# Detecting Phonemes in Fluent Speech

GARY S. DELL AND JEAN E. NEWMAN

*University of Toronto*

It has been shown that the latency to detect a word-initial phoneme in a sentence is sensitive to a word predictability variable (Morton & Long, *Journal of Verbal Learning and Verbal Behavior,* 1976, 15, 43–52) and a phonetic similarity variable (Newman & Dell, *Journal of Verbal Learning and Verbal Behavior,* 1978, 17, 359–374). Two experiments are reported demonstrating that these variables interact in a manner suggesting the action of two processes in phoneme detection—one proceeding top-down from the word level and the other bottom-up from an analysis of the acoustic signal. The implications of these results for a theory of word recognition are discussed.

According to modern linguistic theory, the phoneme, or phonetic segment, is one of the major building blocks of language— phonemes are the constituents of words just as words are the constituents of sentences. The role of the phoneme in speech perception, however, is not so clear-cut. While some theorists hold that the construction of a phonetic or phonological representation is a basic process in the perception of fluent speech, others (e.g., Warren, 1976; Klatt, 1980) claim that the phoneme is largely an unnecessary construct. It is certainly true that the naive listener is usually unaware of the phonemic nature of speech. For him or her, speech contains words and ideas, not phonemes. Because of this fact, any attempt to study the perception of speech sounds must of necessity involve modification of the natural listening situation. One approach has focused on the identification of isolated syllables. In this case the simplification of the stimulus directs the

listener's awareness to the sounds of speech, and judgments that are made reveal the relevant acoustic parameters.

A second approach, and the one that we are concerned with, is to study the detection of phonemes in fluent speech. Again the listener is forced to attend to the phonemic nature of speech but, unlike the previous approach, the task is not to identify the stimulus, but to detect the occurrence of a single, previously specified, target phoneme. This is the paradigm known as phoneme monitoring. It was developed by Foss and his associates (Foss, 1969; Foss & Lynch, 1969) as a technique for measuring on-line processing difficulty during sentence comprehension. The subject listens to a stimulus sentence and is told to press a response key as quickly as possible on hearing a word beginning with the target phoneme, and the measure of interest is the latency of the response. We shall refer to the word containing the target as the *target word,* the preceding word as the *critical word,* and its initial phoneme as the *critical phoneme.*

Although this paradigm has been used primarily as a measuring device (e.g., Foss & Jenkins, 1973; Swinney & Hakes, 1976) our concern is with the more basic question of how the target phoneme is detected; must the subject recognize the target word before the phoneme can be identified (here called *top-down* detection) or can the target be detected prior to target word recognition

(*bottom-up* detection)? This issue has implications beyond those specific to the phoneme monitoring paradigm. If phoneme detection is a top-down process, it suggests that the construction of a segmentally-organized representation is not a basic process in the recognition of words in fluent speech. Phonemes are merely an afterthought; they are synthesized from an already-recognized word, perhaps because the phoneme monitoring task demands it. On the other hand, the finding that a phonemic representation leading to a detection response can be constructed prior to lexical access, would support the hypothesis that such a representation forms a basis for word recognition. Also, if both top-down and bottom-up processes are involved, it would be useful to ascertain how these processes co-occur.

At present, most of the evidence supports the hypothesis that phonemes are detected in a purely top-down fashion. Foss and Swinney (1973) found that reaction times to phoneme targets are substantially longer than those to word targets, and they argued that the detection of a phoneme follows the recognition of the word containing it. Similar results were reported by Warren (1971). Although these results certainly support the top-down hypothesis, objections have been raised. It has been pointed out (Foss & Swinney; 1973; Morton & Long, 1976) that the finding of longer detection times for phonemes does not eliminate the possibility that phonemes are perceived before words. Phonemes may simply take longer to get to consciousness or to an overt identification mechanism.

In a more compelling demonstration that phoneme detection requires lexical access, Morton and Long (1976) had subjects monitor for word-initial phonemes in sentences in which the target word was very predictable (in their terms, had a high transitional probability), such as in (1), and in corresponding sentences, such as (2), where the target word was less predictable (near zero transitional probability). In both cases the target is /b/.

(1) A sparrow sat on the *branch* whistling a few shrill notes to welcome the dawn.
(2) A sparrow sat on the *bed* whistling a few shrill notes to welcome the dawn.

On average, targets in predictable words were detected 70-milliseconds faster than those in less predictable words. Morton and Long suggested that the finding may be entirely due to differences in the time taken to recognize the target words, and thus the phoneme identification process must include the identification of the word containing it.

Evidence for the bottom-up detection of phonemes does exist, however. Foss and Blank (1980) found that the frequency of the target word did not influence the time to detect word-initial phonemes in sentences. If detection were mediated by target word recognition, one would expect variables that influence lexical access, such as frequency, to have at least some effect. More strikingly, Foss and Blank reported that if the target was contained in a nonword, detection latencies were just as short as if it was in a word. Again, it is difficult to reconcile this result with a model that requires recognition of a word before its phonemes become available.

Further support for bottom-up detection comes from studying false alarms (responses to phonemes other than the target) in the phoneme monitoring task. Newman and Dell (1978) found that for the target phoneme initial-/b/, the vast majority of false alarms were to initial-/p/, a phoneme which differs from /b/ only in the voicing feature. Other false alarms were to phonemes that were similar to /b/ but to a lesser extent. These results can be interpreted as support for bottom-up detection for two reasons. First, the words that were incorrectly responded to would not have been words if their initial phonemes had been the target /b/. For example, words such as *policeman, primitive,* and *pilot* attracted responses, but it is difficult to argue that subjects were responding on the basis

of recognizing a /b/-initial word, because the substitution of /b/ for /p/ in these cases results in nonwords. Secondly, the fact that similarity to /b/ was the critical factor in determining false alarms, points to a model of detection where phoneme identity is ascertained by a simple matching of features of the stimulus to features of the target.

So far we have reviewed evidence supporting both top-down and bottom-up detection. The common-sense conclusion would be that both processes are going on because they are certainly not mutually exclusive. This claim has in fact been made (Newman & Dell, 1978; Foss, Harwood, & Blank, 1980; Cutler & Norris, 1979) and has been more fully specified by Foss and Blank (1980) under the name of the dual code hypothesis. Our model of speech processing as it applies to phoneme detection extends this earlier work, and will be described in the next section.

