trained* 3-gram model performs with 100% accuracy on the test set. Admittedly, the test set in this case was identical to the training set, but we should keep in mind that we were keeping only limited history—namely 3-grams—and a significant number of utterances in the input corpus (4023 utterances) were four words or more in length. Thus, it is not completely insignificant that the algorithm was able to perform this well.

## 7. Future work

We are presently working on the incorporation into the model of more complex phoneme distributions, such as the biphone and triphone distributions. Some preliminary results we have obtained in this regard appear to be encouraging.

With regard to estimation of word probabilities, a fruitful avenue we are exploring involves modification of the model to address the sparse data problem using interpolation such that

$$P(w_i | w_{i-2}, w_{i-1}) = \lambda_3 f(w_i | w_{i-2}, w_{i-1}) + \lambda_2 f(w_i | w_{i-1}) + \lambda_1 f(w_i)$$

where the positive coefficients satisfy  $\lambda_1 + \lambda_2 + \lambda_3 = 1$  and can be derived so as to maximize  $P(\mathbf{W})$ .

Taking the lead from Brent (1999), attempts to model more complex distributions for unigrams such as those based on template grammars, as well as the systematic incorporation of prosodic, stress, and phonotactic constraint information into the model, are both subjects of current interest. Unpublished results already obtained suggest that biasing the segmentation such that every word must have at least one vowel in it dramatically increases segmentation precision from 67.7% to 81.8%, and imposing a constraint that words can begin or end only with permitted clusters of consonants increases precision to 80.65%. We are planning experiments to investigate models in which these properties can be learned in the same way as  $n$ -grams.

## Appendix: Inventory of Phonemes

The following tables list the ASCII representations of the phonemes used to transcribe the corpus into a form suitable for processing by the algorithms.

Consonants		Vowels		Vowel + r	
ASCII	Example	ASCII	Example	ASCII	Example
p	pan	I	bit	3	bird
b	ban	E	bet	R	butter
m	man	&	bat	#	arm
t	tan	A	but	%	horn
d	dam	a	hot	*	air
n	nap	O	law	(	ear
k	can	U	put	)	lure
g	go	6	her		
N	sing	i	beet		
f	fan	e	bait		
v	van	u	boot		
T	thin	o	boat		
D	than	9	buy		
s	sand	Q	bout		
z	zap	7	boy		
S	ship				
Z	pleasure				
h	hat				
c	chip				
G	gel				
l	lap				
r	rap				
y	yet				
w	wall				
W	when				
L	bottle				
M	rhythm				
~	button				

### Acknowledgments

The author thanks Michael Brent for stimulating his interest in the area. Thanks are also due to Koryn Grant for verifying the results presented here. Eleanor Olds Batchelder and Andreas Stolcke gave many constructive comments and useful pointers during the preparation of a revised version of this paper. The "Damn British" anecdote is due to Robert Linggard. Judith Lee at SRI edited the manuscript to remove many typographical and stylistic errors. In addition, the author is very grateful to anonymous reviewers of an initial version, particularly "B" and "D", whose encouragement and constructive criticism helped significantly.

### References

- Ando, R. K. and Lillian Lee. 1999. Unsupervised statistical segmentation of Japanese Kanji strings. Technical Report TR99-1756, Cornell University, Ithaca, NY.
- Batchelder, Eleanor Olds. 1997. *Computational evidence for the use of frequency information in discovery of the infant's first lexicon*. Unpublished Ph.D. dissertation, City University of New York, New York, NY.
- Bernstein-Ratner, N. 1987. The phonology of parent child speech. In K. Nelson and A. van Kleeck, editors, *Children's Language*, volume 6. Erlbaum, Hillsdale, NJ.
- Brent, Michael R. 1999. An efficient probabilistically sound algorithm for segmentation and word discovery. *Machine Learning*, 34:71-105.
- Brent, Michael R. and Timothy A. Cartwright. 1996. Distributional regularity and phonotactics are useful for segmentation. *Cognition*, 61:93-125.
- Christiansen, M. H., J. Allen, and M. Seidenberg. 1998. Learning to segment speech using multiple cues: A connectionist model. *Language and Cognitive Processes*, 13:221-268.
- Cutler, Anne and D. M. Carter. 1987. Predominance of strong initial syllables in the English vocabulary. *Computer Speech and Language*, 2:133-142.
- de Marcken, C. 1995. Unsupervised acquisition of a lexicon from continuous speech. Technical Report AI Memo No. 1558, Massachusetts Institute of Technology, Cambridge, MA.

- Elman, J. L. 1990. Finding structure in time. *Cognitive Science*, 14:179–211.
- Hausser, David. 1988. Quantifying inductive bias: AI learning algorithms and Valiant's learning framework. *Artificial Intelligence*, 36:177–221.
- Jelinek, F. 1997. *Statistical Methods for Speech Recognition*. MIT Press, Cambridge, MA.
- Jusczyk, Peter W. 1997. *The Discovery of Spoken Language*. MIT Press, Cambridge, MA.
- Jusczyk, Peter W., Anne Cutler, and N. Redanz. 1993. Preference for predominant stress patterns of English words. *Child Development*, 64:675–687.
- Jusczyk, Peter W. and E. A. Hohne. 1997. Infants' memory for spoken words. *Science*, 27:1984–1986.
- Katz, Slava M. 1987. Estimation of probabilities from sparse data for the language model component of a speech recognizer. *IEEE Transactions on Acoustics, Speech and Signal Processing*, ASSP-35(3):400–401.
- MacWhinney, Brian and C. Snow. 1985. The child language data exchange system. *Journal of Child Language*, 12:271–296.
- Mattys, Sven L. and Peter W. Jusczyk. 1999. Phonotactic and prosodic effects on word segmentation in infants. *Cognitive Psychology*, 38:465–494.
- Olivier, D. C. 1968. *Stochastic grammars and language acquisition mechanisms*. Unpublished Ph.D. dissertation, Harvard University, Cambridge, MA.
- Saffran, Jennifer R., E. L. Newport, and R. N. Aslin. 1996. Word segmentation: The role of distributional cues. *Journal of Memory and Language*, 35:606–621.
- Viterbi, Andrew J. 1967. Error bounds for convolutional codes and an asymptotically optimal decoding algorithm. *IEEE Transactions on Information Theory*, IT-13:260–269.
- Witten, Ian H. and Timothy C. Bell. 1991. The zero-frequency problem: Estimating the probabilities of novel events in adaptive text compression. *IEEE Transactions on Information Theory*, 37(4):1085–1091.
- Zimin, W. and G. Tseng. 1993. Chinese text segmentation for text retrieval problems and achievements. *JASIS*, 44(9):532–542.