

PHONEME RECOGNITION AS A FUNCTION OF TASK AND CONTEXT

R.J.J.H. van Son and Louis C.W. Pols

Institute of Phonetic Sciences (IFA)/ACLC, University of Amsterdam, The Netherlands
Rob.van.Son@hum.uva.nl

ABSTRACT

Phoneme recognition can mean two things, conscious *phoneme-naming* and pre-conscious *phoneme-naming*. Phoneme naming is based on a (learned) label in the mental lexicon. Tasks requiring phoneme awareness will therefore exhibit all the features of retrieving lexical items. Phone categorization is a hypothetical pre-lexical and pre-conscious process that forms the basis of both word recognition and phoneme naming. Evidence from the literature indicates that phone-categorization can be described within a pattern-matching framework with weak links between acoustic cues and phoneme categories. This is illustrated with two experiments. The current evidence favors a lax-phoneme theory, in which all phonemic categories supported by the acoustic evidence and phonemic context are available to access the lexicon. However, the current evidence only supports segment-sized categories. It is inconclusive as to whether these categories are the same in number and content as the phonemes.

1. INTRODUCTION

In the phonetic literature, *phoneme recognition*, is generally used in two distinct senses. In one sense, phoneme recognition refers to "phoneme awareness". In the other sense, it refers to a hypothetical intermediate, symbolic, representation in speech recognition. Phoneme awareness comes with (alphabetic) literacy. Most literate people can easily identify or monitor phonemes in speech. This awareness builds on a "deeper", automatic categorization of sounds into classes of which the listeners are not consciously aware. In the first sense, phoneme recognition is a form of word-recognition. In the second sense, it is a form of data-reduction that is hidden from awareness. We will refer to phoneme recognition in the former, conscious sense as *phoneme-naming* and in the latter, pre-conscious, sense as *phoneme-categorization*.

Phoneme-naming and phone-categorization are not identical. It is clear that in conscious phoneme-naming, labels are attached to the categories found in pre-conscious phone-categorization. However, the whole word first and then extracting the constituent phonemes from the lexicon [1]. The conscious, lexical aspects of phoneme-naming will induce effects in all experiments that rely on it. Obvious differences with phone-categorization are that lexical decisions are known to be competitive (winner-takes-all), frequency dependent, and prime-able. However, there is no reason to assume that the underlying categorization is competitive, the frequency effects reported are intricate at best [2], and prime-ability might even be detrimental to word recognition. Furthermore, conscious awareness of phonemes and the associated attention allow the recruitment of "higher" mental modules that are inaccessible to unconscious processes [3].

It is not clear whether the phoneme is really a "natural" element in recognition. Normally, phonemes are defined as distinctive feature bundles. That is, a phoneme is the smallest unit that will distinguish between words, e.g., [tɛnt] versus [dɛnt] or [kɛnt]. In these examples, /t d k/ are phonemes that differ in the feature voicing (/t/-/d/), place of articulation (/t/-/k/), or both (/d/-/k/). Not all combinations of feature values that theoretically could be combined in a phoneme actually occur in a language. Languages ensure that differences between phonemes are large enough to be kept apart easily in both articulation and identification [4][5]. Of the 600 or more sounds that can be distinguished in the world's languages, English uses less than 50. Furthermore, between languages, features and phonemes, can be defined differently, making for even more differences. For instance, both English [tɛnt] and [dɛnt] are transcribed as Dutch [tɛnt] whereas both Dutch [tɛnt] and [dɛnt] are transcribed as English [dɛnt]. Not all phoneme combinations are possible. [tɛnt] can legally be changed into [tɛnd]. But a change to [tɛnk] results in an invalid English word. To get a valid English word, we have to change the place of articulation of the whole cluster /nt/, e.g., [tɛnk] is a "phonotactic" rule in English that "forbids" /nk/ clusters. All languages have such rules, but they are different for each language (e.g., [tɛnd] is not a valid Dutch word). The phonotactic, and phonological, rules are a second (syntactic) layer that have to be added to the phonemes to get a workable set. In a sense, the phonemes define legal feature combinations and *phonotactic rules* define legal feature sequences. Therefore, it should not come as a surprise that phonemes and phonotactics are complementary in speech recognition. People have difficulty producing and perceiving both phonemes with invalid feature combinations as well as feature (phoneme) sequences that violate phonotactic rules [6]. Speech that violates the combinatorial rules of the features in a language will generally be mapped to the nearest valid phoneme sequence. This is a problem in second language learning as many (most) students never succeed in completely mastering the new phonemes and phonotactic rules. We can capture thinking on phoneme recognition in terms of two extreme positions, which few phoneticians will actually defend. At the one extreme, the *obligatory phoneme hypothesis* states that all speech is internally represented as a string of phonemes and whether we use phonemes or features in this hypothesis is actually immaterial as all legal feature collections can be rewritten as legal phoneme sequences and vice versa. However, note that current theories of word recognition do not need phonemes or features. Most models would just as well work on "normalized" sound traces. This brings us to the other extreme, the *lax phoneme hypothesis*. In this lax phoneme hypothesis, recognizing

