

# A computational analysis of uniqueness points in auditory word recognition

PAUL A. LUCE

Indiana University, Bloomington, Indiana

To determine the extent to which words in isolation may be recognized prior to their offsets, *uniqueness points* or *optimal discrimination points* were computed for all words in a 20,000-word computerized lexicon with more than two phonemes. The results of this analysis revealed that the frequency-weighted probability of a word's diverging from all other words in the lexicon prior to the last phoneme was only .39. This finding suggests that an optimally efficient strategy of word recognition may be severely limited in scope due to structural properties of the mental lexicon.

According to Marslen-Wilson's (1984; Marslen-Wilson & Tyler, 1980; Marslen-Wilson & Welsh, 1978) cohort theory of auditory word recognition, the initial acoustic-phonetic information of a word presented to the listener activates a "cohort" of lexical candidates that share word-initial acoustic-phonetic information. Lexical candidates that are incompatible with ensuing top-down and bottom-up information are successively eliminated from the cohort until only one word remains, at which time that word is recognized. Because a word may be recognized when the initial acoustic-phonetic information of the word is compatible with no other words in the cohort, word recognition within the framework of cohort theory is said to be "optimally efficient" in the sense that a listener may recognize a word prior to hearing the entire word. Thus, a crucial concept in the cohort theory of word recognition is that of the *uniqueness point* or *optimal discrimination point*. The uniqueness point is that point, measured from the beginning of the word, at which a word diverges from all other words in the lexicon. For isolated words, the uniqueness point defines the earliest point, theoretically, at which a word can be recognized, although a word may be recognized prior to its uniqueness point given sufficiently constraining contextual information.

The role of the uniqueness point in the recognition of stimuli presented in isolation has been demonstrated for nonwords in an auditory lexical-decision task (Marslen-Wilson, 1984) and for specially selected words in a gating task (Luce, Pisoni, & Manous, 1984; Tyler & Wesels, 1983). Thus, the concept of a uniqueness point is an empirically justifiable notion. On intuitive grounds, however, it appears that an optimally efficient strategy of recognizing words at the earliest point at which they diverge from all other words is somewhat limited. For

example, words of shorter length (e.g., *car*) may overlap completely with the initial portions of words of longer length (e.g., *card* or *carbohydrate*). Thus, in a large number of cases, shorter words presented without any contextual support may have no uniqueness point whatsoever. Furthermore, given the well-known finding that frequency of usage varies inversely with word length, it is possible that for words most frequently used in the language, an optimally efficient strategy is of little or no importance in the absence of contextual information permitting early word recognition. In short, the possibility exists that optimally efficient strategies in auditory word recognition are severely constrained by structural properties of the lexicon alone.

## METHOD AND RESULTS

To evaluate the possibility that an optimally efficient strategy of word recognition may not be applicable to a large number of words in the listener's lexicon, uniqueness points were computed from phonetic transcriptions for approximately 20,000 words contained in an on-line data base. (These results are part of a larger study concerned with the structural properties of the lexicon. See Luce, 1986.) The results of these analyses, broken down by word length, are shown in Table 1.

The percentages of words of a given length that diverge before the end of the word, at the end of the word, and after the end of the word are shown in Table 1 for each word length. (Words one phoneme long were not included in these analyses.) Words that are said to diverge after the last phoneme are words that are embedded entirely in the initial portions of longer words (e.g., *car* is embedded entirely in *card*, *carbohydrate*, etc.). Also shown in Table 1 are the number of words of each word length and the mean frequency of these words, derived from the Kučera and Francis (1967) norms. For the purposes of the present analysis, words occurring in the data base that were not listed in the Kučera and Francis frequency norms were assumed to have a value of one.

This work was supported in part by NIH research grant NS 12179. I would like to thank David B. Pisoni for his comments and suggestions and Howard Nusbaum for helpful discussions. Requests for reprints should be sent to Paul A. Luce, Speech Research Laboratory, Department of Psychology, Indiana University, Bloomington, IN 47405.