## A Model of Phoneme Detection

*Bottom-up detection.* According to Foss and Blank (1980) bottom-up detection is mediated by the construction of a *phonetic code,* which is a linguistic representation of the speech stream that arises from acoustically represented input, and serves as the basis for lexical access. This representation is composed of phonetic segments or feature "bundles" and preserves much acoustic information that is nondistinctive (e.g., aspiration after initial voiceless stops in English). In addition, it is presumed to be a fragile, rapidly decaying, representation. Detection via the phonetic code involves two processes: the determination of which segment is word initial, and the determination of whether that segment is the target. The first of these processes typically involves the recognition of the previous word. Thus by determining that a certain string of segments constitutes a word, one can label the next segment as word-initial (Foss & Blank, 1980; Cole & Jakimik, 1979). A prediction from this analysis is that factors influencing the time to recognize the word preceding the target (which we call the

critical word), will affect target detection time. This prediction concerning the critical word has been confirmed with respect to its frequency (Foss, 1969), predictability (Blank & Foss, 1978), and lexical status (Foss & Blank, 1980). After a segment has been labeled as word initial, it must be ascertained whether or not that segment is the target. The results of Newman and Dell (1978) discussed above, argue for some sort of feature comparison process in this decision. An initial segment represented as a feature bundle is compared to a featural description of the target, and a response is made when all the features are found to match.

*Top-down detection.* Detection by this route involves the construction of what Foss and Blank (1980) call the *phonological code.* This code is an ordered set of phonemes comprising a word, that becomes available after that word has been recognized, or at least after the word's lexical entry has been activated. The phonological code is both more durable than the phonetic code, and more abstract, in that it does not represent allophonic variation. Detection of a target initial phoneme, once the phonological code of the target word is available, is much simpler than detection through the phonetic code. The subject just reads off the initial phoneme and responds if it is the target. There is no need to determine which is the initial phoneme, because this information is directly represented, and furthermore, we are assuming that the comparison process involving the initial phoneme and the target is wholistic rather than feature-by-feature.

Evidence for the two different comparison processes associated with each of the two detection routes comes from same – different judgment tasks with phoneme pairs. In an experiment by Cole and Scott (1972), subjects were simultaneously presented with two nonsense syllables—either the same or differing in their initial consonants (e.g., /ga/, /ba/). Under such conditions we are forced to assume that the phonetic codes of the two stimuli would be

compared. Cole and Scott found that "different" response times increased with the number of shared features between the initial consonants, thus supporting our contention that comparisions using the phonetic code involve features. A related experiment by Chananie and Tikofsky (1969) shows how the decision process appears to change if the phonological code is involved. Like Cole and Scott, Chananie and Tikofsky used syllable pairs that were identical or differed in their initial consonants, but the stimuli were words (e.g., *pan, tan*) and were presented successively. The results were that "different" response times were not affected by the number of features common to the initial consonants. Because the stimuli were lexical items and were presented one at a time, it is reasonable to assume that lexical access was involved, thus allowing the phonological code to be the medium of comparison. So, the lack of a similarity effect suggests that phonological phonemes are compared as units rather than features. Thus, there is independent evidence for our assumption that the use of the phonetic code involves

featural comparisons whereas the use of the phonological code does not.

At this point, we shall turn to the question of how the top-down and bottom-up detection processes co-occur. Our proposal, the parallel access hypothesis, is that the detection routes operate in parallel, as diagrammed in Figure 1, with the first route to finish determining response time—in other words a "horse-race" model. Such a model can be tested by simultaneously manipulating the rate of processing (the finishing time) of each component of the putative parallel process. Our experiments were designed to provide such a test.

*Testing the Model*

The experiments that we shall present manipulated the speed of the top-down detection route by varying the predictability of the target word, exactly as was done by Morton and Long (1976). This manipulation is seen as affecting the word recognition stage of the process, hence it is specific to the top-down detection route.

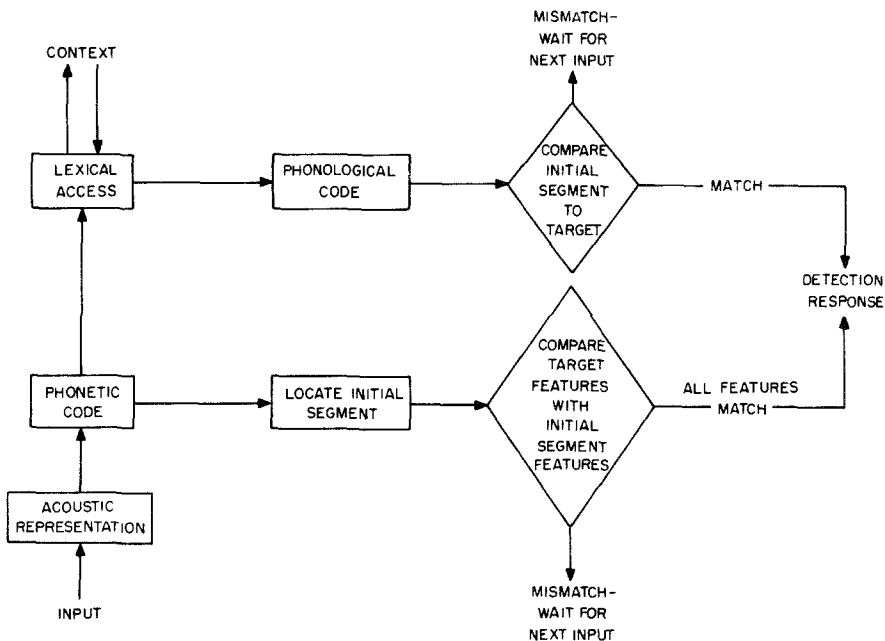In addition to the top-down manipulation, we attempted to independently ma-



FIG. 1. The parallel access model of phoneme detection.

nipulate the speed of the bottom-up detection route. This was accomplished by varying the similarity of the critical phoneme (the initial phoneme of the critical word) to the target. This variable was found by Newman and Dell (1978) to affect detection times, such that responses to the target initial /b/ in a sentence like (3) below would be on the average 90-milliseconds longer than those to the /b/ in (4), because the /p/ in *platter* is very similar to /b/, while the /s/ in *saucer* is not at all similar.

(3) The large *platter* broke on the floor.
(4) The large *saucer* broke on the floor.

(Unlike Newman and Dell who named this the phonological similarity effect, we shall refer to it as the *phonetic* similarity effect, because of the distinction that was drawn above between the two codes and their properties.)

It is our contention that similar critical phonemes affect detection time by delaying the bottom-up detection process. More specifically, we will attempt to show that the delay occurs in that component of the bottom-up route which determines whether or not the phoneme in question is the target—in other words, the feature-matching and decision component. In our model, an initial phoneme similar to the target will be treated exactly as a dissimilar one at every stage of analysis, except for the feature comparison stage of the bottom-up process. Here the similarity will be recognized. (Note that similarity to the target would not be recognized by the top-down decision component because of our assumption that this decision is wholistic rather than feature-by-feature.) The effect of similarity at the feature comparison stage is to create a difficult decision. If we assume that a decision to respond is made if every feature of the stimulus matches a feature of the target, and that the comparison is terminated (with no response) if a mismatch is found, then the correct decision for a similar phoneme will take extra 'time or effort. Evidence for the difficulty of

processing similar phonemes is provided by the finding mentioned earlier, that subjects often respond (incorrectly) to similar phonemes differing from the target by a single feature.