2. THE UNITS OF SPEECH

An experiment we performed some years ago illustrates the problems of theories relying on static or dynamic specification [9][10][11]. In our experiment we compared the responses of Dutch subjects to isolated synthetic vowel tokens with *curved* formant tracks (F1 and F2) with their responses to corresponding tokens with *stationary* (level) formant tracks. We also investigated the effects of presenting these vowel tokens in a synthetic context (/nVf/, /Vn/).

Nine formant "target" pairs (F1, F2) were defined using published values for Dutch vowels. These pairs

5. EXPERIMENTAL ILLUSTRATION

The second type of approach is *dynamic*. It assumes that the dynamics of speech generation predetermine deviations from the canonical target realizations. These deviations can be "undone" by the extrapolation of the appropriate parameter tracks (dynamic specification, see [10][11]) or by some detailed modeling of the mechanical behavior of the articulators (Motor theory). Experimental evidence for any of these theories has been hotly disputed. As Narey [8] rightfully remarks: Proponents of both approaches make such a good case of disproving the other side, that we should believe them both and consider both disproved.

The previous discussion is "phonological" in nature in that no reversences were made to the acoustics, articulation, or perception of speech sounds. Features and phonemes are symbolic entities that have to be linked to acoustic categories to be of any use in speech communication. Two classical approaches to the perceptual categorization problem can be distinguished. First, are the *static clustering theories*. These theories assume that each phoneme is a simple perceptual category. This category is defined as a unit cluster in some perceptual space. Some, rather complicated, transformation is performed on the speech signal after which the kernel (center) of each phoneme realization will map to a point inside the boundaries of the perceptual area designated for that phoneme. The best known example of this kind of approach is the Qunatal theory of speech [18].

4. THE ACOUSTICS OF PHONEMES

The lax phoneme hypothesis might at first not seem to require labeling each phone with a "master" phoneme label. However, for lexical access, each phone has to be reevaluated to determine its proper place in the utterance. For instance, /hO-l-@-fem/ must be resyllabified to /hO-l Of fem/ to be recognized as a three word phrase. The identity of the pre- and post-vocal /l/ and /f/ sounds is not trivial. At some level, even a lax-phoneme model should facilitate this exchange of allphones.

along the lines of maximal communicative *efficiency* in both production and perception [4][5]. This will favor "simple" inventories and rules. Still, each language-community can "choose" freely what variation it does or does not permit [4]. Our proposition is that phonemes are not only characterized by some perceptual "canonical form", but that phonotactical constraints and phonological rules are an integral part of phoneme identity. *A phone is the realization of a phoneme only in a certain context.* This is well illustrated by the fact that contexts that violate phonotactics hamper phoneme recognition [6].

Divergent allophones of a phoneme do not have to share any perceptual properties. Their unity at the phoneme level could, in principle, be completely arbitrary. That the allophones of a phoneme almost always do share some fundamental properties can be explained from the fact that phoneme inventories and the

Very often, only the context of a phone allows one to select the intended phoneme label. That these complex collections of "context dependent" phones are genuine objects and not artifacts of a procrustean theory is clear from the fact that both speakers and listeners can seamlessly handle the rather complex transformations to "undo" reduction, coarticulation, and resyllabification (e.g., "hall-of-fame" as /hO-l-@-fem/ or "wreck a nice beach" as /fE-k-@-n-AI-sptIS/).