**Table 1**  
**Percentages of Words (Lengths 2 Through 17) That Diverge Before, At, and After the Last Phoneme**

Word Length	Number of Words	Mean Frequency	Percentage of Words of a Given Length That Diverge		
			Before Last Phoneme	At Last Phoneme	After Last Phoneme
2	263	1,096.28	0.00	5.70	94.30
3	1,839	126.28	0.76	25.07	74.17
4	3,025	35.32	12.99	51.37	35.64
5	3,172	16.34	49.65	33.70	16.36
6	3,063	11.38	73.23	18.32	8.56
7	2,735	9.05	83.40	11.01	5.59
8	2,210	7.00	85.48	9.77	4.75
9	1,540	7.12	90.32	7.21	2.47
10	1,023	5.46	94.13	4.20	1.66
11	527	4.50	96.77	3.04	0.19
12	268	3.86	97.39	1.87	0.75
13	106	2.58	100.00	0.00	0.00
14	37	4.57	100.00	0.00	0.00
15	13	1.46	100.00	0.00	0.00
16	6	1.00	100.00	0.00	0.00
17	3	1.00	100.00	0.00	0.00

*Note*—The number of words and their mean frequencies are also shown for each word length

As shown in Table 1, the percentages of words of a given length that diverge prior to the last phoneme are relatively low for words two, three, and four phonemes long. Even for words five phonemes long, only approximately half of the words have uniqueness points prior to the last phoneme. Note also that the mean frequencies at each of these word lengths are quite high, indicating that among words most frequently used, very few have uniqueness points prior to the end of the word. Also of interest in this table are the percentages of words of a given length that overlap entirely with the initial portions of longer words. Again, for words two, three, and four phonemes long, very high percentages of these words are not discriminable from longer words in the lexicon prior to the end of the word. In isolation, of course, these words will be discriminable from longer words at their offsets (i.e., when a “null” phoneme is encountered). However, these results suggest that when shorter words are embedded in sentence contexts where word boundaries may not be clearly marked, discrimination of shorter words from longer words in the lexicon may be accomplished only when a portion of the initial acoustic-phonetic information in the following word has been analyzed. Such a result suggests that words in a cohort may have to be activated retroactively when their initial portions combine with a preceding word to form a possible word. However, this suggestion is, at present, difficult to countenance within cohort theory.

The results of the uniqueness-point analysis are shown graphically in Figure 1. In this figure, cumulative probability of occurrence in the lexicon is plotted against word length in phonemes. As shown in this graph, the probability of a word that diverges prior to the last phoneme is .59, whereas the probabilities of a word that diverges at the last phoneme and after the last phoneme are .22 and .19, respectively. Note also that for words diverging

at or after the last phoneme, the functions reach asymptote much earlier than those for words diverging before the last phoneme. This shows, not surprisingly, that words of shorter length tend not to have uniqueness points prior to the end of the word.

Figure 2 shows the cumulative probabilities of occurrence in the lexicon weighted by log frequency. When weighted by frequency, the probability of a word that diverges before the last phoneme drops to .39, whereas the probabilities of a word that diverges at or after the last phoneme increase to .23 and .38, respectively. These results corroborate the earlier observation that words of higher frequency, being shorter, tend not to have uniqueness points prior to the end of the word. In addition, the large increase in cumulative probability for words diverging after the end of the word indicates that a large proportion of high-frequency words overlap in their entirety with the initial portions of longer words in the lexicon.

Mean discrimination points as a function of word length are plotted in Figure 3. The dashed line represents the actual mean discrimination points. The solid line, plotted for purposes of comparison, represents hypothetical discrimination points occurring at the ends of each word at each word length. The divergence of these two lines displays the degree to which words at each length actually diverge prior to the end of the word. As shown in Figures 1 and 2 and in Table 1, mean discrimination points tend to diverge from the ends of the words only for words of longer length. However, it is of interest to note that the line for the mean isolation points tends to level off at longer word lengths, indicating that the longer the word, the proportionally earlier the isolation point. This suggests that the mean overlap of word-initial acoustic-phonetic information tends to be limited to less than eight or nine phonemes, regardless of the length of the word.