We contend that the effect of the difficulty associated with similar critical phonemes is to retard the bottom-up processing of the *target* phoneme, thus resulting in the phonetic similarity effect. During the processing of a critical phoneme similar to the target, the bottom-up decision will be delayed, for the reasons outlined above, except in cases where the result is a false alarm. The bottom-up analysis of the subsequently occurring target phoneme is then delayed by the amount of the original difficulty in rejecting the critical phoneme. When the critical phoneme shares few features with the prespecified target, it can be rejected quickly and the detection of the actual target can be made on schedule. Note that the phonetic similarity effect is due, in this view, to a delay in processing the target phoneme as a result of the feature comparison process needed to reject the critical phoneme. The same feature comparison process is used to detect the target—the resulting increase in reaction time is not due to "adaptation" or fatigued feature detectors, but rather to the sequential nature of the decision processes in a limited capacity system. It would be hard to argue for an adaptation explanation of these results, as such effects typically take many trials to develop (e.g., Eimas & Corbit, 1973), rather than the single trial mechanism that would have to be invoked for this task.

Therefore in our framework, the phonetic similarity variable acts solely on the bottom-up detection route, and target word predictability affects the word recognition stage of the top-down detection route. By manipulating both these variables, the way in which top-down and bottom-up processes co-occur can be determined.

Subjects in our experiments were directed to monitor for words beginning with

the target phoneme /b/ in sentences where the critical phoneme was either similar to the target (/p/), or dissimilar (/s/), and where the target word was predictable from preceding context (High), or not (Low). Sentences (5)–(8) illustrate the four conditions:

(5) The surfers drove to a *private beach* to try out the waves. (/p/-High)

(6) The surfers drove to a *private bay* to try out the waves. (/p/-Low)

(7) The surfers drove to a *secret beach* to try out the waves. (/s/-High)

(8) The surfers drove to a *secret bay* to try out the waves. (/s/-Low)

Given that phonetic similarity and target predictability act as we have assumed, the predictions derived from the parallel access model are straightforward. When the target bearing word is unpredictable (Low conditions), the time to identify the target by means of lexical access will be slow, and because of the hypothesized parallel process, effects of the phonetic similarity variable will show up in detection times. More specifically, the /p/-Low condition should result in substantially longer reaction times than the /s/-Low condition. When lexical access is more rapid (High predictability conditions), the probability of top-down detection is greater and the interfering bottom-up effects of the critical phoneme should be washed out or at least attenuated. Similarly, effects of target predictability should be most apparent when the critical phoneme is /p/. To sum up, the parallel access model predicts an overall advantage of predictable over unpredictable targets, faster reaction times with /s/ as the critical phoneme than with /p/, and more importantly, an interaction such that the effects of each variable are attenuated when combined with the "faster" level of the other variable.

The pure top-down hypothesis proposed by Morton and Long (1976) predicts that the two variables should not interact. According to Morton and Long, phonemes are detected by a two-stage process; the first stage recognizes the target word, and the second identifies the phonemes of that word. Phonetic similarity can be assumed to act on the second stage, and target word predictability on the first. Because the stages are serially ordered, the variables should have only additive effects. Thus, according to this hypothesis, the main effect of target predictability should be as strong with /s/-initial critical words as with /p/-initial critical words.

## EXPERIMENT I

### Method

*Pretest.* The sentences used in this experiment were designed to permit the simultaneous manipulation of phonetic similarity and target word predictability. In order to minimize possible sources of variation, a within-sentence design was chosen. Two levels of phonetic similarity were crossed with two levels of target word predictability, resulting in four versions of a basic sentence frame. The critical word began with either a /p/ or an /s/, and the target word was either predictable (High) or unpredictable (Low), from the preceding context. A further constraint was that both the /p/- and /s/-initial critical words did not differentially influence the predictability of the target word. Additionally, in all four conditions the target words always began with the same target phoneme (/b/) so they were equally confusable (or not) with the preceding critical phoneme.

While it is possible to determine phonetic similarity mechanically by the use of a feature system, it is not possible to determine target word predictability without recourse to subjects' judgments. Thus 38 sentence frames were constructed and submitted to a pretest. Eighty-four subjects (members of a psycholinguistics class) participated. Each subject was given all the sentence frames, up to and including the critical word. In half the sentence frames the critical word began with /p/ and in the other half the critical word began with /s/. The subjects' task was

to supply the word that they thought would be most likely to come next in the sentence. In order to ensure that subjects did not notice that a large percentage of the words that they would fill in began with /b/, there were also 20 filler sentences to be completed. Thus all subjects performed the fill-in task with all 38 candidate sentence frames, but for a given sentence only *one* of the /p/- or /s/-initial words was used as the last word in the sentence fragment. To illustrate, using the example given above, half of the subjects did the fill-in task with *The surfers drove to a private . . .* and the other half were asked to fill in the next word of *The surfers drove to a secret. . . .* It was expected that *beach,* which was the target word that we had intended for the High conditions using this sentence frame, would be filled in fairly often resulting in a high value for the average transition probability between the preceding sentential context and the target word. Correspondingly, it was expected that *bay* (the designated target for the Low conditions) would be filled in rarely, if at all, resulting in a very low average transition probability.

The results of the pretest were as follows: Thirty-two of the pretested sentences were chosen as stimulus materials for the reaction time study. The mean fill-in rates for intended High transition probability target words were 46.2% following /p/-initial critical words and 45.9% following /s/-initial critical words. The mean fill-in rate for intended Low transition probability target words was 0.8% following both /p/- and /s/-initial critical words. Thus the mean difference in transition probability between the High and Low conditions was very large. This difference was expected to produce a substantial context effect, certainly comparable with that of Morton and Long (1976) who had mean transition probabilities of 36.5 and 0.1% for their High and Low conditions, respectively. Further, it should be stressed that the size of the difference between High and Low transition probability was equal for /p/- and /s/-initial

critical words in our materials. Thus the two experimental variables were manipulated in an orthogonal fashion.