One central presupposition that many theories on phoneme recognition share is that each phoneme has a unified (and unique) canonical *target* to which a realization can be matched (see [7] for evidence that this is a *perceptual target*). In this view, phone-categorization and phoneme-naming use the same "labels". However, many *phones* of a given *phoneme* do not overlap in any perceptual representation, e.g., pre- and postvocalic liquids, glides, and plosives (c.f., aspirated and unaspirated allophones of voiceless plosives). Whereas other phonemes share the same phones (but in different contexts), e.g., long and short vowels. This can be most clearly seen when phonotactically defined allophones in one language are distinct phonemes in another (dark and light /l/ in English or Dutch are two separate phonemes in Catalan).

3. WHAT MAKES A PHONEME

To summarize the obligatory and lax phoneme hypotheses. The obligatory hypothesis states that all words (utterances) are mentally represented as phoneme strings. The lax phoneme hypothesis states that phonemes and phonotactics are features of a data-reduction process that selects and organizes relevant information, but doesn't force decisions on segmental identity. In the lax hypothesis, missing information is ignored and strict categorization is deferred if necessary. In the obligatory phoneme hypothesis, missing information has to be provided (invented) during a forced phoneme categorization.

A simple way of coding words in a robust way is to correlate acoustic events for sequential order and co-occurrence. Essentially, this is a "context-sensitive" clustering analysis in which the speech stream is (partially) categorized in a normalization and data-reduction step. After this data-reduction step, the remaining information can be processed further to fit the requirements of the lexicon.

Humans are able to repeat new words, as all processing and, not least, allow reproduction, sequential noisy evidence, the words have to be coded in ways that allow normalization, error correction, sequential unlimited number of speakers, on-line with incomplete, ambiguous. To be able to recognize ~10⁶ words from an track reading, the tracks will be incomplete and underlying words (or pronounceable sounds). As in all speech is tracking acoustic events for evidence of

and regularization of utterances as a precursor for lexical access.

No tracks were constructed that would cross other formant tracks or F_0 . All tracks were synthesized with durations of 25, 50, 100, and 150 ms (see for more details: [9][10][11]). Stationary tokens with level formant tracks (i.e., $\Delta F_1 = \Delta F_2 = 0$) were also synthesized with durations of 6.3 and 12.5 ms. Of the other tokens (with either $\Delta F_1 = +/ -225$ Hz or $\Delta F_2 = +/ -375$ Hz), the first and second half of the tracks, i.e., on- and off-glide-only, were also synthesized with half the duration of the "parent" token (12.5, 25, 50, and 75 ms). Some other tokens with smaller excursion sizes were used too, these will not be discussed here (but see [9][10][11]). In experiment 1, tokens were presented in a pseudo-random order to 29 Dutch subjects who had to mark the orthographic symbol on an answering sheet with all 12 Dutch monophthongs (forced choice).

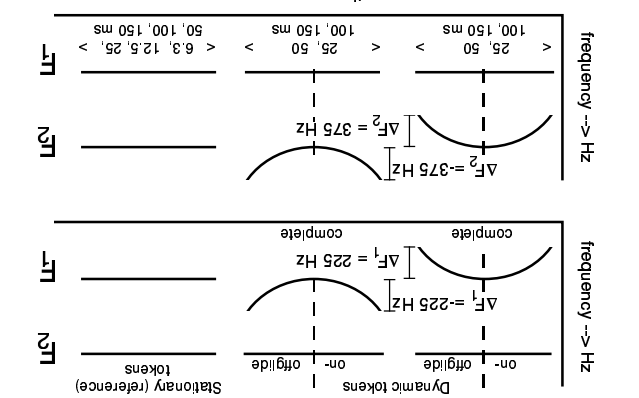
In experiment 2, a single realization each of 95 ms synthetic /n/ and /f/ sounds were used in mixed pseudo-syllabic stimuli. Static and dynamic vowel tokens from the first experiment with durations of 50 and 100 ms and mid-point formant frequencies corresponding to /I E A o/ and /Vn/ pseudo-syllables. The corresponding vowel tokens with only the on- or off-glide part of parabolic formant tracks (50 ms durations only) were used in CV and VC structures respectively. For comparison, corresponding stationary vowel tokens with 50 ms duration were also used in CV and VC pseudo-syllables. Each vowel token, both in isolation and in these pseudo-syllables, was presented twice to 15 Dutch subjects who

$$F_n(t) = Target - \Delta F_n \cdot (4 \cdot (t/D)^2 - 4 \cdot t/D + 1)$$

in which:
 $F_n(t)$: Value of formant n (i.e., F_1 or F_2) at time t.
 ΔF_n : Excursion size, $F_n(\text{mid-point}) - F_n(\text{on/offset})$.
 $\Delta F_1 = 0, +225$ or -225 ; $\Delta F_2 = 0, +375$ or -375 (Hz).
 $Target$: Formant target frequency.
 D : Total token duration ($0 < t < D$).