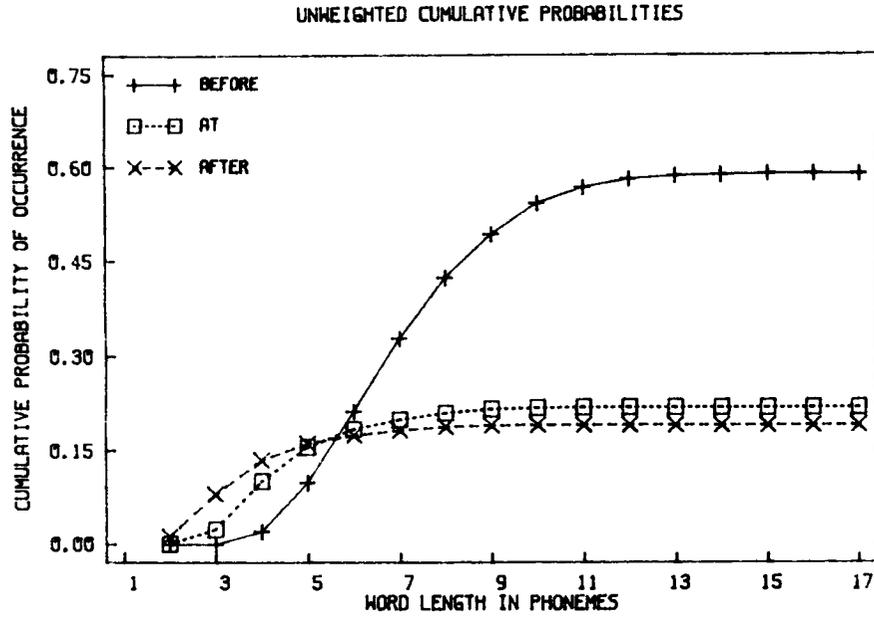


Figure 1. Unweighted cumulative probabilities of occurrence for words having uniqueness points prior to the end of the word, at the end of the word, and after the end of the word. Solid line with pluses represents words that diverge from all other words in the lexicon prior to the last phoneme. Dotted line with squares represents words that diverge at the last phoneme. Dashed line with Xs represents words that diverge after the last phoneme.

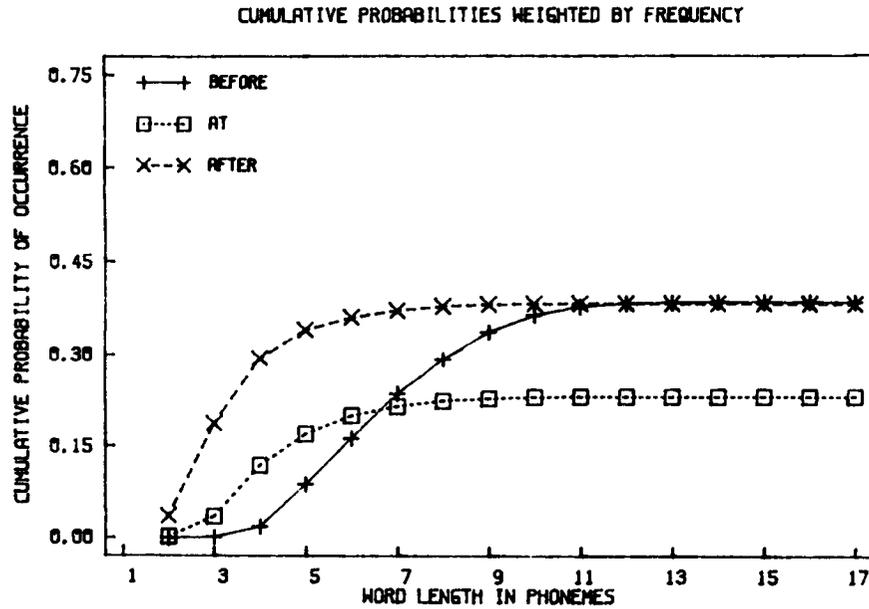


Figure 2. Cumulative probabilities of occurrence weighted by log frequency for words having uniqueness points prior to the end of the word, at the end of the word, and after the end of the word.