*Materials.* The 32 sentence frames used in this experiment are given in the Appendix. In each sentence the critical and target words constituted either an adjective–noun or noun–verb pair. Because it has been shown that the length of the critical word affects phoneme monitoring latencies (Mehler, Segui, & Carey, 1978; Newman & Dell, 1978; Vipond, Note 2), this variable was controlled. Each critical word was either one or two syllables long, with the mean length of /p/-initial critical words being 1.59 syllables and the mean length of /s/-initial critical words equal to 1.63 syllables. The frequency of critical words is also known to affect phoneme monitoring latencies (Foss, 1969) and was equated for /p/- and /s/-initial critical words. The mean frequency of /p/-initial critical words was 48.7/1.014 million, and /s/-initial critical words had a mean frequency of 42.3/1.014 million (Kučera & Francis, 1967). The lengths and frequencies of the target words were also controlled: for High transition probability target words, mean length = 1.34 syllables and mean frequency = 35.2/1.014 million; for Low transition probability target words, mean length = 1.38 syllables and mean frequency = 33.0/1.014 million. Additionally, since it is known that the stress value of the target phoneme also affects phoneme monitoring (Shields, McHugh, & Martin, 1974), all target phonemes occurred in stressed syllables.

Each sentence frame occurred in four possible forms: (1) with a /p/-initial critical word; (a) followed by a High transition probability target word (/p/-High), (b) followed by a Low transition probability target word (/p/-Low) or (2) with an /s/-initial critical word (a) followed by a High transition probability target word (/s/-High), (b) followed by a Low transition probability target word (/s/-Low). The sentence frames were arbitrarily divided into four groups of eight. Four lists of 32

critical sentences were assembled; on each list the four conditions were equally represented, each sentence group occurring in one of the four conditions. The particular sentence group associated with a condition was rotated across the lists. In addition, there were 11 filler sentences that either contained targets early in the sentence or in adjectives, in order to prevent subjects from generating expectations about target locations in the critical sentences. There were also 17 catch trials (sentences which did not contain a target). The same fillers and catch trials were used on all four lists, and in the same order across lists. The order of the critical sentences was determined randomly but with the constraint that equal numbers of the four conditions occur in each quarter of the list, and that no more than three sentences of the same condition follow one another. The same 'random' order was used for the critical sentences across the four lists.

The four lists were recorded by the female experimenter, using normal intonation, on a Uher Royal de Luxe tape recorder. The tapes were presented to subjects binaurally through Uher W 671 headphones. An ITC digital clock timer was wired to the Dia-pilot output connection of the Uher tape recorder. The timer was started whenever a pulse was picked up by the Dia-pilot tape head. The pulse was recorded on a separate track from the voice. Each pulse was put on the tape manually for each sentence so as to coincide with the onset of the /b/; determined by pulling the tape slowly over the tape head. It is reasonable to assume some slight error in placement but, as the pulses were put on the tape without knowledge of the experimental condition, this error should have been randomly distributed across conditions.

*Subjects and procedure.* Twenty-eight students at the University of Toronto were paid $1.75 for participating in the half-hour experiment. All subjects were native speakers of English. Subjects were assigned sequentially to lists as they were tested: Subject 1 to List 1, Subject 2 to List 2, and so on.

Subjects were told that they were participating in an experiment on sentence comprehension but that the experimenter was interested in monitoring attention as comprehension took place. They were instructed to press a response key as quickly, but also as accurately, as possible whenever they heard a word beginning with /b/. They were told that there would be no more than one initial /b/ per sentence. Subjects were instructed that there would be catch trials and that they would be given a later comprehension test. Subjects read the instructions and then did six practice trials, two of which were catch trials, before starting the trials from the appropriate list. The later comprehension test consisted of subjects being allowed 4 minutes to recall as many of the target words as possible. The test followed immediately upon completion of the monitoring task; subjects were unaware of the nature of the comprehension test in advance. The results of this test were not analyzed for this experiment, other than to ensure that subjects remembered at least six of the target words, thus demonstrating that they had in fact paid attention to the words in the sentences.

*Results and Discussion*

Analysis of the response time data was carried out following a procedure suggested by Forster and Dickinson (1976) for language experiments. In this procedure three analyses of variance are performed: one using condition by subject means ($F_1$ statistic), a second with condition by sentence frame means ($F_2$ statistic), and a third using the entire data matrix with missing observations filled in. This last analysis allows for tests of significance of the error terms of the first two analyses. According to Forster and Dickinson, the appropriate test for a given effect is min $F'$ (or $F'$) only when the corresponding $F_1$ and $F_2$ error terms are significant; otherwise $F_1$ and $F_2$ are appropriate. For each effect we shall

cite $F_1$, $F_2$, and their error terms, and shall indicate whether these error terms are significant by the presence of an asterisk superscript. When both are significant ($\alpha <$ .05), min $F'$ will be given as well.

The results are summarized in Figure 2. Misses include trials in which subjects did not respond, or responded only after 2 seconds. When subjects responded before the occurrence of the target, the trial was not counted as a miss, nor was any subsequent response after the target included.

As predicted, response time was substantially longer when /p/ was the critical phoneme as opposed to when it was /s/; $F_1(1,24) = 26.36$, $p < .01$, $MS_e = 3062.7$; $F_2(1,28) = 26.47$, $p < .01$, $MS_e = 4295.1^*$. Consistent with this finding is the fact that subjects made false alarms to the critical phoneme /p/ on 4.2% of the trials, but made none to /s/. Targets in predictable words were detected faster than those in unpredictable words, but this effect was only marginally significant; $F_1(1,24) = 9.51$, $p <$ .01, $MS_e = 5667.4^*$; $F_2(1,28) = 6.46$, $p <$ .05, $MS_e = 8158.7^*$; min $F'(1,51) = 3.85$, $p < .10$. The interaction predicted by the par-

allel access hypothesis was in the expected direction but not reliably so; $F_1(1,24) = 2.57$, $p > .10$, $MS_e = 2032.6$; $F_2(1,28) = 1.94$, $p > .10$, $MS_e = 3545.0^*$.

One possible explanation of performance in this task is that phonemes can be detected either top-down or bottom-up, but the detection route employed is a result of a strategy decision on the part of individual subjects. For a given stimulus sentence, some subjects would detect phonemes via word recognition, others by listening for them directly, but both strategies would not be employed by the same subject for the same trial. Such a view would predict that the phonetic similarity main effect ($RT_{/p/} - RT_{/s/}$) for a given subject should be negatively correlated with the subject's target word predictability effect ($RT_{\mathrm{Low}} - RT_{\mathrm{High}}$). There was, however, no relationship; the linear correlation was $-.01$.