Figure 1. Formant track shapes as used in the experiments discussed in section 5. The dynamic tokens were synthesized with durations of 25, 50, 100, and 150 ms. The stationary tokens were synthesized with durations of 6.3, 12.5, 25, 50, 100, and 150 ms. The dynamic tokens were also synthesized as on-glide- and off-glide-only tokens, i.e., respectively the parts to the left and right of the dashed lines.

corresponded approximately to the vowels /u/yo/EAy/ and were tuned to give slightly ambiguous percepts. For these nine targets, smooth formant tracks were constructed for F_1 and F_2 that were either level or parabolic curves according to the following equation (see figure 1):

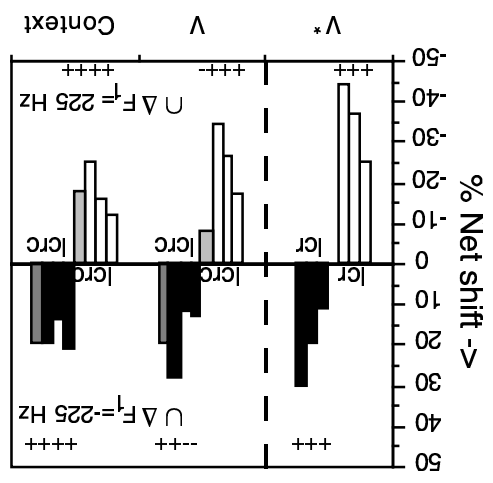


Contrary to the predictions of the *dynamic* models of speech recognition, there was *no extrapolation* found. Contrary to the predictions of the *static clustering theories*, the *kernel* was not exclusively used for identification. None of the theories predicted, or can even explain, the prevalence of averaging responses to dynamic stimuli. No segment internal cues seem to be

For each response to a dynamic token, the position in formant space with respect to the static token was determined. For instance, an /E/ response to a dynamic token was considered to indicate a higher F_1 perception and a lower F_2 perception than an /I/ response to the corresponding static token. By subtracting the number of lower dynamic responses from the number of higher dynamic responses, we could get a *net-shift* due to the dynamic formant shape (testable with a sign test). Analysis of all material clearly showed a very simple pattern over all durations: Responses *averaged* over the trailing part of the F_1 tracks (figure 2). The same was found for the curved F_2 tracks (not shown), although here the effects were somewhat weaker and not always statistically significant (see [9][10][11] for details). The use of vowels in /n/, /f/ context had no appreciable effect except for a decreased number of long vowel responses in open "syllables" (not shown). In accord with the Dutch phonotactical rule against short vowels in open syllables. However, this lack of effect could be an artifact from an unnatural quality of the pseudo-syllables.

The speech recognition theories discussed above make clear predictions about the behavior of our listeners. Static theories predict that vowel identity is largely unaffected by formant dynamics. Dynamic theories predict some compensation for reduction in dynamic stimuli. In our case, all dynamic theories would predict formant track extrapolation in some form (perceptual overshoot [9][10][11]).

Figure 2. Net shift in responses as a result of curvature of the F_1 . 'V*' are the results of the first experiment (all tokens pooled on duration, $n >= 696$). 'V' and 'Context' are the results of the second experiment with vowel tokens presented in isolation ('V': $n=120$, left; $n=90$, right), or in context, CV, CVC, VC; C one of /n/ f/ tokens, white/black bars: 50 ms tokens, I=on-glide-only, c=complete, r=off-glide-only tokens. +: significant ($p < 0.001$, sign test), -: not significant. Results for F_2 were comparable but weaker.



The pattern-matching framework can be illustrated by [14] and by a study of our own [15]. In the latter study, we constructed gated tokens from 120 CVC speech fragments taken from a long text reading where the sentence accent on the vowel was noted (see [15] for details). The tokens were divided into vowel kernel (kernel, the central 50 ms), vowel transition (T, everything outside the kernel), a short consonant part (C, 10 ms), and an overlapping longer consonant part (CC, 25 ms). Gated tokens were constructed from these segments according to figure 3. These tokens were randomized for vowel identification (Kernel, V, CV, VC, CVC tokens), pre-vocalic consonant identification (CT, CCT, CV, CCV tokens), and post-vocalic consonant identification (TC, TCC, VC, VCC tokens). Listeners were asked to identify vowels (17 subjects), and pre- or post-vocalic consonants (15 subjects for both) by picking the relevant orthographic symbol on a CRT screen. Subjects could pick any legal phoneme, except as well as usual consonants (/Jg/, affricates). For more details, see [15].