DISCUSSION

The results of the present statistical analyses of the uniqueness points of words in a large, computer-based lexicon suggest that an optimally efficient strategy of word recognition, as proposed by Marslen-Wilson (1984), is severely constrained by structural properties of the lex-

con. In particular, when frequency of usage is taken into account, the probability of a word diverging from all other words in the lexicon prior to the end of the word itself is only .39. In addition, it was shown that of shorter words, which tend in the long run to be high in frequency, a high percentage not only fail to diverge prior to the end of the word, but overlap in their entirety with the initial

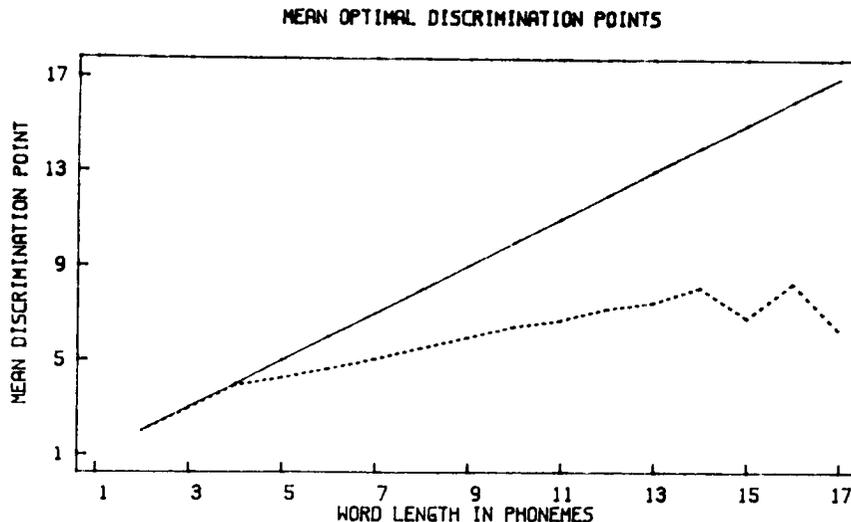


Figure 3. Mean uniqueness points as a function of word length (dashed line). The solid line, plotted for purposes of comparison, represents hypothetical uniqueness points at the ends of words of each length.

portions of longer words. This finding suggests that in fluent speech, many of the most frequent words in the language will not be recognized until some portion of the word-initial acoustic-phonetic information of the following word is processed, given, of course, minimal word-boundary cues and contextual information relevant to the recognition of the target word.

Although these results do not dispute Marslen-Wilson's (1984) notion of an optimally efficient processing strategy in auditory word recognition, they do call into question the scope of such a strategy and the role it plays in the real-time processing of fluent continuous speech. In a large majority of cases, it will be virtually impossible to recognize a word without having heard the entire word (or indeed, a portion of the *following* word, when produced in fluent speech). Thus, the utility of a uniqueness point may be minimal at best or may be restricted to the processing of longer, low-frequency words. These results imply that the notion of a uniqueness point as embodied in the cohort theory of auditory word recognition may not be as important a concept as has thus far been thought. Indeed, the present results suggest that problems of more crucial interest to researchers in auditory word recognition center around the issues of segmentation of words in the acoustic stream and recovery from garden-path (i.e., misleading) analyses arising from the overlap of shorter words with the initial portions of longer words.

Finally, these results demonstrate the utility of using large computerized data bases in testing the plausibility and applicability of concepts in theories of auditory word recognition.

#### REFERENCES

- KUČERA, F., & FRANCIS, W. (1967). *Computational analysis of present day American English*. Providence, RI: Brown University Press.
- LUCE, P. A. (1986). *Structural distinctions between high and low frequency words in auditory word recognition*. Unpublished doctoral dissertation, Indiana University, Bloomington.
- LUCE, P. A., PISONI, D. B., & MANOUS, L. M. (1984). Isolation points and frequency effects in the gating paradigm: Predictions from an on-line data-base. In *Research on speech perception* (Progress Report No. 10, pp. 303-310). Bloomington: Indiana University.
- MARSLÉN-WILSON, W. (1984). Function and process in spoken word recognition: A tutorial review. In H. Bouma & D. Bouwhuis (Eds.), *Attention and performance: Vol. 10. Control of language processes* (pp. 125-148). Hillsdale, NJ: Erlbaum.
- MARSLÉN-WILSON, W., & TYLER, L. K. (1980). The temporal structure of spoken language understanding. *Cognition*, **8**, 1-71.
- MARSLÉN-WILSON, W., & WELSH, A. (1978). Processing interactions and lexical access during word recognition in continuous speech. *Cognitive Psychology*, **10**, 29-63.
- TYLER, L. K., & WESSELS, J. (1983). Quantifying contextual contributions to word-recognition processes. *Perception & Psychophysics*, **34**, 409-420.

(Manuscript received August 15, 1985;  
revision accepted for publication February 21, 1986.)