The lack of a significant interaction between phonetic similarity and target word predictability can be taken as support for the top-down hypothesis of Morton and Long (1976) over the parallel access hypothesis but the data are not compelling. Individual comparisons among condition means suggested a pattern of results consistent with the parallel access hypothesis. The difference in reaction time between predictable and unpredictable targets was reliable when the critical phoneme was /p/, $F_1(1,48) = 12.04$, $p < .01$; $F_2(1,56) = 8.35$, $p < .01$, min $F'(1,103) = 4.93$, $p < .05$; but not when /s/ was the critical phoneme, $F_1(1,48) = 3.32$; $F_2(1,56) = 1.83$; min $F'(1,99) = 1.18$. These statistical tests, do not, of course, demonstrate the presence of an interaction. Rather, they suggest the need for a more powerful experiment—one in which the additivity or nonadditivity of the two variables will be apparent. In the present experiment phonetic similarity consistently affected reaction time, while the target word predictability effect was less consistent. This can be seen in the relative sizes of the error terms for subjects and sentences for the two variables. Our view is



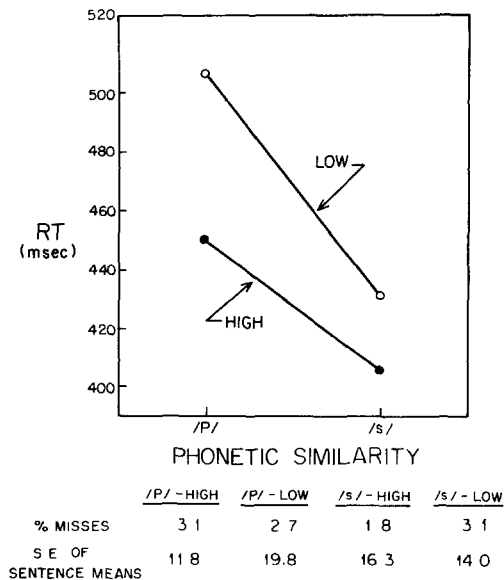| | /P/ -HIGH | /P/ - LOW | /s / - HIGH | /s / - LOW |
|---|---|---|---|---|
| % MISSES | 3 1 | 2 7 | 1 8 | 3 1 |
| S E OF SENTENCE MEANS | 11 8 | 19.8 | 16 3 | 14 0 |

FIG. 2. Detection time as a function of target word predictability and phonetic similarity, Experiment I. (Standard errors expressed in milliseconds.)

that the relative weakness of the target word predictability variable led to the equivocal nature of the present results.

From the point of view of the top-down hypothesis (Morton & Long, 1976), the lack of a significant effect of the target word when the critical phoneme was /s/ was merely due to the overall weakness of the target word predictability variable. Therefore, an experiment achieving a stronger effect of target predictability would provide the critical test.

Such an experiment is also necessitated from the point of view of the parallel access hypothesis. According to the parallel access hypothesis, reaction time on a given trial is determined by the fastest of the top-down and bottom-up detection routes. Thus large variation in the speed of both top-down processing (by varying target predictability) and in the speed of bottom-up processing (by varying the critical phoneme) is required in order to observe the predicted interaction. The interaction may not have been obtained because not all subjects were sensitive to the variation in target word predictability. Such a view would predict that subjects who did show the interaction would tend to be those who had a main effect of target word predictability. That is, there should be a positive correlation between subjects' target predictability effect ($RT_{Low} - RT_{High}$) and the interaction predicted by the parallel access model ($RT_{/p/-Low} + RT_{/s/-High} - RT_{/p/-High} - RT_{/s/-Low}$). This was indeed the case, in the present experiment $r = +.344$, $p < .05$. Although the correlation is not strong, it suggests that the predictions of the parallel access model are borne out when subjects are sensitive to target predictability. Furthermore, an overall interaction should be obtained in an experiment which more effectively induces top-down processing.

So far we have argued that a decision regarding the tenability of the parallel access and pure top-down explanations of phoneme detection must await the outcome of an experiment in which both target word predict-

ability and phonetic similarity consistently affect detection time. Our second experiment attempted to increase the target predictability effect by inducing subjects to pay greater attention to the meaning of the sentences. This was accomplished by directing subjects to paraphrase each sentence after hearing it. We felt that these instructions, more so than those of the first experiment, would cause subjects to make on-line predictions regarding individual words, thus increasing the effectiveness of the target word predictability variable in determining detection times. The paraphrase task has the additional feature that it probably induces a processing strategy more similar to "normal" comprehension, which is more concerned with meaning than with initial-/b/s. But regardless of this fact, the decision times obtained should provide a clear test of the competing hypotheses.

EXPERIMENT II

*Method*

The materials were identical to those used in the first experiment. Twenty-four new subjects (six to a list) from the same population, provided the data. As in the first experiment, subjects were instructed to pay attention to the meaning of each sentence and press the response key as quickly as possible upon hearing a word-initial /b/. In addition, they were instructed to paraphrase each sentence after hearing it. Subjects were told to put some thought into their paraphrases. Rearranging words, such as changing actives to passives, was not acceptable; instead they were told to supply synonyms, or near synonyms, for the content words in the sentences whenever possible. The tape recorder was stopped after each sentence in order to allow subjects ample time to paraphrase. In addition, subjects' paraphrases were recorded on a cassette tape recorder placed beside them. This procedure resulted in subjects being convinced (as shown in postexperimental interviews) that the experiment concerned their ability to para-

phrase, and judging from the quality of their paraphrases, produced the greater attention to meaning that we had hoped for. These instructions had the added benefit that the supposedly secondary phoneme monitoring task really did appear to be less central to the experiment. Finally, unlike the first experiment, there was no recall test.

## Results and Discussion

The data analysis was carried out as described before. Unlike the first experiment, target word predictability and phonetic similarity equally affected response times. Responses following /p/-initial critical words were 52-milliseconds slower on the average than those following /s/-initial critical words; $F_1(1,20) = 6.75, p < .05, MS_e = 9559.0$; $F_2(1,28) = 7.38, p < .05, MS_e = 10,705.4*$; and predictable target words led to faster responses (51 milliseconds) than unpredictable ones, $F_1(1,20) = 5.68, p < .05, MS_e = 10,895.8*$; $F_2(1,28) = 9.14, p < .01, MS_e = 9082.6$. Given that the main effects are reliable and equally strong, the parallel access hypothesis predicts that an overall interaction should be obtained as well. As suggested by the mean detection times (Figure 3), the predicted interaction did result, $F_1(1,20) = 7.02, p < .05, MS_e = 4407.9$; $F_2(1,28) = 4.77, p < .05, MS_e = 8058.3*$. When the critical phoneme was /p/, predictable targets were detected much faster than unpredictable ones; $F_1(1,40) = 11.78, p < .01$; $F_2(1,56) = 13.68, p < .01$; min $F'(1,91) = 6.33, p < .05$. This difference almost completely disappeared when /s/ was the critical phoneme; $F_1 < 1; F_2 < 1$. The false alarms paralleled those of the first experiment, with subjects responding to /p/ on 3.9% of the trials, but never to /s/. It should be noted that the percentages of misses were approximately the same for the four experimental conditions, and somewhat higher than those of the first experiment (Figures 2 and 3). The greater percentage of misses in Experiment II can be taken as evidence for the subjects' greater attention to meaning in this experiment.