The results were analyzed in terms of the $\log_2(\text{response-perplexity})$ which is a measure of the missing information (in bits) and is measured on a ratio scale. Contrary to its complement, the mutual information or transmission rate, the missing information is insensitive to the size of the stimulus set. The results are summarized in Figure 4.

Two conclusions can be drawn from figure 4. First, phoneme identification benefits from extra speech, even if it originates from outside the segment proper, e.g., adding a vowel kernel to a transition-consonant fragment. Second, adding contextual speech in front of a phoneme improves identification more than appending it at the back. Our results could be best explained by assuming that, as predicted by the weak theory of speech perception, listeners use speech cues from far outside the segment to identify it. Furthermore, the prevalence of early (pre-gated) speech over late (appended) cues indicates that phoneme-identification (a naming task) is a fast process in which *label-decisions* are made as early as possible, disregarding subsequent cues.

7. PHONEMES IN CONTEXT

Phoneme production (articulation) is highly context dependent. Coarticulation is one of the prime sources of variation in phonemes. It changes all aspects of speech to such an extent that no genuine acoustic invariant has been found in 50 years of research. So it would be no surprise if phoneme recognition itself would be context dependent.

A reanalysis of classical studies on dynamic theories of phoneme recognition [10][11], showed that all of these studies could be interpreted in terms of purely phonemic-context effects: Only if the appropriate context was identified, was there compensation for coarticulation. Furthermore, the amount of perceptual compensation depended *only* on the context, and was independent of the size of any dynamic aspect of the speech. In an extensive in-depth analysis of earlier experiments, Narey gives very convincing arguments for the use of phoneme-sized, symbolic context in phoneme recognition [16]. If we summarize these studies, listeners seem to interpret acoustic cues not with respect to their *acoustic* context, but instead with respect to their *phonemic* context [8][9][10][11][12][16].

This phonemic-context effect can be illustrated with results from our own study (Figure 5, [15]). Both vowels

build up over time [20][21]. The lattice would contain "activations" (or sublabeling [21]). These activations like lattice of possible "phone(m)-categories" and their "activations" as the construction of an ASR-

This can be visualized as the construction of an ASR-performed as possible without discarding relevant cues during categorization as much data reduction is expected, and phoneme restoration. It seems as if categorical boundaries shift in response to lexical demonstrations in the Ganong effect [17], where that ambiguities are preserved [20][21], as is e.g., the McGurk effect [19]. Categorization is also lax in visual information is used in the categorization process, and positive voicing [16]. Furthermore, when present, phoneme, e.g., vowel duration in vowel identification each individual cue can be used for more than one process. First of all, acoustic cues are "recycled" and a little is already known of the categorization *phonemic* (symbolic) context [10][11][15][16].

normalization and compensation depends on the evidence is compatible with a view that the compensate for coarticulation and reduction. Most combine speech cues into phoneme-sized categories that What is still needed is a mechanism to normalize and categories used to access the lexicon or articulation [7].

phome size [8][12][13]. The output are phoneme-sized speech signal that map directly onto symbolic entities of be a combination of acoustic and visual "events" in the pattern-matching framework [8][12][13]. The input will points can already be made. A good case has been several theories on phoneme recognition. However, several

It is probably too early to present a real synthesis of

8. A SYNTHESIS?

and consonants in CV tokens were identified better when the other member of the pair was identified correctly than when it was identified incorrectly. The fact that nothing was found for the VC tokens makes it less likely that this effect was only due to the fact that a better articulated vowel implicated a better articulated consonant. However, we do not have an explanation of the difference between CV and VC tokens.

Figure 5. Error rates for vowel and consonant identification in CV- and VC-type tokens with respect to the correct and incorrect identification of the other segment in the same token. Voiceless errors in consonants and Long/Short errors in vowels were ignored. Differences are statistically significant for the CV tokens only (Chi-square = 28.7, n = 1, p < 0.01).