The results are entirely consistent with the parallel access (or dual code) hypothesis. When top-down processing is disrupted (Low conditions), effects attributable to the bottom-up variable are observed, and when bottom-up processing is slow (/p/ conditions) detection times are seen to be more a function of top-down processing. The broader implications of these results are considered in the following section.

## GENERAL DISCUSSION

Our research began with the simple question of how phonemes are detected in fluent speech, and specifically, whether detection is necessarily mediated by target word recognition as has been claimed by a number of researchers (Foss & Swinney, 1973; Morton & Long, 1976; Rubin, Turvey, & van Gelder, 1976; Treisman & Squire, 1974). The results of our experiments lead us to question this claim. The compelling top-down effect of target word predictability reported by Morton and Long was found to be limited in its domain. Targets in predictable words were not detected significantly more rapidly than those in unpredictable words, when the phonetic environment was not confusable with the target. In this case a variable that should have strongly affected the speed of word recognition did not influence phoneme detection latency. Our interpretation of this result is that target word recognition is not a prerequisite to phoneme detection. This assumes that the predictability of a particular word in context always influences the speed of recognition of that word. By questioning this assumption, one can interpret the data differently. If one starts from the premise that phoneme detection is purely top-down, then our results seem to say that word recognition is not always affected by predictability. The absence of a predictability effect in an unambiguous phonetic context (/s/-initial conditions) merely means that predictable words are not recognized faster than unpredictable ones in such conditions. Because the latter in-

terpretation assumes that detection is only top-down, however, it is not in accord with other evidence for bottom-up detection (Foss & Blank, 1980; Newman & Dell, 1978). Thus we find that the simplest interpretation of the present results is that bottom-up and top-down detection routes are attempted simultaneously, with the most efficient (i.e., rapid) route determining the response time, which we termed the parallel access hypothesis.

Given that we failed to replicate Morton and Long's (1976) effect in the /s/-initial conditions, it becomes important to look at differences between the two studies. According to our model, the magnitude of the predictability effect in Morton and Long's study (70 milliseconds) indicates that the top-down route was, on the average, faster than the bottom-up route under the conditions of their experiment. Does this imply that Morton and Long's sentences had confusable initial phonemes preceding the target words (analogous to our /p/-initial conditions)? Not necessarily.[1] It is reasonable to assume that many factors affect the relative speeds of the two proposed detection routes, any number of which could have differed between Morton and Long's and the present experiment. Some noteworthy differences between the experiments include the syntactic context of the target words (Morton and Long's targets were nearly always preceded by one or more function words; our targets always followed another content word) and the nature of the target phonemes (our target was always /b/; Morton and Long's targets changed each trial). For each of these differences one can offer hypotheses regarding their effects on the relative speeds of top-down and bottom-up detection in the two studies. For example, the on-line effect

---

[1] In 33 of Morton and Long's 40 sentences, the critical word was an article or possessive pronoun whose initial phoneme never resembled the target. The remaining sentences had targets preceded by content words only two of which had initial phonemes that differed from the target by only one phonemic feature.
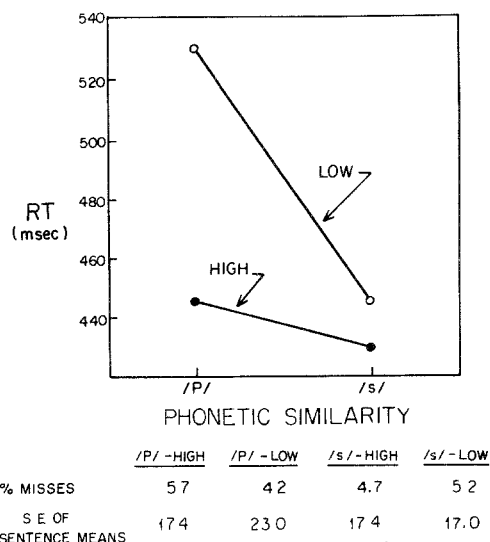


FIG. 3. Detection time as a function of target word predictability and phonetic similarity, Experiment II. (Standard errors expressed in milliseconds.)

| | /P/ -HIGH | /P/ -LOW | /s/ -HIGH | /s/ -LOW |
|---|---|---|---|---|
| % MISSES | 57 | 42 | 4.7 | 52 |
| S E OF SENTENCE MEANS | 17.4 | 23.0 | 17.4 | 17.0 |

of predictability of a word may act more rapidly if that word occurs in a syntactic position that is effectively cued by preceding function words. So the word *chair* might be equally predictable by a Cloze procedure in the sentences *John sat on the chair* and *John sat on the plush chair*, but the word recognition system may concentrate its expectations at the position right after *the*, leading to a stronger on-line predictability effect in the former sentence. Given that Morton and Long's sentences were mostly of the first type, one might expect that their predictability effect would be large and would also be difficult to attenuate. (See Foss and Blank (1980) for a discussion of the robustness of the predictability effect in Morton and Long's sentences.) The top-down route is simply too efficient for the bottom-up route to play any role. Since the target words in our sentences were always preceded by other content words, analogous to the second type of sentence presented above, then one can speculate that the highly predictable target words were not recognized as quickly as would be expected from their high fill-in rates in the pretest. If this were true, then the top-down

route for highly predictable targets would not be as efficient in our sentences as compared to those of Morton and Long, which would explain why we found that the top-down effect was almost completely attenuated under conditions of rapid bottom-up detection.

In sum, we feel that the absence of a predictability effect in the /s/ conditions is not inconsistent with the results of Morton and Long (1976). According to our view of phoneme monitoring, a predictability effect should be obtained if (a) there are differences in the speed of top-down processing as a function of the levels of the predictability variable and (b) the speed of top-down processing for at least one of these levels is faster than that of bottom-up processing. We claim that condition (b) did not hold for our experiments when the critical phoneme was /s/.

So far we have not discussed at length the function in normal comprehension of the phonetic and phonological codes, which are hypothesized to mediate bottom-up and top-down detection. As mentioned before, the phonetic code has the important function of serving as a basis for lexical access. Of what use is the phonological code? Foss and Blank (1980) have suggested that the phonological code can "fill in" holes in the phonetic code that arise from noisy input. This filling-in process, however, does not seem to be particularly useful if one assumes that phonological phonemes of a word only become available after that word is recognized. Given that the function of the phonetic code is to access a word, why bother to fill in the phonetic code if one already knows what the word is? We feel that the phonological code can be given an elevated status if it is assumed that phonological phonemes become available not only as a result of word recognition, but also whenever the existence of a word is hypothesized. That is, activation of a lexical entry from any source activates its phonological phonemes. The phonological phonemes function as hypotheses as to which speech sounds are actually present in the speech signal. The confirmation or disconfirmation of these hypotheses can take place by matching the phonological codes with the developing phonetic code. This idea is certainly not new; it is a restatement of analysis-by-synthesis (e.g., Halle & Stevens, 1964) but it specifies hypothesis generation as the activation of a set of phonological phonemes, and hypothesis testing as the comparison of phonological and phonetic codes. The original analysis-by-synthesis theory was more concerned with the problem of the lack of segment invariance, hence it postulated a matching process at the acoustic level. The research reported here, by supporting the existence of two phonemic codes, and the notion that access to the codes occurs in parallel, is consistent with the view being developed. It should be noted that this view differs markedly from that of researchers who see phonological phonemes as nonfunctional or "abstract" (Savin & Bever, 1970) or as functioning only to create perceptual closure as occurs in the phoneme restoration effect (Warren, 1970).

Because we are hypothesizing that the top-down aspects of phoneme perception are actually involved in the word recognition process rather than resulting from it, we would like to discuss the role that phoneme representations might play in a current theory of word recognition. The active direct-access, or cohort, theory of word recognition (Marslen-Wilson & Welsh, 1978; Marslen-Wilson, Note 1) is unique in that contextual and stimulus information are used in an efficient manner to achieve recognition at the earliest opportunity. Unlike a logogen approach (Morton, 1969) in which recognition occurs when a threshold is reached, the active direct-access theory holds that recognition of a word occurs when all other possibilities have been eliminated. The elimination of possible word candidates is accomplished largely through an unspecified process that detects mismatches between word candidates and

incoming stimulus information. Our hypothesis is that the detection of mismatches could be carried out through comparison of the phonological codes from the word candidate set with the phonetic code representing current stimulus information. Although we have no direct support from our experiments for this hypothesized matching process, we are able to make some specific conclusions regarding the implications of our results for a theory of word recognition.

First, our research, together with that of Newman and Dell (1978) and Foss and Blank (1980) provides clear evidence that a target phoneme can be detected before the recognition of the word containing it. This conclusion gives new life to the idea that some sort of linguistic representation (phonemes or features) is constructed prior to lexical access, and that this representation provides the basis for lexical access. This idea has been unpopular with some researchers (e.g., Warren, 1976) precisely because of the supposed exclusively top-down nature of the phoneme.

Second, the fact that two distinct sources of information enter into the phoneme detection process has obvious parallels with the literature on visual word recognition. In that research there is also considerable debate concerning the unit of perception. Some of the most persuasive evidence (Johnson, 1975) shows that, with singly presented stimulus words, it is possible to detect a word faster than it is possible to detect letters within it. This result is analogous to that of Foss and Swinney (1973), who found that it is possible to detect an auditorily presented word faster than it is possible to detect a phoneme within the word. The discussion of the role of letters in word perception seems to mirror the arguments used in the "abstractness" of the phoneme debate. For these reasons, we feel that much the same arguments that we have used concerning the role of phonemes in speech perception can be applied to the role of letters in word perception. It is entirely possible that there are two (or more) routes

to letters, and that it is this fact that has led to the debate over the usefulness of letters in reading. In fact, Chambers and Forster (1975) have proposed that performance in a visual matching task involves three levels of analysis, which operate simultaneously.

To summarize, we have shown that there are two important variables which determine the means by which phonemes are detected in the phoneme monitoring task. Each of these variables has been shown elsewhere to affect phoneme monitoring latencies. Morton and Long (1976) demonstrated that target word predictability, a variable that involves lexical access, affects the time to detect a target phoneme, thereby suggesting that phonemes are constructed after the word containing them has been recognized. This result seemed to show that phonemes are detected in a purely top-down manner, supporting the idea that phonemes are abstract entities. In contrast, Newman and Dell (1978) showed that phoneme monitoring latencies are strongly affected by the phonetic similarity between the critical and target phonemes; the effects of this variable are hypothesized to occur before lexical access, and reflect the bottom-up processing of the target phoneme. In the current experiments, we have reconciled these two, apparently contradictory, results by showing that there are two routes to phonemes; the effects of target word predictability were eliminated when the critical phoneme was not confusable with the target; when the target and critical phonemes were similar, the target was detected through lexical access. These results prompted us to propose a new model of phoneme monitoring involving parallel access of phonemes in this task. Our results have broader implications than just those proposed for phoneme monitoring; we have provided evidence that the phoneme (which has often been dismissed as redundant or abstract) is in fact used in word recognition. We have extended our model and shown how our parallel access hypothesis can complement a current

theory of word recognition (Marslen-Wilson & Welsh, 1978). Our suggestion allows for the generation and confirmation of hypotheses concerning current stimulus information during on-line word recognition. In conclusion, our results argue for the reinstatement of the phoneme as an important component of speech perception.

## APPENDIX

The critical and target words associated with the four experimental conditions are given in italics, the /p/-initial critical word first, followed by, in order, the /s/-initial critical word, the highly predictable target word, and the unpredictable target word.

1. The surfers drove to a (private/secret beach/bay) to try out the waves.
2. When Mary kissed him (Paul/Sam blushed/blinked) and turned away.
3. The fisherman's lead sinker dropped to the (pond's/stream's bottom/bank) and was lost.
4. After filling the tub (Peter/Sally bathed/basked) in the water.
5. The sparrow landed on the (pine/cedar branch/bed) and rested awhile.
6. For winter weather the shoe department featured (plastic/suede boots/belts) as well as leather ones.
7. Musical instruments such as the trumpet are made of (pure/solid brass/bronze).
8. While her toast was still warm (Polly/Susan buttered/brushed) it with dairy spread.
9. After pouring water over the kneeling woman the (priest/saint blessed/bound) her and directed her to pray.
10. The muscular Mr. Canada contestants had (perfect/strong biceps/blood) which indicated good health.
11. Moving high in the sky the (pretty/swift bird/ball) was easily seen.
12. The mosquito left a (painful/stinging bite/blister) on her arm.
13. On her wrist (Paula's/Sarah's bracelet/brooch) looked rather strange.
14. On his shirt the man sewed a (plain/single button/badge).
15. When Henry dropped it the (platter/saucer broke/bounced) on the floor.
16. The red cape placed over his horns enfuriated the (proud/savage bull/buck).
17. Strong forces make (platinum/steel bend/buckle) very easily.
18. The juicy worm caught in the (pigeon's/swallow's beak/belly) provided a nourishing meal.
19. The cook mixed the flour and eggs in the (purple/silver bowl/basket).
20. Johnny is a very (polite/silly boy/baby) which is unusual for someone his age.
21. The very long and tedious (play/sermon bored/baffled) its audience.
22. That well-constructed house is the one that (Patrick/Steven built/bombed).
23. The insect with the handsome wings was a (passing/sleeping butterfly/bee).
24. Always looking for a scapegoat the (press/schools blamed/blasted) the federal government.
25. The driver asked the passengers to get off the (parked/stalled bus/bicycle) for their own safety.
26. The animal gnawing on the tree was a (plump/sluggish beaver/bear).
27. In order to hide their treasure the (pirates/sailors buried/boarded) it in the cave.
28. At the tavern the longshoreman enjoyed a (pleasant/soothing beer/brandy) and good company.
29. The dynamite in the (pack/sack blew up/blocked) the entrance to the castle.
30. After reading the (picture/sad book/bill) the whole family discussed it.
31. At the lending library the (pupil/student borrowed/bought) the damaged text.
32. The very hot oven made the (pie/sauce burn/bubble).

## REFERENCES

BLANK, M. A., & FOSS, D. J. Semantic facilitation and lexical access during sentence processing. *Memory & Cognition*, 1978, 6, 644–652.

CHAMBERS, S. M., & FORSTER, K. I. Evidence for lexical access in a simultaneous matching task. *Memory & Cognition*, 1975, 3(5), 549–559.

CHANANIE, J. D., & TIKOFSKY, R. S. Choice response time and distinctive features in speech discrimination. *Journal of Experimental Psychology*, 1969, 81, 161–163.

COLE, R. A., & JAKIMIK, J. Understanding speech: How words are heard. In G. Underwood (ed.), *Strategies of information processing*. London: Academic Press, 1979.

COLE, R. A., & SCOTT, B. Distinctive feature control of decision time: Same–different judgments of simultaneously heard phonemes. *Perception & Psychophysics*, 1972, 12 (1B), 91–94.

CUTLER, A., & NORRIS, D. Monitoring sentence comprehension. In W. E. Cooper & E. C. T. Walker (eds.), *Sentence processing: Studies in honor of Merrill Garrett*, Hillsdale, N.J.: Erlbaum, 1979.

EIMAS, P. D., & CORBIT, J. D. Selective adaptation of linguistic feature detectors. *Cognitive Psychology*, 1973, 4, 99–109.

FORSTER, K. I., & DICKINSON, R. G. More on the language-as-fixed-effect fallacy: Monte Carlo estimates of error rates for $F_1$, $F_2$, $F'$, and min $F'$.

*Journal of Verbal Learning and Verbal Behavior,* 1976, *15,* 135–142.

Foss, D. J. Decision processes during sentence comprehension: Effects of lexical item difficulty and position upon decision times. *Journal of Verbal Learning and Verbal Behavior,* 1969, *8,* 457–462.

Foss, D. J., & Blank, M. A. Identifying the speech codes. *Cognitive Psychology,* 1980, *12,* 1–31.

Foss, D. J., Harwood, D. A., & Blank, M. A. Deciphering decoding decisions: Data and devices. In R. Cole (ed.), *Perception and production of fluent speech,* Hillsdale, N.J.: Erlbaum, 1980.

Foss, D. J., & Jenkins, C. M. Some effects of context on the comprehension of ambiguous sentences. *Journal of Verbal Learning and Verbal Behavior,* 1973, *12,* 577–589.

Foss, D. J., & Lynch, R. Decision processes during sentence comprehension: Effects of surface structure on decision times. *Perception & Psychophysics,* 1969, *5,* 145–148.

Foss, D. J., & Swinney, D. On the psychological reality of the phoneme: Perception, identification, and consciousness. *Journal of Verbal Learning and Verbal Behavior,* 1973, *12,* 246–257.

Halle, M., & Stevens, K. N. Speech recognition: A model and a program for research. In J. A. Fodor & J. J. Katz (eds.), *The structure of language: Readings in the philosophy of language.* Englewood Cliffs, N.J.: Prentice-Hall, 1964.

Johnson, N. F. On the function of letters in word identification: Some data and a preliminary model. *Journal of Verbal Learning and Verbal Behavior,* 1975, *14,* 17–29.

Klatt, D. A new look at the problem of lexical access. In R. Cole (ed.), *Perception and production of fluent speech.* Hillsdale, N.J.: Erlbaum, 1980.

Kučera, H., & Francis, W. N. *Computational analysis of present-day American English.* Providence, Rhode Island: Brown University Press, 1967.

Marslen-Wilson, W. D., & Welsh, A. Processing interactions and lexical access during word recognition in continuous speech. *Cognitive Psychology,* 1978, *10,* 29–63.

Mehler, J., Sequi, J., & Carey, P. Tails of words: Monitoring ambiguity. *Journal of Verbal Learning and Verbal Behavior,* 1978, *17,* 29–35.

Morton, J. Interaction of information in word recognition. *Psychological Review,* 1969, *76,* 165–178.

Morton, J., & Long, J. Effect of word transition probability on phoneme identification. *Journal of Verbal Learning and Verbal Behavior,* 1976, *15,* 43–52.

Newman, J. E., & Dell, G. S. The phonological nature of phoneme monitoring: A critique of some ambiguity studies. *Journal of Verbal Learning and Verbal Behavior,* 1978, *17,* 359–374.

Rubin, P., Turvey, M. T., & van Gelder, P. Initial phonemes are detected faster in spoken words than in non-words. *Perception & Psychophysics,* 1976, *19,* 394–398.

Savin, H. B., & Bever, T. G. The nonperceptual reality of the phoneme. *Journal of Verbal Learning and Verbal Behavior,* 1970, *9,* 295–302.

Shields, J. L., McHugh, A., & Martin, J. G. Reaction time to phoneme targets as a function of rhythmic cues in continuous speech. *Journal of Experimental Psychology,* 1974, *102,* 250–255.

Swinney, D. A., & Hakes, D. T. Effects of prior context upon lexical access during sentence comprehension. *Journal of Verbal Learning and Verbal Behavior,* 1976, *15,* 681–690.

Treisman, A., & Squire, R. Listening to speech at two levels at once. *Quarterly Journal of Experimental Psychology,* 1974, *26,* 82–97.

Warren, R. M. Perceptual restoration of missing speech sounds. *Science,* 1970, *167,* 392–393.

Warren, R. M. Identification times for phoneme components of graded complexity and for spelling of speech. *Perception & Psychophysics,* 1971, *9,* 345–349.

Warren, R. M. Auditory illusions and perceptual processes. In N. J. Lass (ed.), *Contemporary issues in experimental phonetics.* New York: Academic Press, 1976. Pp. 389–417.

## REFERENCE NOTES

1. Marslen-Wilson, W. D. *Sequential decision processes during spoken word recognition.* Paper presented at the Annual Meeting of the Psychonomic Society, San Antonio, Texas, November 1978.

2. Vipond, D. *Influence of pretarget word length on phoneme monitoring reaction times.* Unpublished manuscript, University of Colorado, January 1977.